# CS162
# Operating Systems and
# Systems Programming
# Lecture 17

## General I/O, Storage Devices
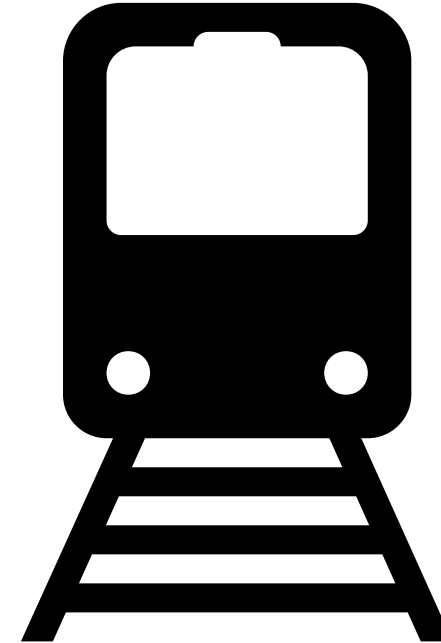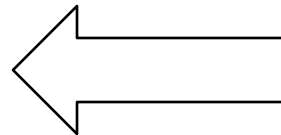
Professor Natacha Crooks & Matei Zaharia
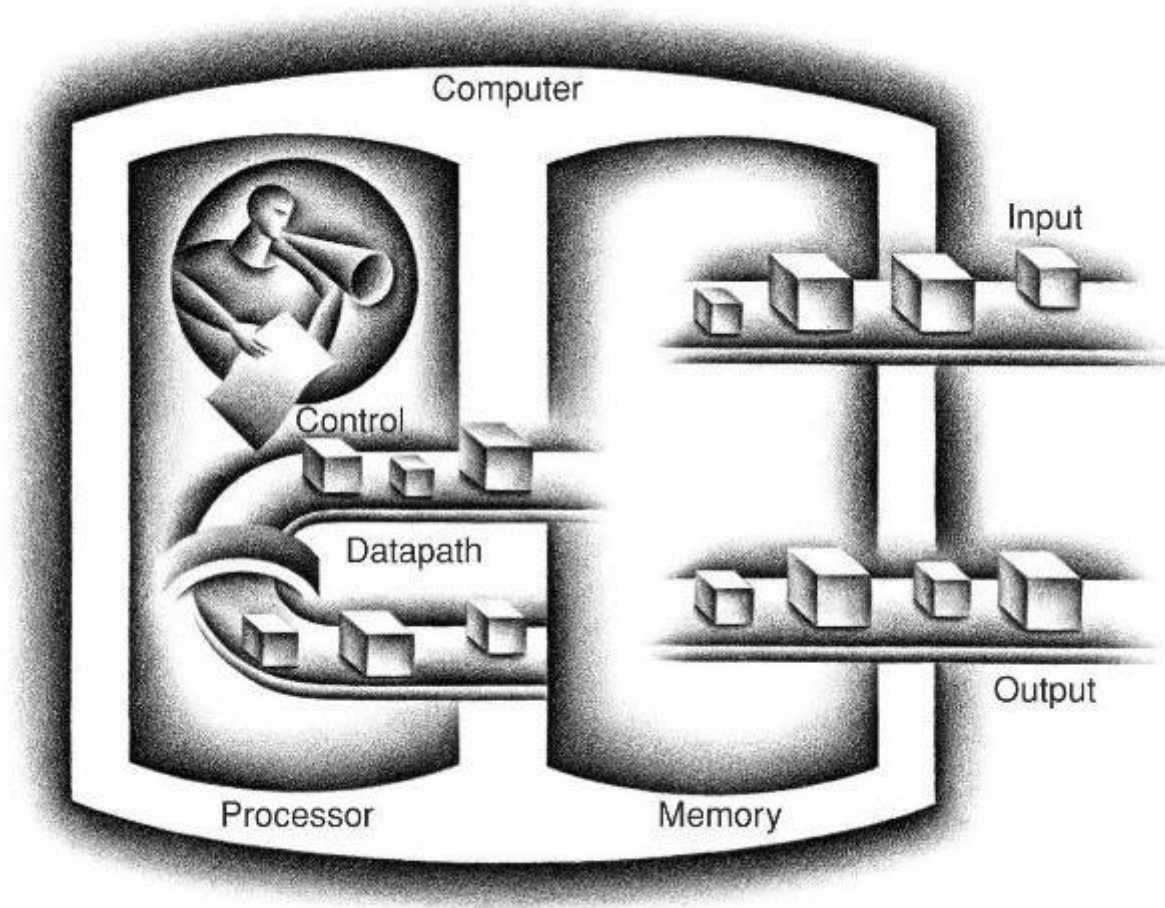
https://cs162.org/

