# CS 318 Principles of Operating Systems

## Fall 2017

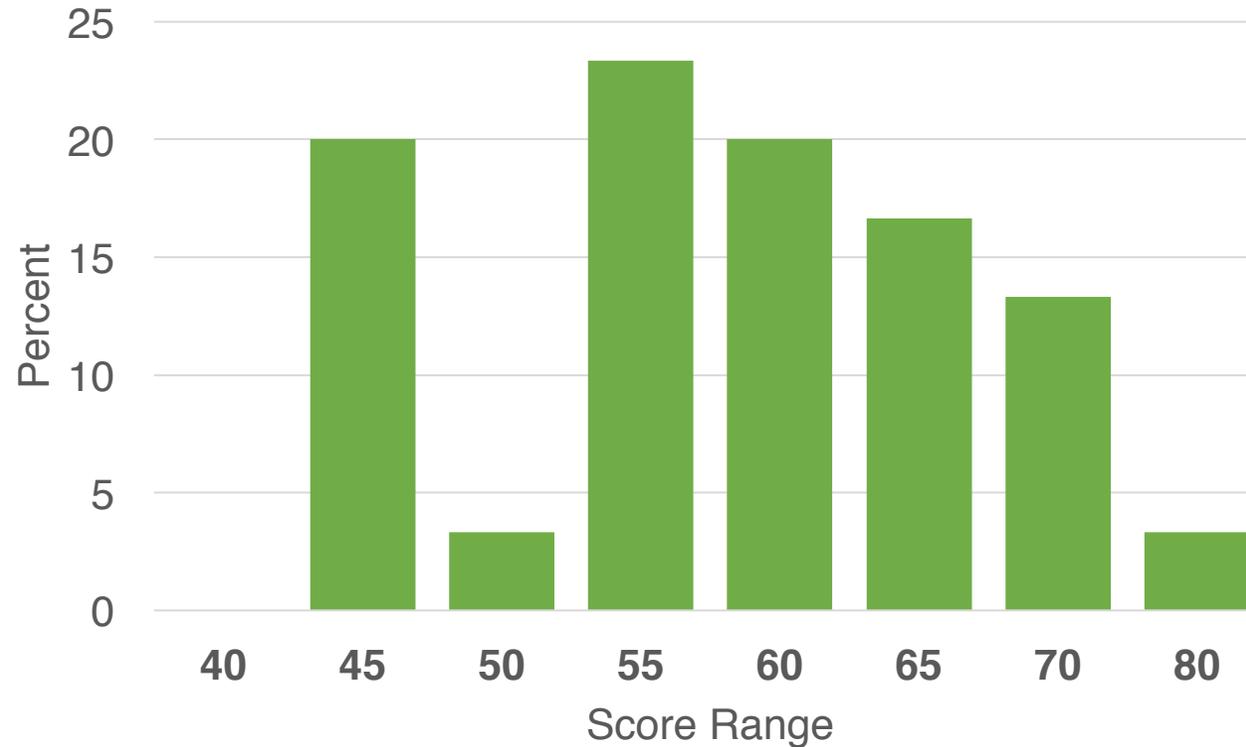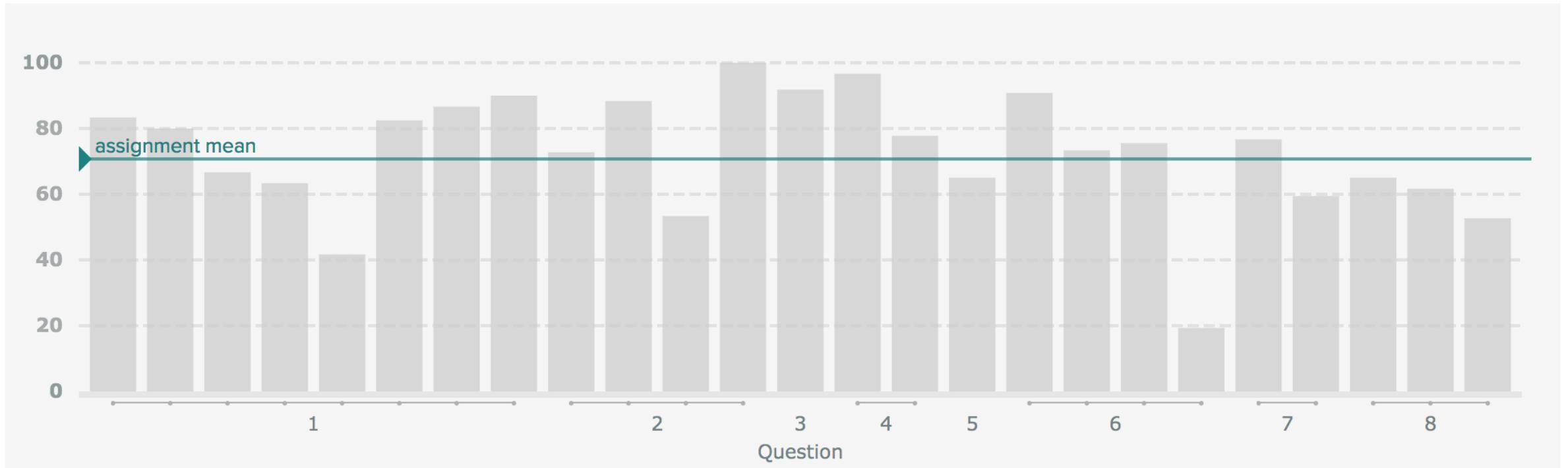## Lecture 14: I/O & Disks

Ryan Huang

JOHNS HOPKINS

WHITING SCHOOL
*of* ENGINEERING

# Midterm Results



- **Mean: 56.5, Max: 72.5, STD Dev 8.7**
  - 318 Section: Mean 58.5 (amazing!), Max 70, STD Dev 9.2

# Midterm Results

# Midterm Results

- **Tend to overthink a problem**
  - The synchronization and dining graduate problem are directly adapted from homework and textbook
  - A variety of overly complex (incorrect) answers

- **Some serious misconception**
  - E.g., syscall makes user-level threads faster (Q4 also from homework)

- **Don't panic if you didn't do well on midterm**
  - Still a lot of chance to make up, e.g., do Lab 3 well
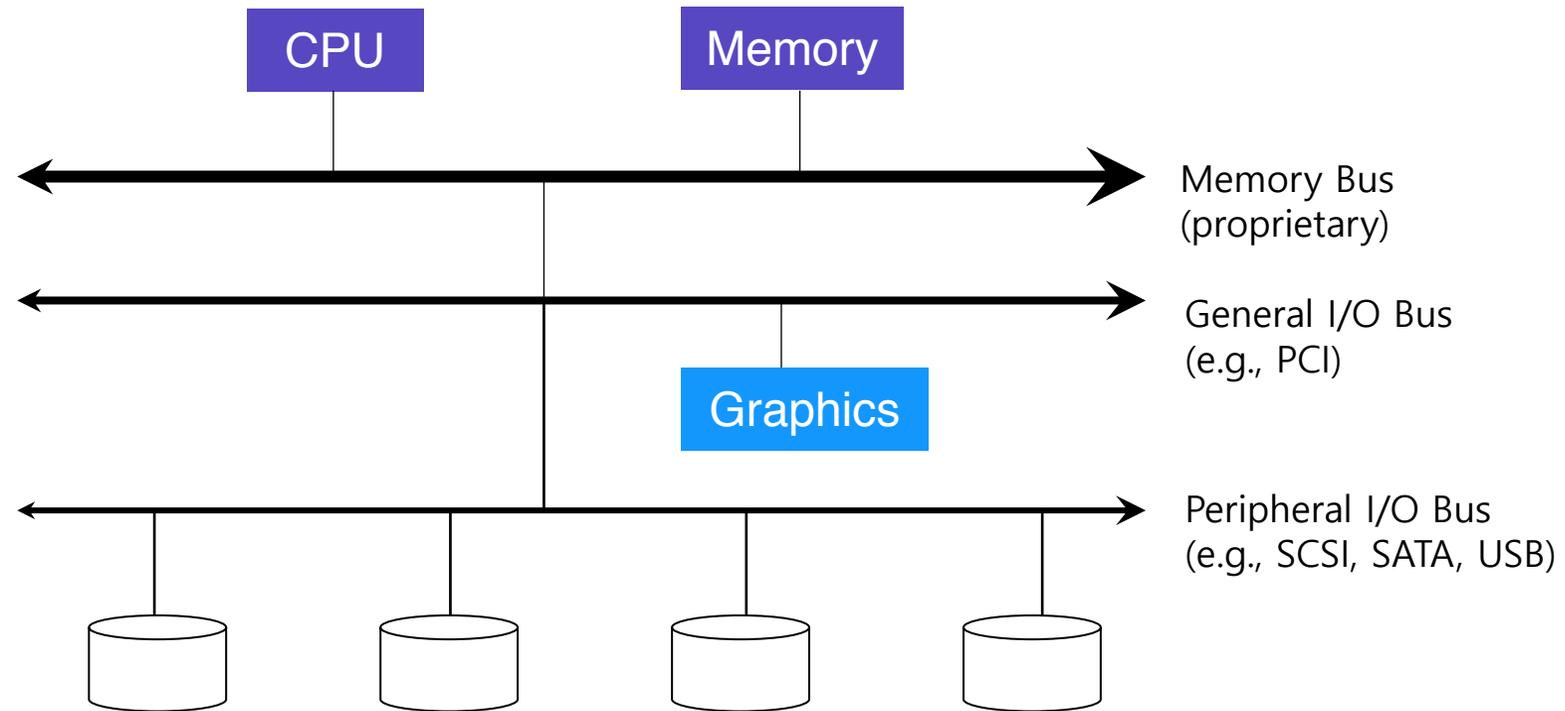  - But do make sure you understand all the questions and answers now

# Administrivia

- **Midterm solution will not be directly posted online**
  - The Rubric items are published in gradescope
  - Can request to see the copy of sample solution in my office hour or the CAs'

- **If there is issue with grading, email the staff list or request through gradescope**

- **If you want to talk about midterm, don't hesitate to contact me**

- **Lab 3 is out, please start early**
  - workload increasing
  - absolute late penalty increasing
  - suggest checking design with the staff first

# I/O Devices

- **I/O is critical to computer system to interact with systems.**

- **Issue:**
  - How should I/O be integrated into systems?
  - What are the general mechanisms?
  - How can we make the efficiently?

# Structure of Input/Output (I/O) Device



CPU

Memory

Memory Bus
(proprietary)

General I/O Bus
(e.g., PCI)

Graphics

Peripheral I/O Bus
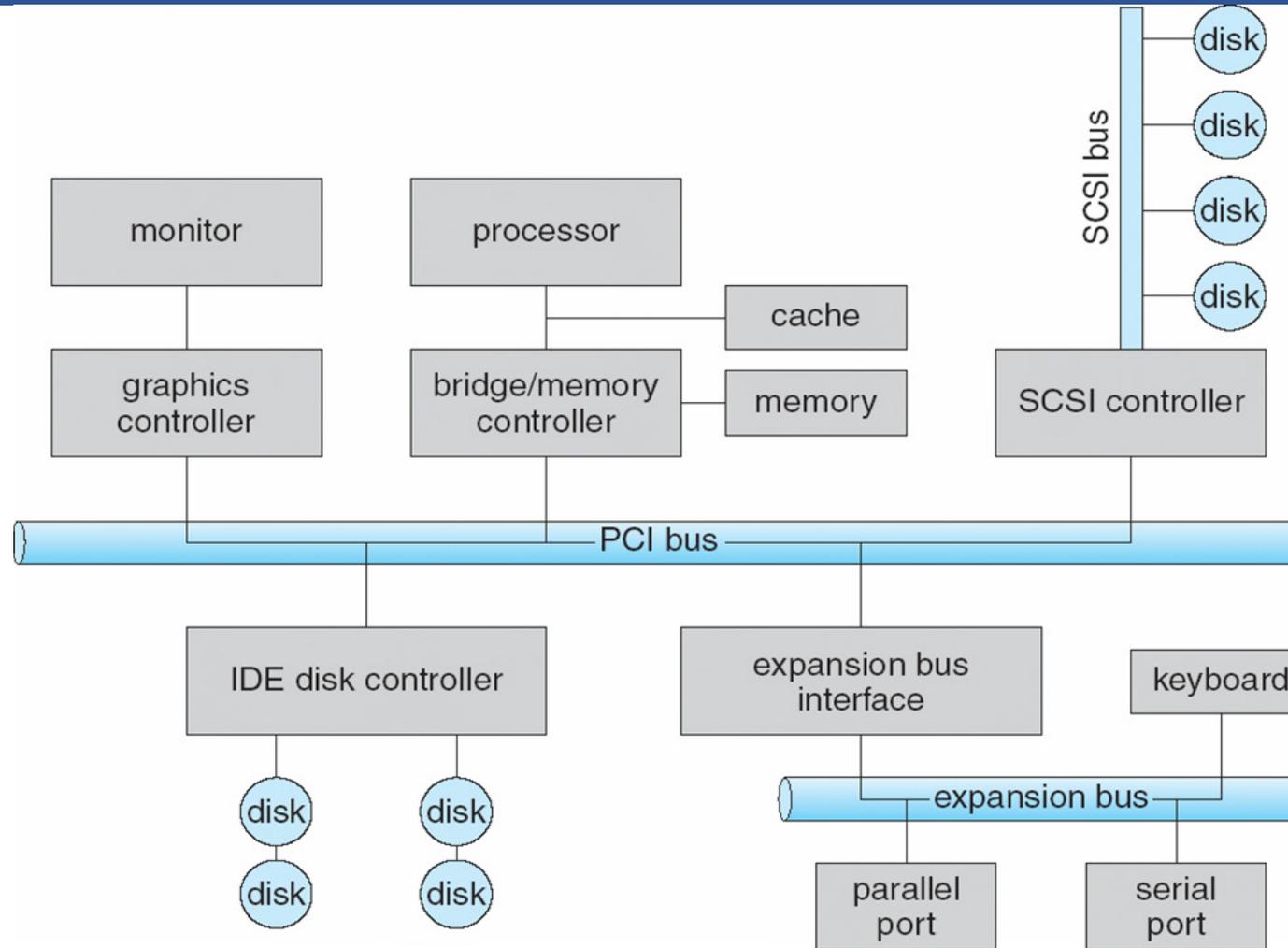(e.g., SCSI, SATA, USB)

# Buses

- **Buses**
  - Data paths that provided to enable information between CPU(s), RAM, and I/O devices.
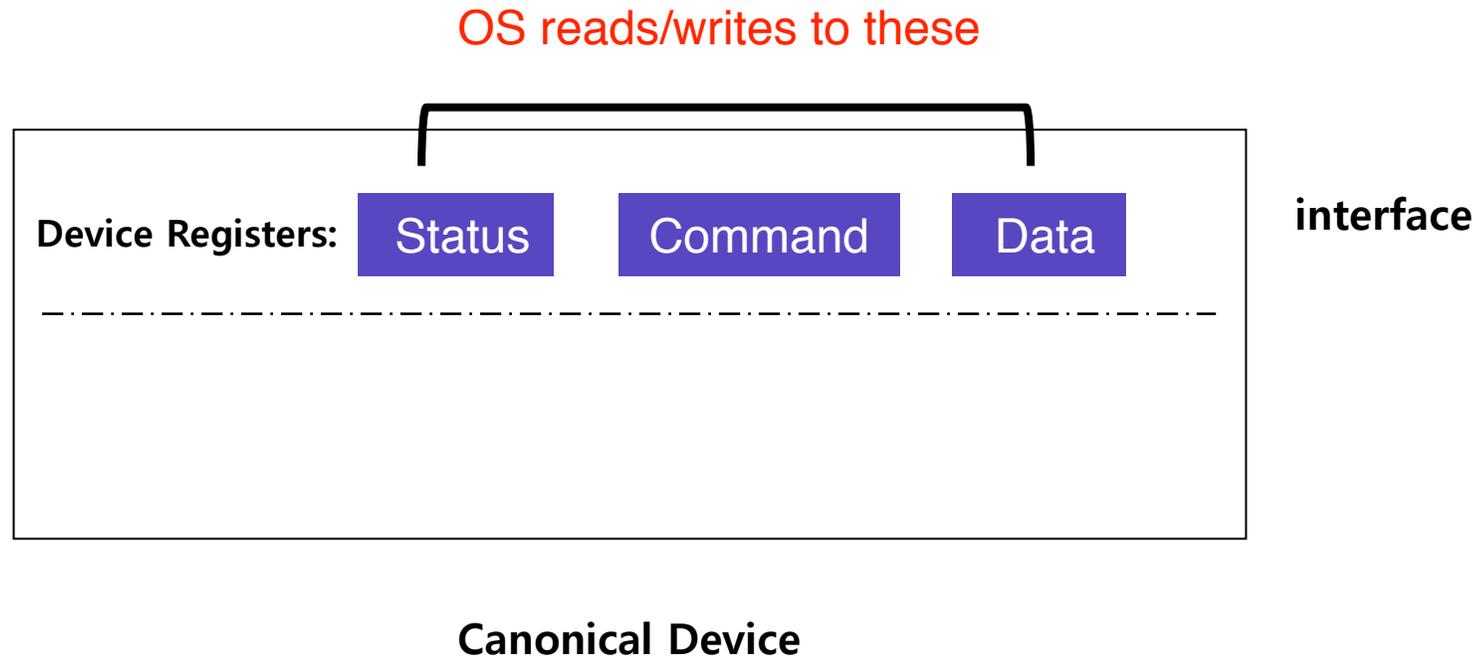
- **I/O bus**
  - Data path that connects a CPU to an I/O device.
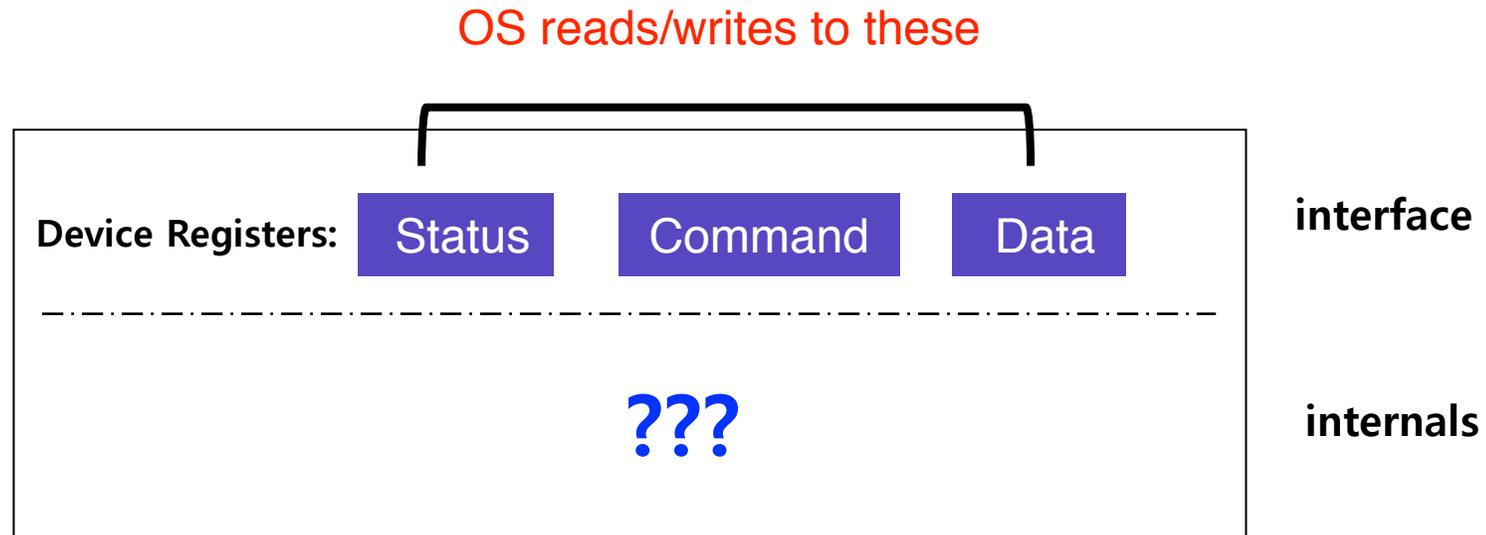  - I/O bus is connected to I/O device by three hardware components: I/O ports, interfaces and device controllers.