# Course Map

- Introduction

- OS Concepts

- Concurrency

- Scheduling

- Memory Management

- Devices and file systems ⬅

- Reliability, networking and cloud

# Recall: Five Components of a Computer
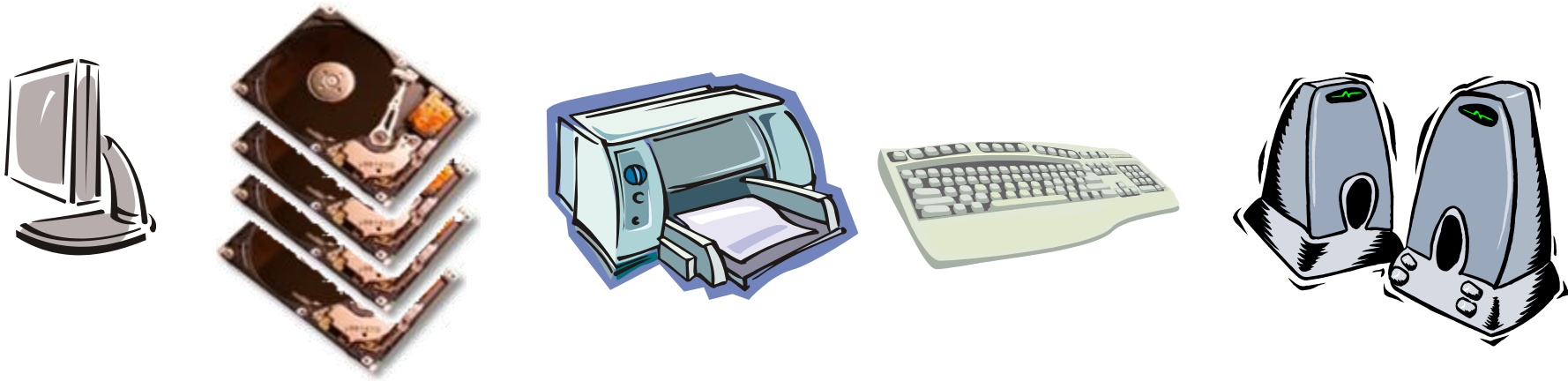


From "Computer Organization and Design" by Hennesy & Patterson

# CPU: You need to get out more!

Input/output is the mechanism through which the
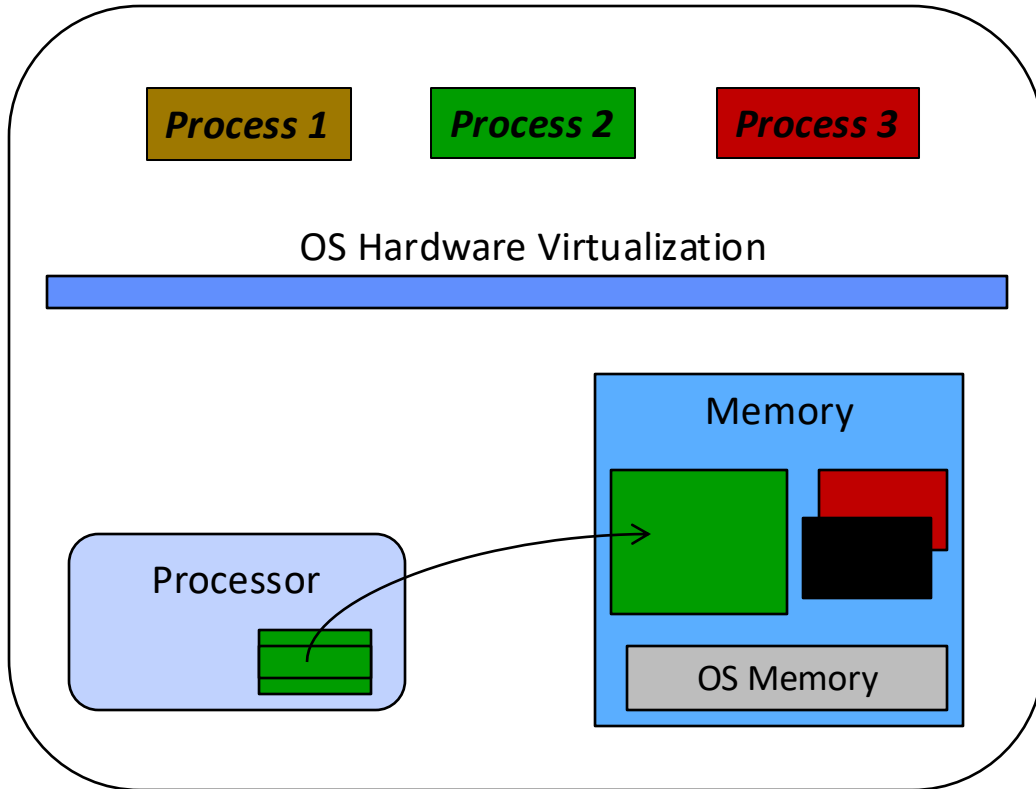
computer communicates with the outside world
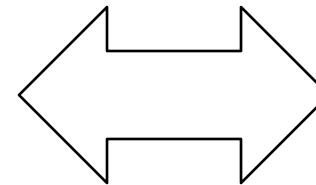