# What Is I/O Bus? E.g., PCI

# Canonical Device

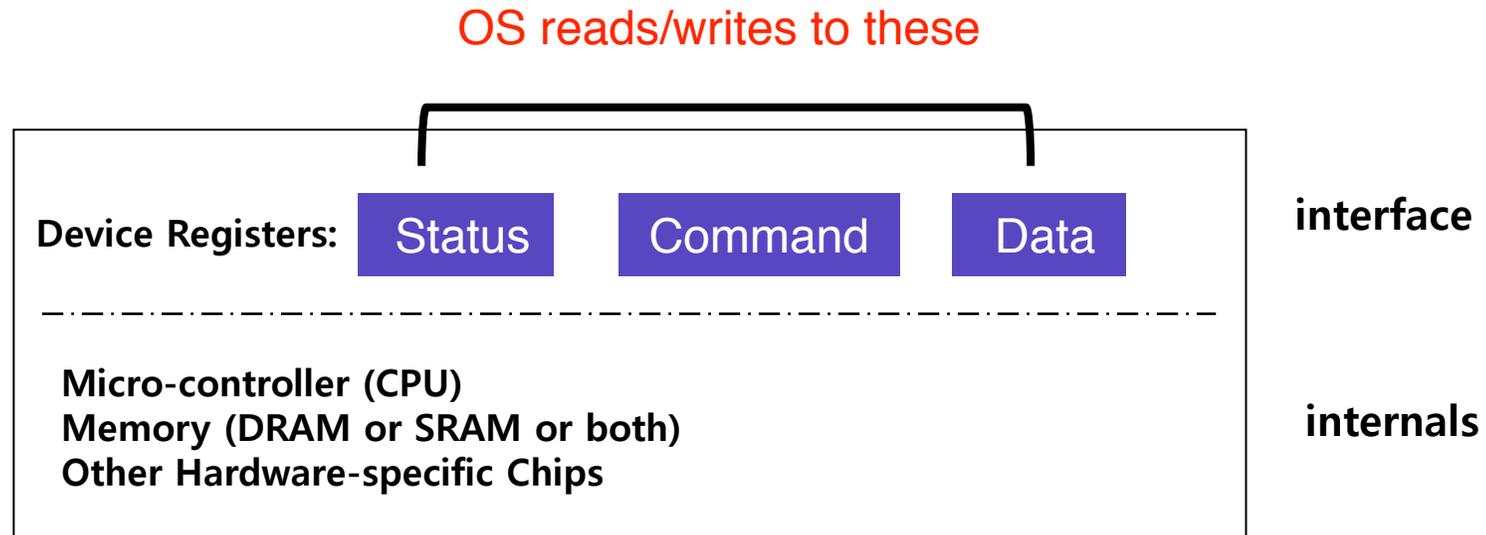OS reads/writes to these

Device Registers:  Status   Command   Data     **interface**

**Canonical Device**

# Canonical Device

OS reads/writes to these

Device Registers: Status | Command | Data

interface

??? 

internals

**Canonical Device**

# Canonical Device

OS reads/writes to these

Device Registers: **Status** **Command** **Data**  interface

----·--·--·--·--·--·--·--·--·--·--·--·--·--·--·--·--·--·

**Micro-controller (CPU)**
**Memory (DRAM or SRAM or both)**
**Other Hardware-specific Chips**

internals

**Canonical Device**

# Hardware Interface Of Canonical Device

- **status register**
  - See the current status of the device

- **command register**
  - Tell the device to perform a certain task

- **data register**
  - Pass data to the device, or get data from the device

- **By reading or writing the above three registers, the OS controls device behavior.**

# Hardware Interface Of Canonical Device

- **Typical interaction example**

```
while ( STATUS == BUSY)

    ; //wait until device is not busy

write data to data register

write command to command register

    Doing so starts the device and executes the command

while ( STATUS == BUSY)

    ; //wait until device is done with your request
```

# Polling

- **OS waits until the device is ready by repeatedly reading the status register**
  - Positive aspect is simple and working.
  - **However, it wastes CPU time just waiting for the device.**
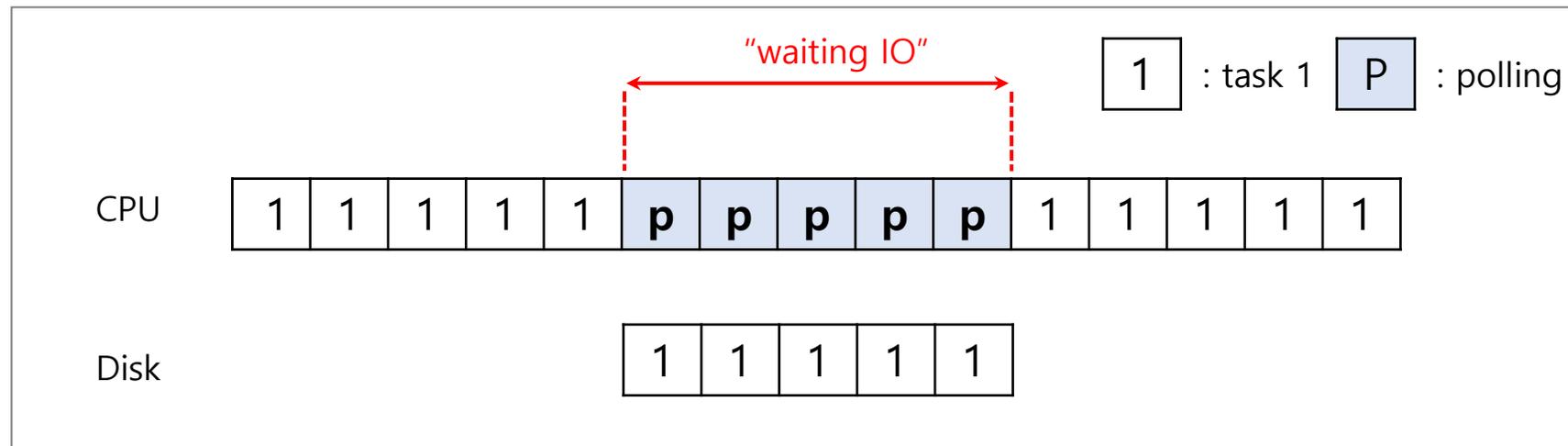    - Switching to another ready process is better utilizing the CPU.



**Diagram of CPU utilization by polling**

# Interrupts

- **Put the I/O request process to sleep and context switch to another**

- **When the device is finished, wake the process waiting for the I/O by interrupt**
  - Positive aspect is allow to **CPU and the disk are properly utilized.**

| 1 | : task 1 | 2 | : task 2 |
|---|----------|---|----------|

| CPU | 1 | 1 | 1 | 1 | 1 | 2 | 2 | 2 | 2 | 2 | 1 | 1 | 1 | 1 | 1 |

| Disk | 1 | 1 | 1 | 1 | 1 |

**Diagram of CPU utilization by interrupt**

# Polling vs Interrupts

- *However,* "**interrupts is not always the best solution**"
    - If, device performs very quickly, interrupt will "slow down" the system.

- **E.g., high network packet arrival rate**
    - Packets can arrive faster than OS can process them
    - Interrupts are very expensive (context switch)
    - Interrupt handlers have high priority
    - In worst case, can spend 100% of time in interrupt handler and never make any progress – *receive livelock*
    - Best: Adaptive switching between interrupts and polling

> If a device is fast → **poll** is best.
> If it is slow → **interrupt**s is better.