# Want Standard Interfaces to Devices

- **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

- **Character Devices:** *e.g.* keyboards, mice, serial ports, some USB devices
  - Single characters at a time
  - Commands include `get()`, `put()`

- **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 polling and connection management functionality
  - Usage: pipes, FIFOs, streams, queues, mailboxes
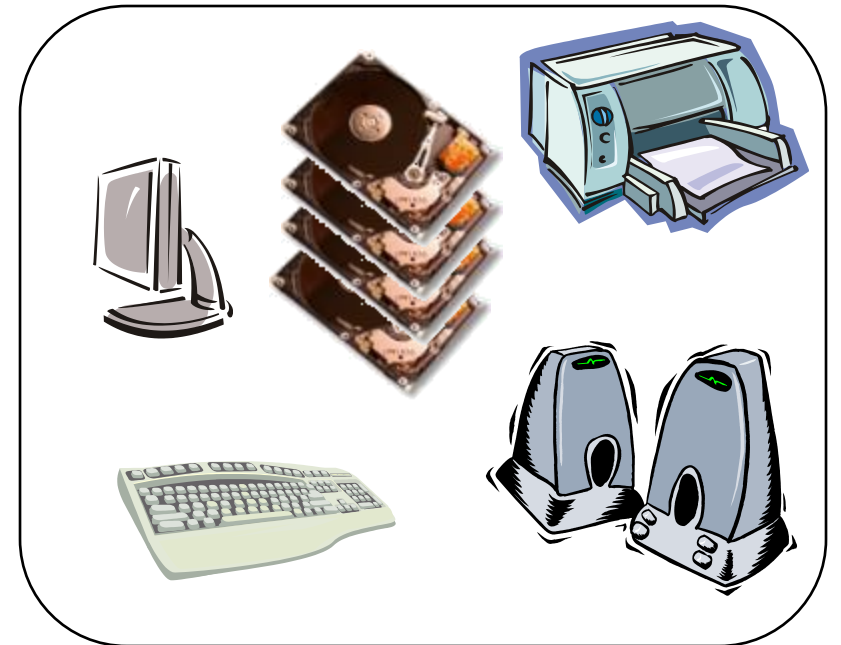
# IO Subsystem: Abstraction, abstraction, abstraction

**Process 1**    **Process 2**    **Process 3**

OS Hardware Virtualization

Memory

Processor

OS Memory

IO Layer

**Virtual Machine Abstraction**

**IO Devices**

# IO Subsystem: Abstraction, abstraction, abstraction

- This code works for pretty much any device

```
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 this lecture

  – Can only scratch surface!

# Requirements of I/O layer

- But… thousands of devices, each slightly different
  - » OS: How can we **standardize** the interfaces to these devices?

- Devices unreliable: media failures and transmission errors
  - » OS: How can we make them **reliable**???

- Devices unpredictable and/or slow
  - » OS: How can we **manage** them if we don't know what they will do or how they will perform?

# Simplified IO architecture



Follows a hierarchical structure because of cost: the faster the bus, the more expensive it is

# Intel's Z270 Chipset

## 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)

# Recall: 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 |

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

Device rates vary over 12 orders of magnitude!!!

# Two questions

- What is a bus?

- How does the processor talk to the devices?

# What's a bus?



- Common set of wires for communication among multiple hardware devices plus protocols for carrying out data transfer transactions

- Split into three parts: control lines, address lines, and data lines

- Protocol: initiator requests access, arbitration to grant, identification of recipient, handshake to convey address, length, data

- High bandwidth close to processor, and slower but more flexible farther out

# Why a Bus?

- Buses let us connect $n$ devices over a single set of wires, connections, and protocols
  - $O(n^2)$ relationships with 1 set of wires (!)

- Downside: Only one transaction at a time
  - The rest must wait
  - "Arbitration" aspect of bus protocol ensures the rest wait

# PCI Bus Evolution



- PCI started life out as a parallel bus (send bits on many wires in parallel)

- But a parallel bus has many limitations
  - Hard to keep all the wires sending/receiving in sync
  - Slowest devices must be able to tell what's happening (e.g., for arbitration)
    - » Bus speed is set to that of the slowest device!

# PCI Express (PCIe) "Bus"

- No longer a parallel bus

- Really a **collection of fast serial channels** or "lanes"

- Devices can use as many as they need to achieve a desired bandwidth

- Slow devices don't have to share with fast ones

- One of the successes of device abstraction in Linux was the ability to migrate from PCI to PCI Express
  - The physical interconnect changed completely, but the old API still worked

# Example PCI Architecture

# How does the Processor Talk to Devices?



- CPU interacts with a *Device Controller*
  - Contains a set of *registers* that can be read and written
  - May contain memory for request queues, etc.
- Processor accesses registers in two ways:
  - Port-Mapped I/O: 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

# How does the processor talk to devices?

- Remember, it's all about abstractions!

Device Controller

Interface
(What the OS sees)

| Status registers | Command regs. | Data |

Hardware interface device presents to OS

Internals
(What is needed to implement the device abstraction)

| Microcontroller | Memory | Other chips |

Device implementation

# How does the processor talk to devices?

Interface
(What the OS sees)

Registers

| Status | Command | Data |
|---|---|---|

**Port-Mapped I/O:**
Privileged in/out instructions

Example in Intel assembly: 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: Port-Mapped I/O in Pintos Speaker Driver

**Pintos: devices/speaker.c**

```c
13   /* Sets the PC speaker to emit a tone at the given FREQUENCY, in
14      Hz. */
15   void
16   speaker_on (int frequency)
17   {
18     if (frequency >= 20 && frequency <= 20000)
19       {
20         /* Set the timer channel that's connected to the speaker to
21            output a square wave at the given FREQUENCY, then
22            connect the timer channel output to the speaker. */
23         enum intr_level old_level = intr_disable ();
24         pit_configure_channel (2, 3, frequency);
25         outb (SPEAKER_PORT_GATE, inb (SPEAKER_PORT_GATE) | SPEAKER_GATE_ENABLE);
26         intr_set_level (old_level);
27       }
28     else
29       {
30         /* FREQUENCY is outside the range of normal human hearing.
31            Just turn off the speaker. */
32         speaker_off ();
33       }
34   }
35
36   /* Turn off the PC speaker, by disconnecting the timer channel's
37      output from the speaker. */
38   void
39   speaker_off (void)
40   {
41     enum intr_level old_level = intr_disable ();
42     outb (SPEAKER_PORT_GATE, inb (SPEAKER_PORT_GATE) & ~SPEAKER_GATE_ENABLE);
43     intr_set_level (old_level);
44   }
```

**Pintos: threads/io.h**