# One More Problem: Data Copying

- **CPU wastes a lot of time in copying *a large chunk of data* from memory to the device.**
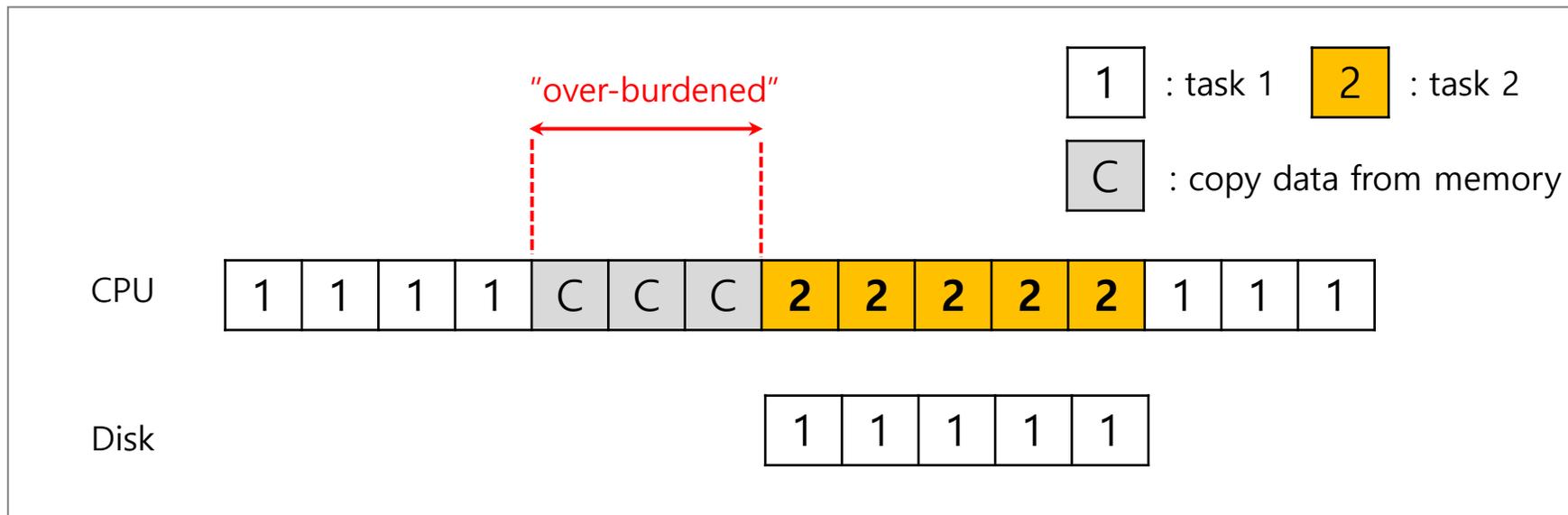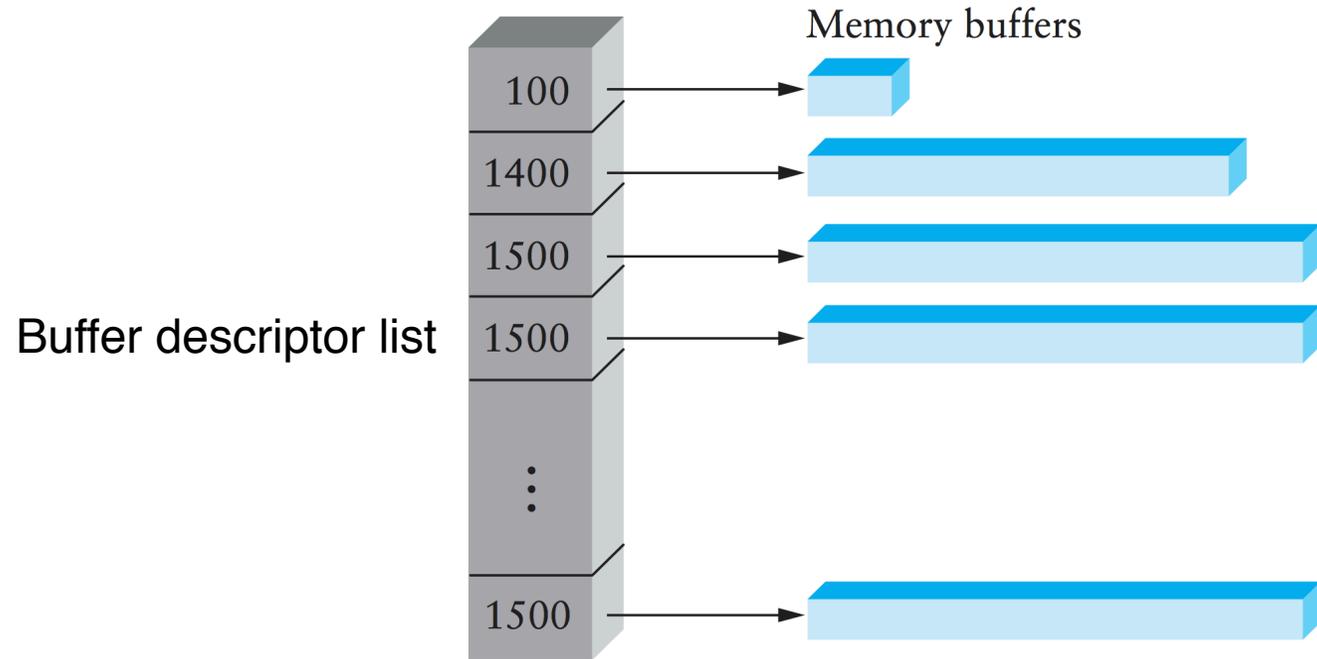


**Diagram of CPU utilization**

# DMA (Direct Memory Access)

Memory buffers

Buffer descriptor list

| 100 |
| 1400 |
| 1500 |
| 1500 |
| ⋮ |
| 1500 |

- **Idea: only use CPU to transfer control requests, not data**

- **Include list of buffer locations in main memory**
  - Device reads list and accesses buffers through DMA
  - Descriptions sometimes allow for scatter/gather I/O

# DMA (Direct Memory Access) Cont.

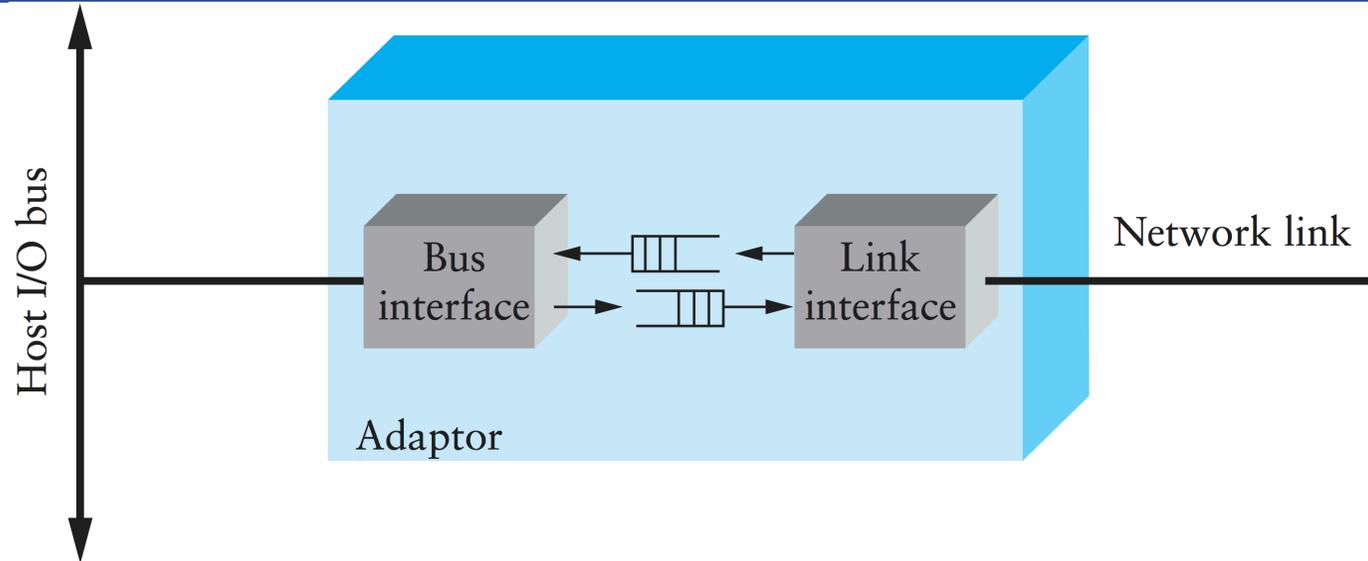- **When completed, DMA raises an interrupt, I/O begins on Disk.**



**Diagram of CPU utilization by DMA**

# Example: Network Interface Card



Host I/O bus

Bus interface

Link interface

Network link

Adaptor

- **Link interface talks to wire/fiber/antenna**
  - Typically does framing, link-layer CRC

- **FIFOs on card provide small amount of buffering**

- **Bus interface logic uses DMA to move packets to and from buffers in main memory**

# Device Interaction

- **How the OS communicates with the device?**

- **Solutions**
  - I/O instructions: a way for the OS to send data to specific device registers.
    - Ex) `in` and `out` instructions on x86
  - memory-mapped I/O
    - Device registers available as if they were memory locations.
    - The OS `load` (to read) or `store` (to write) to the device instead of main memory.

# x86 I/O instructions

```c
static inline uint8_t inb (uint16_t port)
{
  uint8_t data;
  asm volatile ("inb %w1, %b0" : "=a" (data) : "Nd" (port));
  return data;
}

static inline void outb (uint16_t port, uint8_t data)
{
  asm volatile ("outb %b0, %w1" : : "a" (data), "Nd" (port));
}

static inline void insw (uint16_t port, void *addr, size_t cnt)
{
  asm volatile ("rep insw" : "+D" (addr), "+c" (cnt)
                : "d" (port) : "memory");
}
```

**CS 318 – Lecture 14 – I/O & Disks**   **Pintos** `threads/io.h`

# IDE Disk Driver

```c
void IDE_ReadSector(int disk, int off,
      void *buf)
{
  // Select Drive
  outb(0x1F6, disk == 0 ? 0xE0 : 0xF0);
  IDEWait();
  // Read length (1 sector = 512 B)
  outb(0x1F2, 1);
  outb(0x1F3, off); // LBA low
  outb(0x1F4, off >> 8); // LBA mid
  outb(0x1F5, off >> 16); // LBA high
  outb(0x1F7, 0x20); // Read command
  insw(0x1F0, buf, 256); // Read 256 words
}
```

```c
void IDEWait()
{
  // Discard status 4 times
  inb(0x1F7); inb(0x1F7);
  inb(0x1F7); inb(0x1F7);
  // Wait for status BUSY flag to clear
  while ((inb(0x1F7) & 0x80) != 0);
}
```

# Memory-mapped IO

- `in/out` **instructions slow and clunky**
  - Instruction format restricts what registers you can use
  - Only allows $2^{16}$ different port numbers
  - Per-port access control turns out not to be useful (any port access allows you to disable all interrupts)

- **Devices can achieve same effect with physical addresses, e.g.:**

  ```
  volatile int32_t *device_control
      = (int32_t *) (0xc0100 + PHYS_BASE);
  *device_control = 0x80;
  int32_t status = *device_control;
  ```

  - OS must map physical to virtual addresses, ensure non-cachable

# Protocol Variants

Device Registers:   Status    Command    Data

Micro-controller (CPU)
Memory (DRAM or SRAM or both)
Other Hardware-specific Chips

- **Status checks: polling vs. interrupts**

- **Data: PIO vs. DMA**

- **Control: special instructions vs. memory-mapped I/O**

# Variety Is a Challenge

- **Problem:**
  - many, many devices
  - each has its own protocol

- **How can we avoid writing a slightly different OS for each H/W combination?**

# Variety Is a Challenge

- **Problem:**
  - many, many devices
  - each has its own protocol

- **How can we avoid writing a slightly different OS for each H/W combination?**

- **Solution: Abstraction!**
  - Build a common interface
  - Write device driver for each device
  - Drivers are 70% of Linux source code

# File System Abstraction

- **File system specifics of which disk class it is using.**
  - Ex) It issues **block read** and **write** request to the generic block layer.



| | |
|---|---|
| Application | user |
| POSIX API [open, read, write, close, etc] | |
| File System | kernel |
| Generic Block Interface [block read/write] | |
| Generic Block Layer | |
| Specific Block Interface [protocol-specific read/write] | |
| Device Driver [SCSI, ATA, etc] | **The File System Stack** |
| **Hard Drive** | |

# Hard Disks

# Hard Disks

# Hard Disks

# Basic Interface

- **Disk interface presents linear array of sectors**
  - Historically <span style="color:red">512 Bytes</span>
  - Written atomically (even if there is a power failure)
  - 4 KiB in "advanced format" disks
    - Torn write: If an untimely power loss occurs, only a portion of a larger write may complete

- **Disk maps logical sector #s to physical sectors**

- **OS doesn't know logical to physical sector mapping**

# Basic Geometry



- **Platter (Aluminum coated with a thin magnetic layer)**
  - A circular hard surface
  - Data is stored persistently by inducing magnetic changes to it.
  - Each platter has 2 sides, each of which is called a **surface**.