```c
7    /* Reads and returns a byte from PORT. */
8    static inline uint8_t
9    inb (uint16_t port)
10   {
11     /* See [IA32-v2a] "IN". */
12     uint8_t data;
13     asm volatile ("inb %w1, %b0" : "=a" (data) : "Nd" (port));
14     return data;
15   }


64   /* Writes byte DATA to PORT. */
65   static inline void
66   outb (uint16_t port, uint8_t data)
67   {
68     /* See [IA32-v2b] "OUT". */
69     asm volatile ("outb %b0, %w1" : : "a" (data), "Nd" (port));
70   }
```

# 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

**0x80020000** — Graphics Command Queue

**0x80010000** — Display Memory

**0x8000F000**

**0x0007F004** — Command

**0x0007F000** — Status

**Physical Address Space**

# A Simple Protocol for Talking to a Device

```
While (STATUS == BUSY)
    ; // wait until device is not busy
Write data to DATA register
Write command to COMMAND register
    (starts the device and executes the command)
While (STATUS == BUSY)
    ; // wait until device is done with your request
```

Protocol does a lot of **polling**

CPU is responsible for moving data

How can we lower this overhead?

# Interrupt-driven I/O vs Polling

- Use hardware interrupts to avoid busy polling:
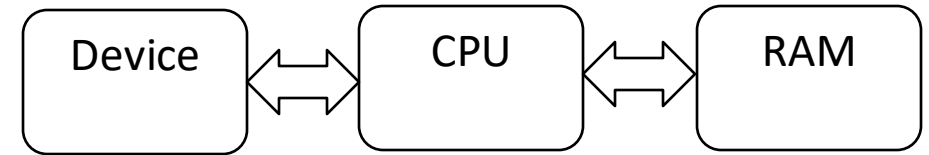  - Allows CPU to process another task. Will get notified when task is done
  - Interrupt handler will read data & error code

- Is it always better to use interrupts?

- Actual devices often support both polling and interrupts: e.g. wait for an interrupt from the network card to read the first packet, then poll its input queue memory space to look for other received packets

# From programmed I/O to direct memory access

- With programmed I/O (our simple protocol):
  - CPU issues read request
  - Device interrupts CPU with data
  - CPU writes data to memory
  - Pros: simple hardware. Cons: Poor CPU is always busy!

```
[Device] <--> [CPU] <--> [RAM]
```

- With direct-memory-access (DMA):
  - CPU sets up DMA request
    - » Gives controller access to memory bus
  - Device puts data on bus & RAM accepts it
  - Device interrupts CPU when done