# Basic Geometry (Cont.)

- **Spindle**
  - Spindle is connected to a motor that spins the platters around.
  - The rate of rotations is measured in **RPM** (Rotations Per Minute).
    - Typical modern values : 7,200 RPM to 15,000 RPM.

- **Track**
  - Concentric circles of sectors
  - Data is encoded on each surface in a track.
  - A single surface contains many thousands and thousands of tracks.

- **Cylinder**
  - A stack of tracks of fixed radius
  - Heads record and sense data along cylinders
  - Generally only one head active at a time

# Cylinders, Tracks, & Sectors

# A Simple Disk Drive



**Rotates this way**

8  9
7        10
**head**  6    **spindle**    11
5        0
4        1
3  2

**A Single Track Plus A Head**

- **Disk head (One head per surface of the drive)**
  - The process of *reading* and *writing* is accomplished by the **disk head**.
  - Attached to a single disk arm, which moves across the surface.

# Single-track Latency: The Rotational Delay



**Rotates this way**

8  9
7        10
**head**  6   **spindle**  11
**arm**  5        0
4        1
3  2

**A Single Track Plus A Head**

- **Rotational delay: Time for the desired sector to rotate**
  - Ex) Full rotational delay is `R` and we start at `sector 6`
    - Read sector 0: Rotational delay = $\frac{R}{2}$
    - Read sector 5: Rotational delay = `R-1` (worst case.)

# Multiple Tracks

- **Let's Read 12!**

# Multiple Tracks: Seek To Right Track

- **Let's Read 12!**

# Multiple Tracks: Seek To Right Track

- **Let's Read 12!**

# Multiple Tracks: Seek To Right Track
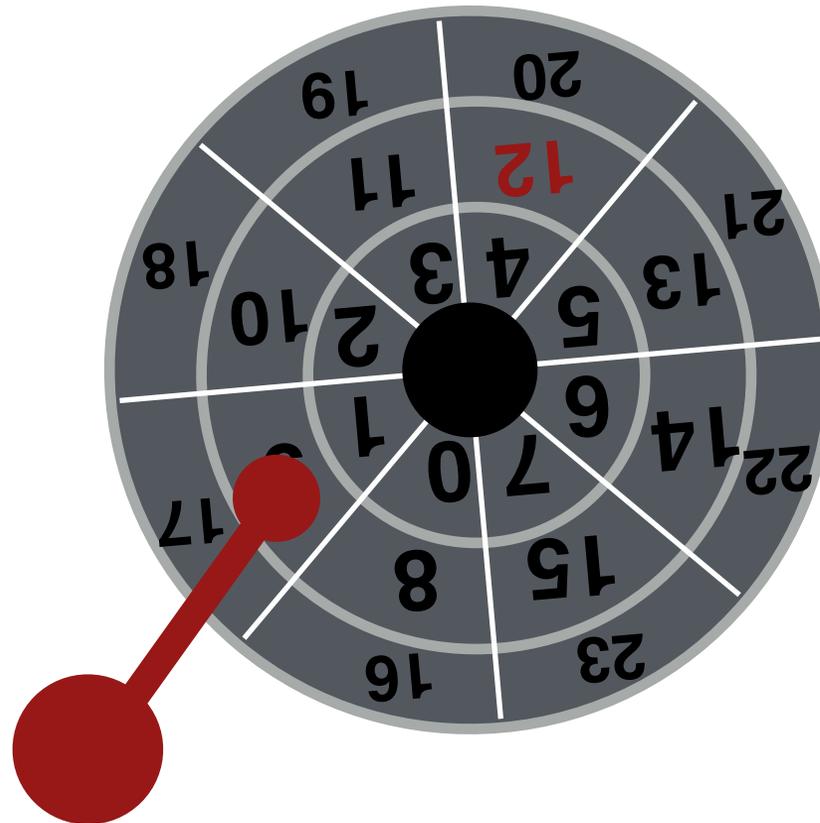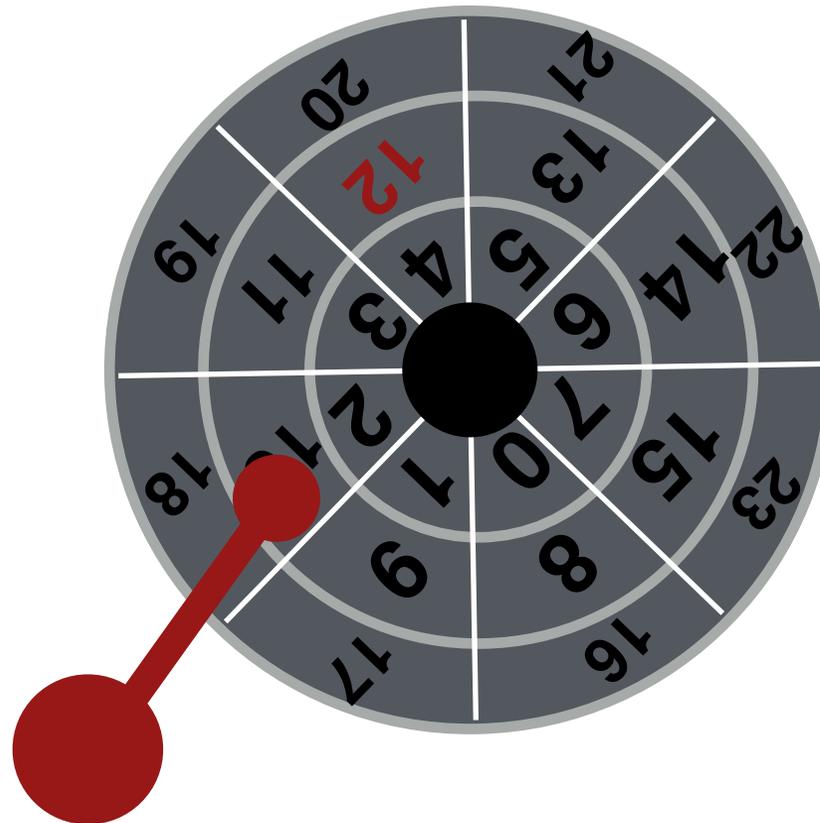
- **Let's Read 12!**

# Multiple Tracks: Wait for Rotation

- **Let's Read 12!**

# Multiple Tracks: Wait for Rotation

- **Let's Read 12!**

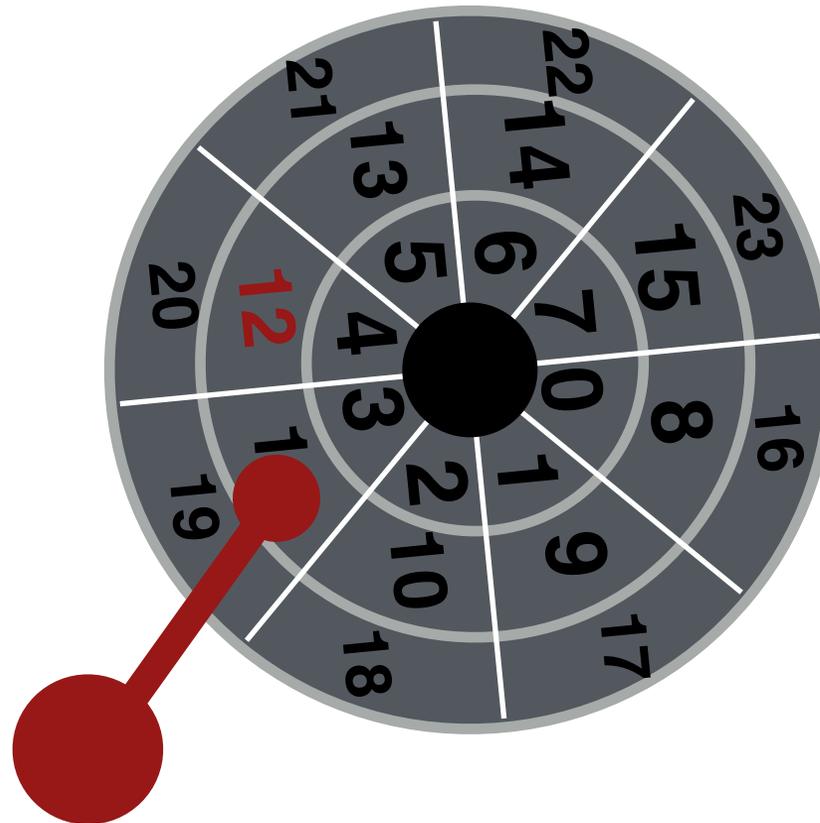# Multiple Tracks: Wait for Rotation

- **Let's Read 12!**

# Multiple Tracks: Wait for Rotation

- **Let's Read 12!**

# Multiple Tracks: Wait for Rotation
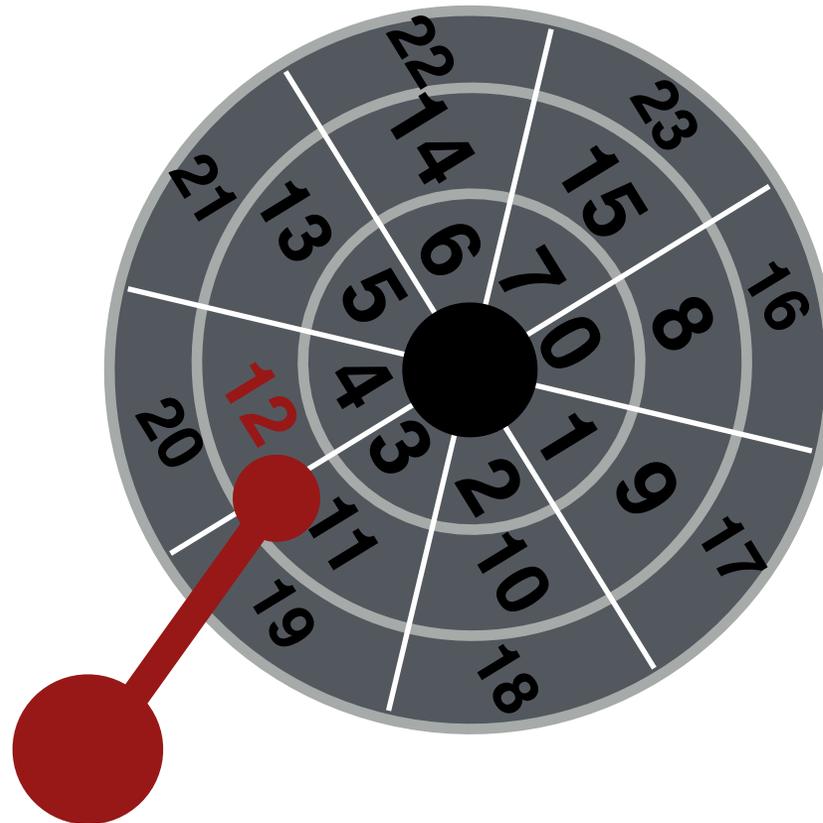
- **Let's Read 12!**

# Multiple Tracks: Wait for Rotation

- **Let's Read 12!**
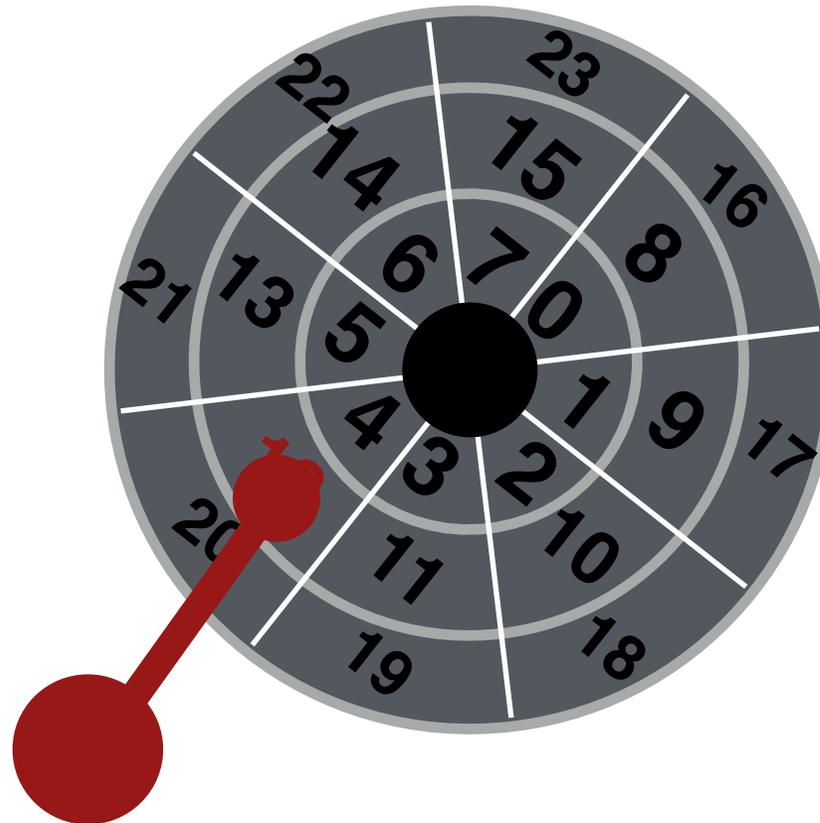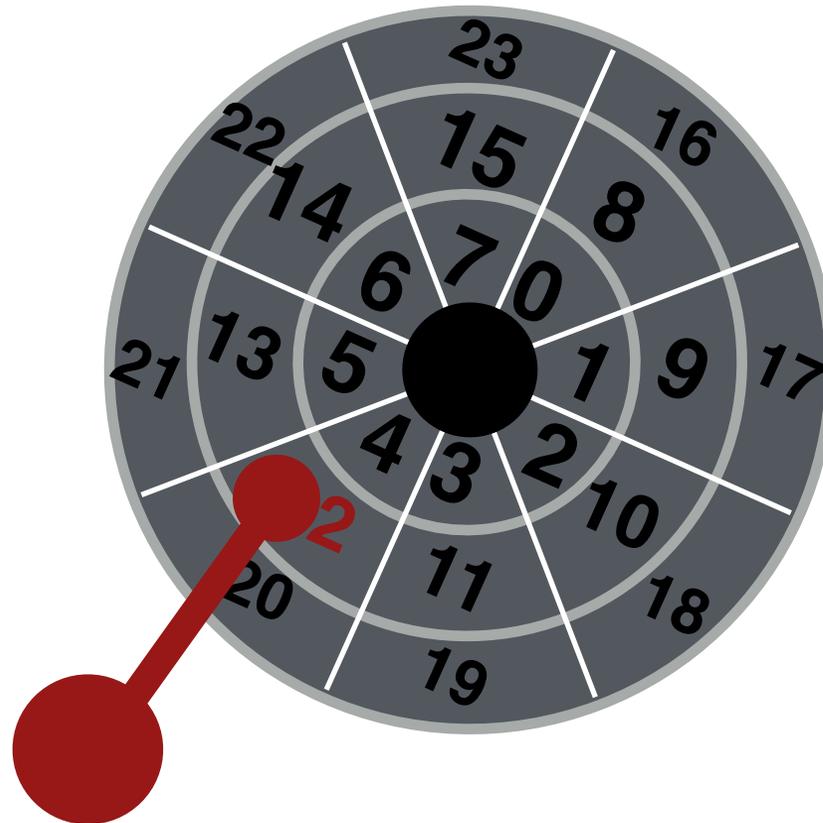
# Multiple Tracks: Transfer Data

- **Let's Read 12!**

# Multiple Tracks: Transfer Data
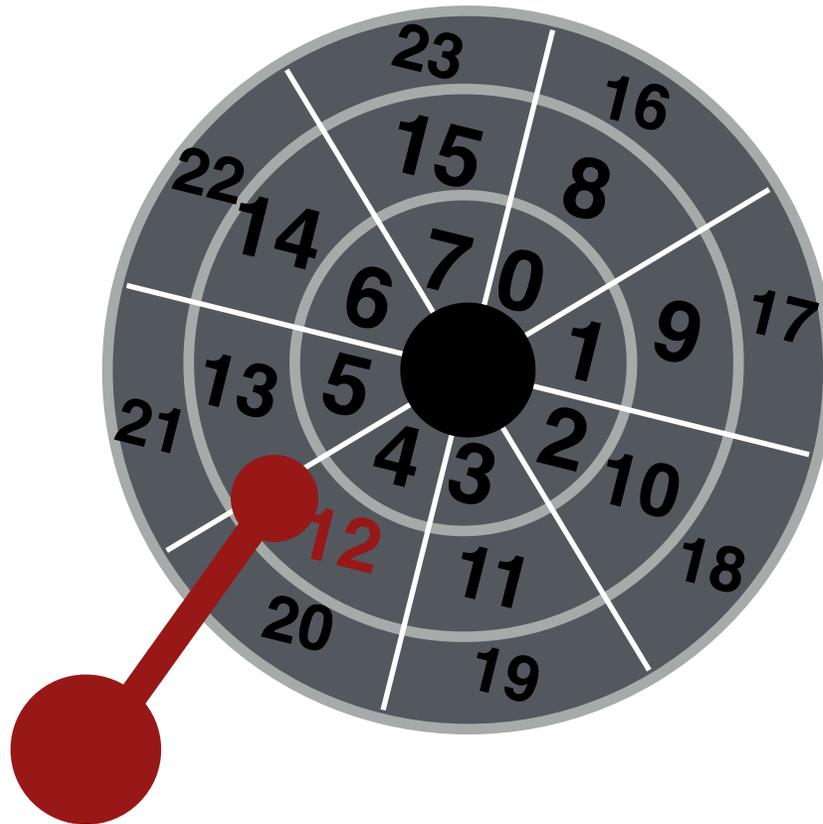
- **Let's Read 12!**
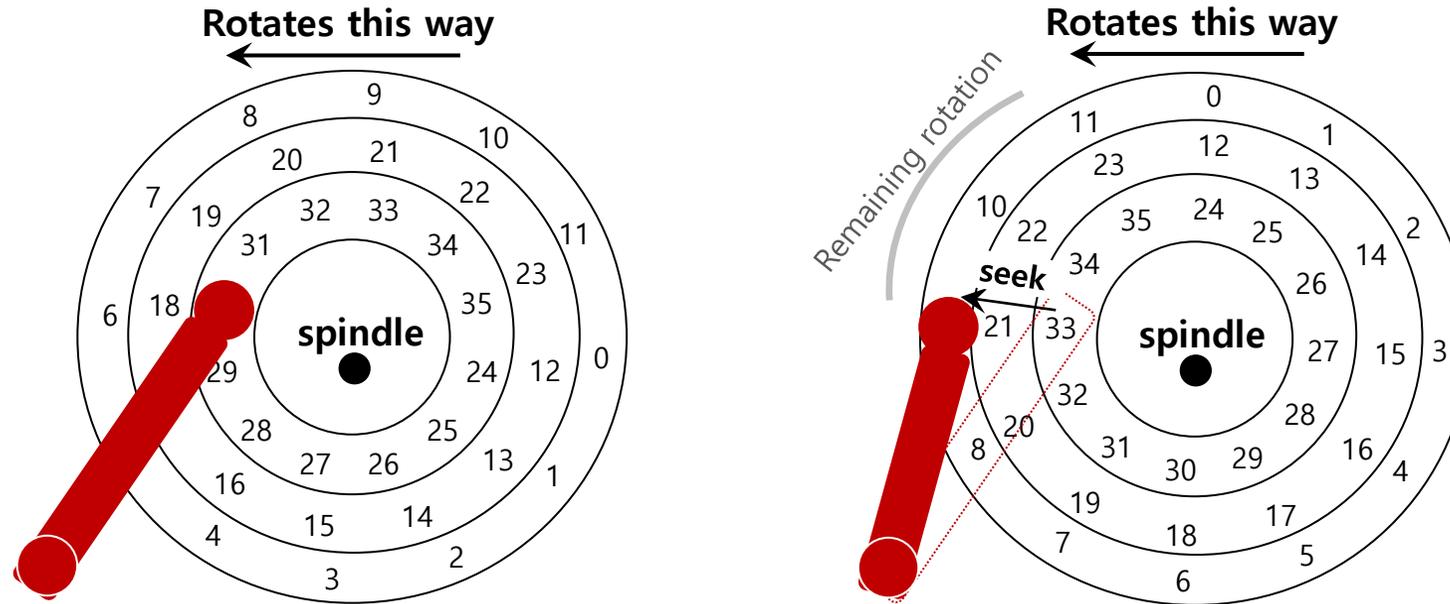
# Multiple Tracks: Transfer Data