```
[Device] <--> [RAM]
```

# DMA in more detail



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

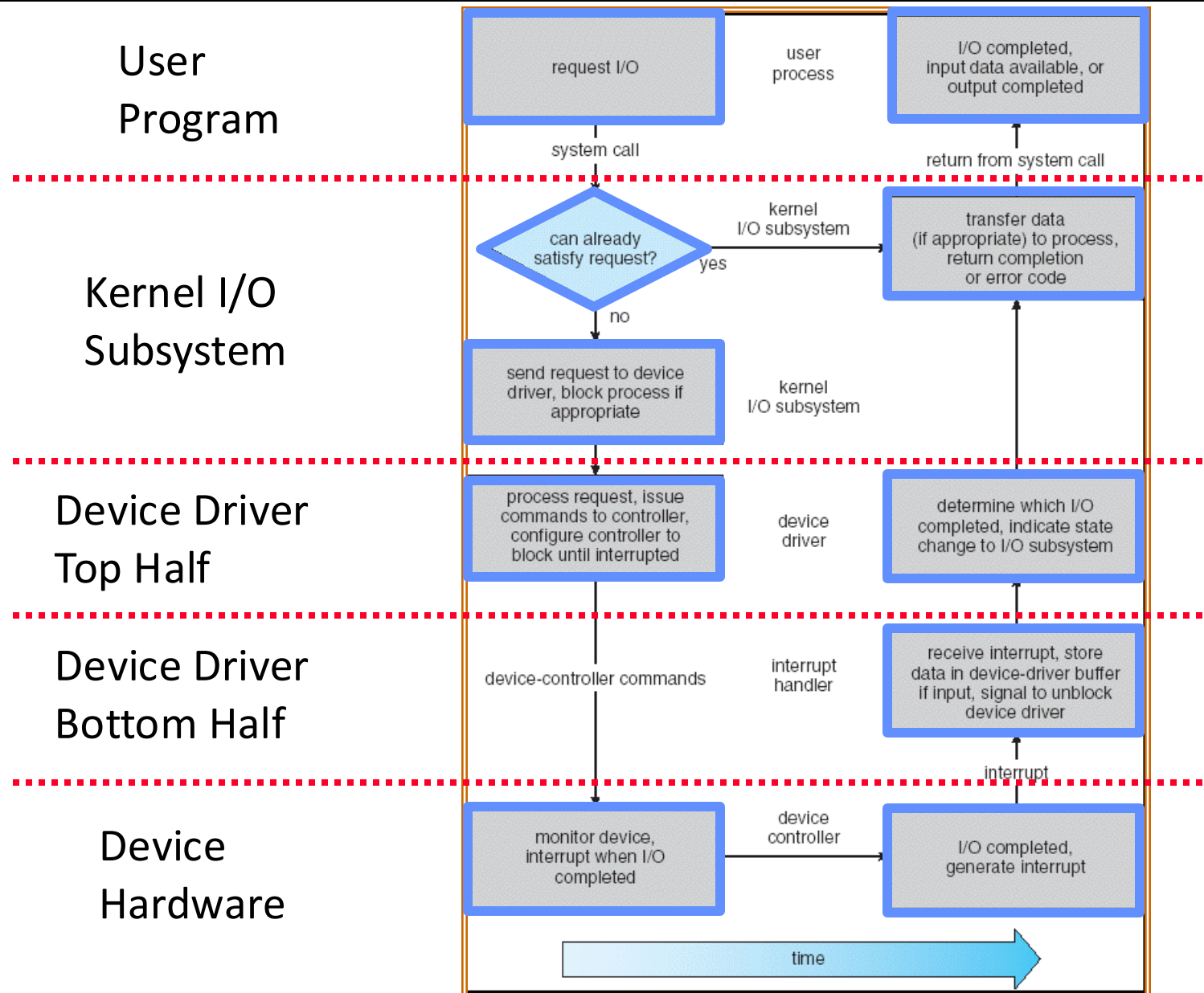# How can the OS handle ~~one~~ all devices

- How do we fit devices with specific interfaces into the OS, which should remain general?
  - Build a "device neutral" OS and hide details of devices from most of OS

- Abstraction to the rescue!
  - **Device Drivers** encapsulate all specifics of device interaction
  - Implement device neutral interfaces

# Device Drivers

- Device Driver: Device-specific code in the kernel that interacts directly with that device hardware
  - Supports a standard, internal interface
  - Special device-specific configuration supported with the `ioctl()` system call

- Device Drivers are 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

- Your body is 90% water, your OS is 70% device-drivers

# Putting it together: Life Cycle of An I/O Request

# Conclusion

- I/O Devices Types:
  - Many different speeds (0.1 bytes/sec to 50+ GBytes/sec)
  - Different Access Patterns:
    - » Block Devices, Character Devices, Network Devices

- Device Controllers: Hardware that controls actual device
  - Processor Accesses through I/O instructions (port-mapped I/O) or load/store to special physical memory addresses (memory-mapped I/O)

- 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

# 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 (or lets user check/poll)
  - When send data, take pointer to user's buffer, return immediately; later kernel takes data and notifies user (or lets user check/poll)

# 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 sequential access

## Flash memory

– Storage that rarely becomes corrupted

– Capacity at intermediate cost (5-20x disk)

– Block level random access

– Good performance for reads; worse for random writes

– Wear patterns issue