- **Let's Read 12!**

# Yay!

# Multiple Tracks: Seek Time



- **Seek: Move the disk arm to the correct track**
  - **Seek time**: Time to move head to the track contain the desired sector.
  - One of the most costly disk operations.

# Seek, Rotate, Transfer

- **Acceleration → Coasting → Deceleration → Settling**
  - **Acceleration**: The disk arm gets moving.
  - **Coasting**: The arm is moving at full speed.
  - **Deceleration**: The arm slows down.
  - **Settling**: The head is *carefully positioned* over the correct track.

- **Seeks often take several milliseconds!**
  - settling alone can take 0.5 to 2ms.
  - entire seek often takes 4 - 10 ms.

# Seek, Rotate, Transfer

- **Depends on rotations per minute (RPM)**
  - 7200 RPM is common, 1500 RPM is high end.

- **With 7200 RPM, how long to rotate around?**
  - 1 / 7200 RPM = 1 minute / 7200 rotations = 1 second / 120 rotations = 8.3 ms / rotation
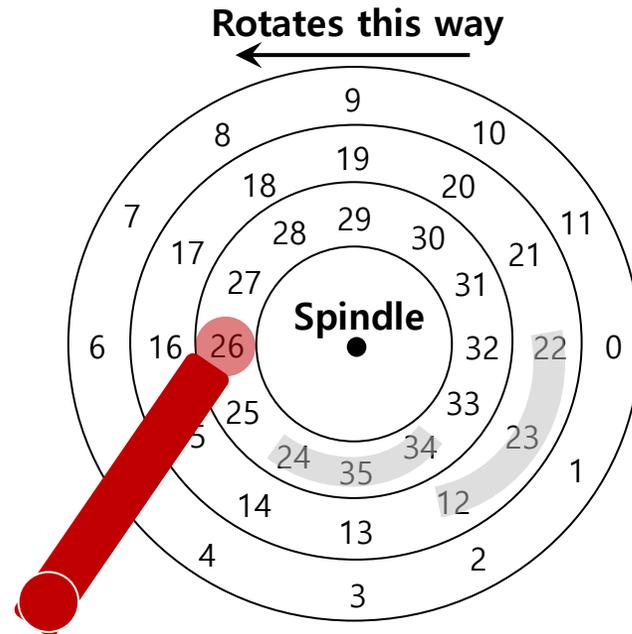
- **Average rotation?**
  - 8.3 ms / 2 = 4.15 ms

# Seek, Rotate, Transfer

- **The final phase of I/O**
  - Data is either *read from* or *written* to the surface.

- **Pretty fast — depends on RPM and ector density**

- **100+ MB/s is typical for maximum transfer rate**

- **How long to transfer 512-bytes?**
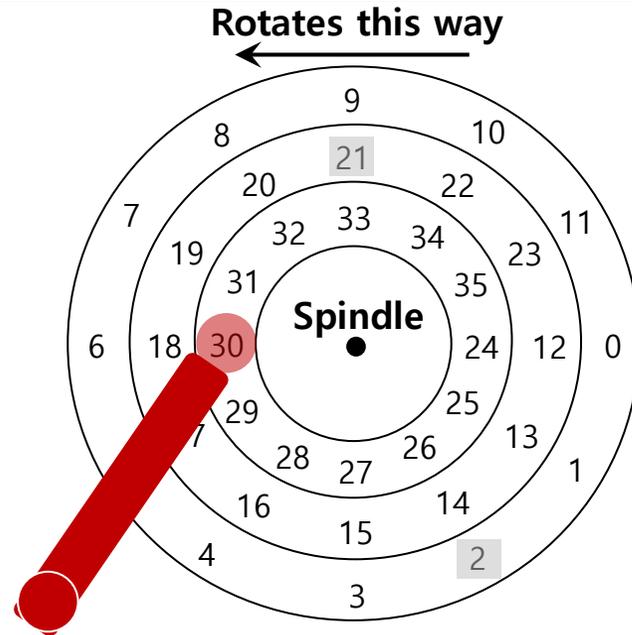  - 512 bytes * (1s / 100 MB) = 5 $\mu s$

# Workload

- **So…**
  - seeks are slow
  - rotations are slow
  - transfers are fast

- **What kind of workload is fastest for disks?**
  - **Sequential**: access sectors in order (transfer dominated)
  - **Random**: access sectors arbitrarily (seek+rotation dominated)

# Track Skew



**Rotates this way**

9
8    10
    19
18      20
7          11
17  28  29  30  21
27          31
16  **Spindle**  32  0
6   26   •   22
25          33
             23
24  35  34  12
14          1
4    13    2
      3

**Three Tracks: Track Skew Of 2**

- **Make sure that sequential reads can be properly serviced even when crossing track boundaries**
  - *Without track skew*, the head would be moved to the next track but the desired next block would have already rotated under the head.

# Disk Scheduling



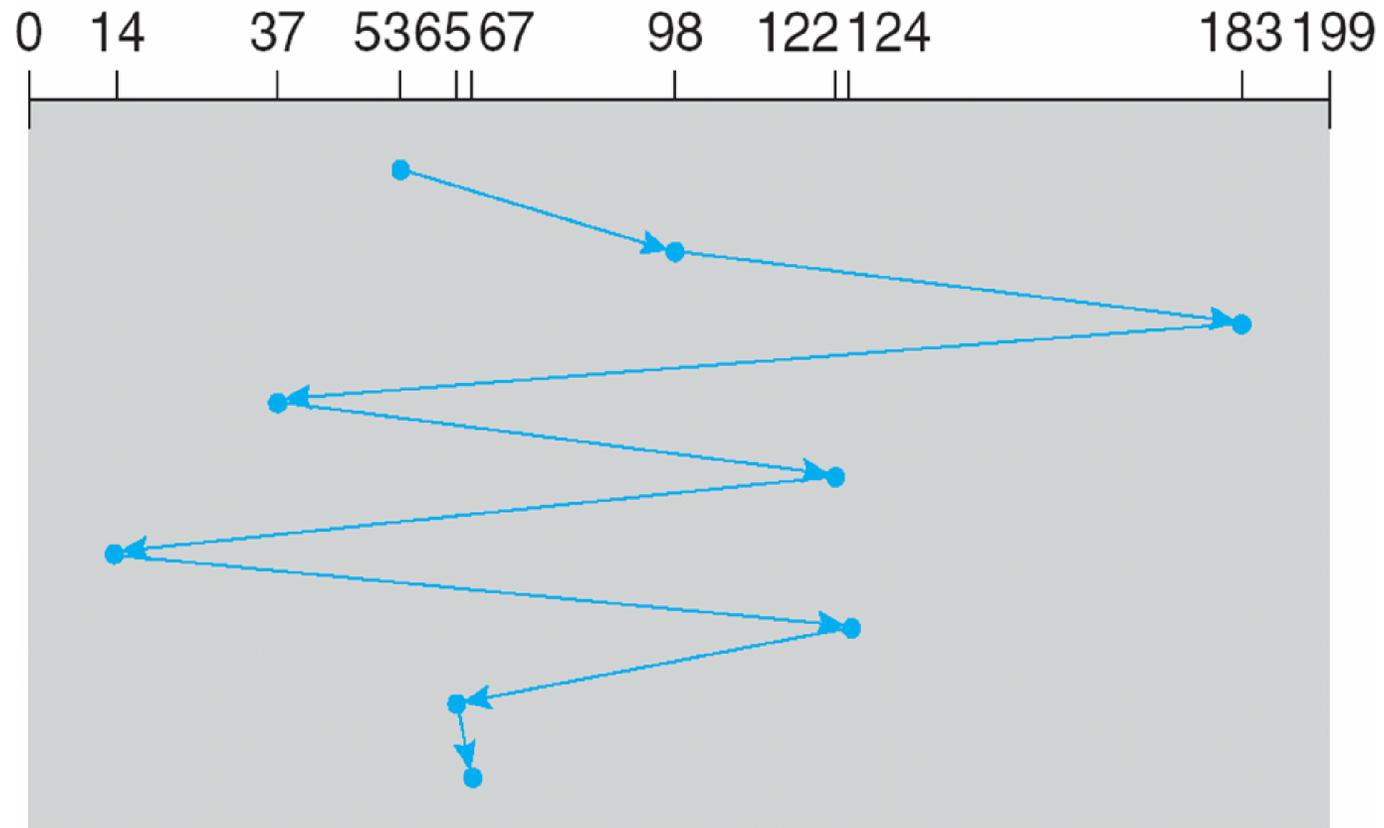- **Disk Scheduler decides which I/O request to schedule next**

# Disk Scheduling: FCFS

- **"First Come First Served"**
  - Process disk requests in the order they are received

- **Advantages**
  - Easy to implement
  - Good fairness

- **Disadvantages**
  - Cannot exploit request locality
  - Increases average latency, decreasing throughput

# FCFS Example



queue = 98, 183, 37, 122, 14, 124, 65, 67
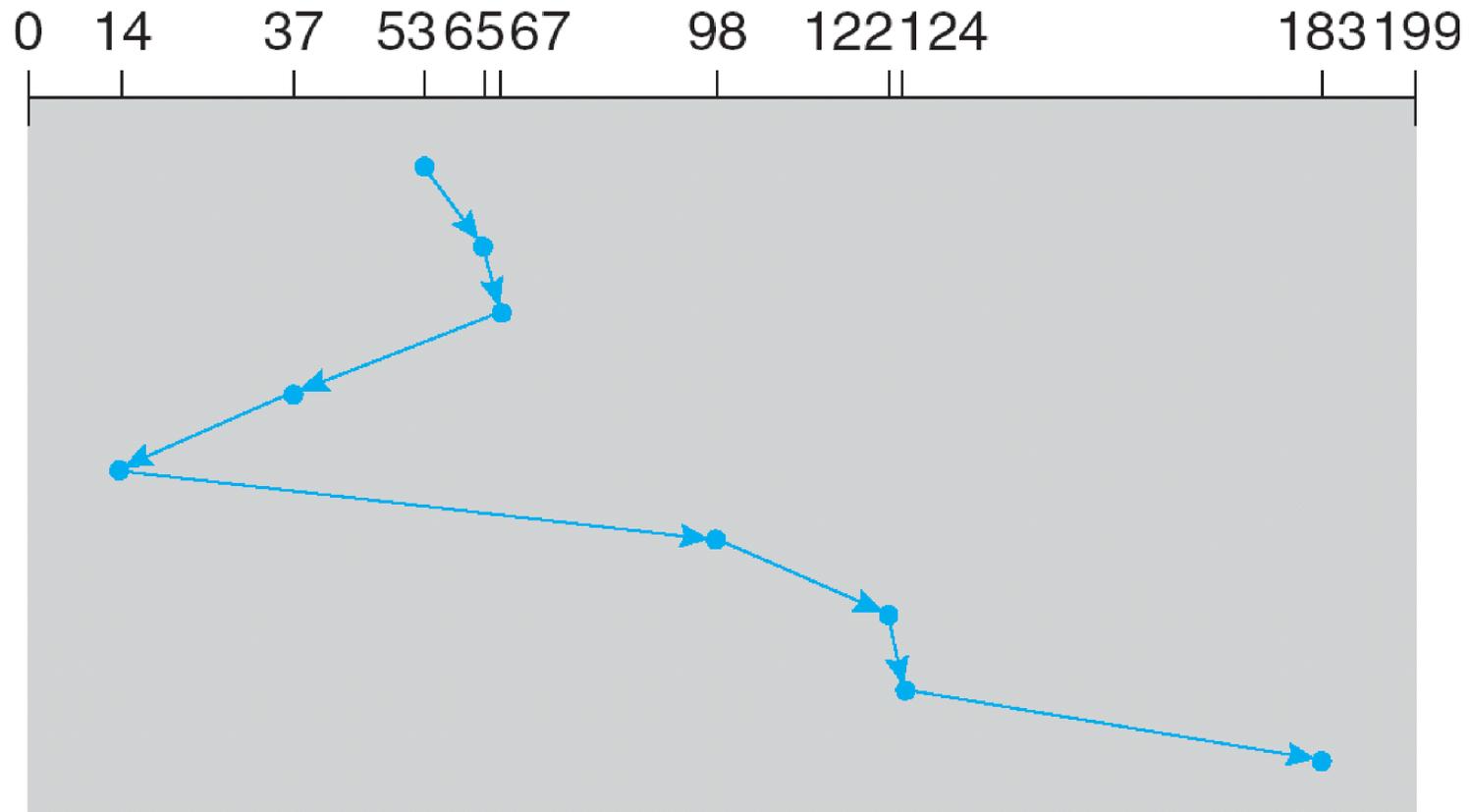head starts at 53

# SSTF (Shortest Seek Time First)

- **Order the queue of I/O request by track**

- **Pick requests on the nearest track to complete first**
  - Also called shortest positioning time first (SPTF)

- **Advantages**
  - Exploits locality of disk requests
  - Higher throughput

- **Disadvantages**
  - Starvation
  - Don't always know what request will be fastest

# SSTF Example



queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

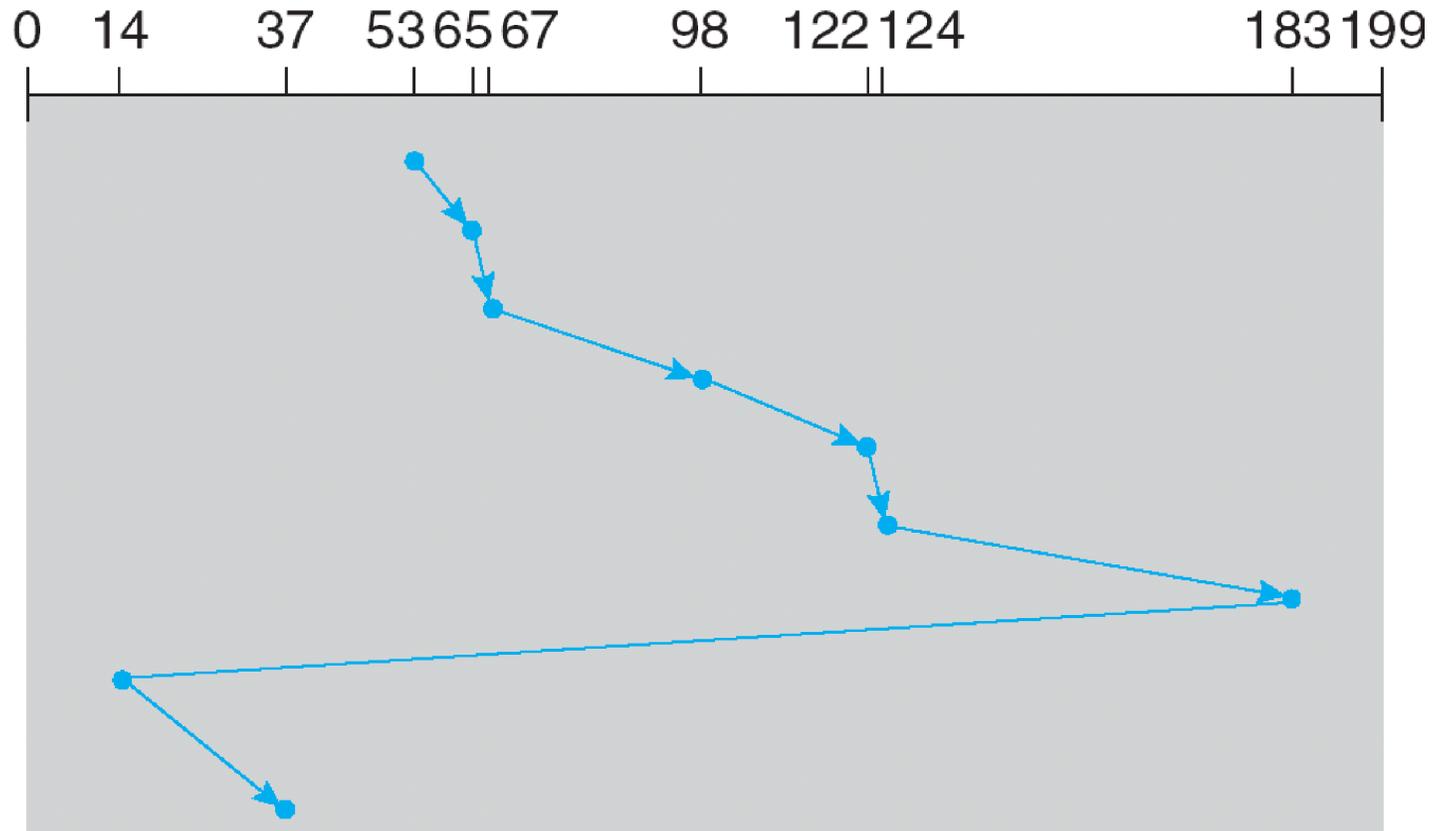0   14        37   5365 67        98   122 124                    183 199

# "Elevator" Scheduling (SCAN)

- **Sweep across disk, servicing all requests passed**
  - Like SSTF, but next seek must be in same direction
  - Switch directions only if no further requests

- **Advantages**
  - Takes advantage of locality
  - Bounded waiting

- **Disadvantages**
  - Cylinders in the middle get better service
  - Might miss locality SSTF could exploit

- **CSCAN: Only sweep in one direction**
  - **Very commonly used algorithm in Unix**

# CSCAN example

queue     98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

# Flash Memory

- **Today, people increasingly using flash memory**

- **Completely solid state (no moving parts)**
  - Remembers data by storing charge
  - Lower power consumption and heat
  - No mechanical seek times to worry about

- **Limited # overwrites possible**
  - Blocks wear out after 10,000 (MLC) – 100,000 (SLC) erases
  - Requires flash translation layer (FTL) to provide wear leveling, so repeated writes to logical block don't wear out physical block
  - FTL can seriously impact performance

- **Limited durability**
  - Charge wears out over time
  - Turn off device for a year, you can potentially lose data!

# Next Time…

- **Read Chapter 39, 40**