# I/O and Devices: H/W S/W Interaction

## Overview

> Basic System Structure
  + device types, roles
> Devices: H/W (and S/W) overview
> Controlling Devices: Problems + Solutions
> Case Study: Disks
  overview, scheduling, trends

## Themes

> Overlap + Parallelism
  Buffering, CPU cost
> Abstractions
> Fairness v. Performance

**Review:** Why should OS manage I/O devices? (What are the roles of OS?)

**Abstractions:** OS as virtual machine
  ⇒ consistent interface to many different devices
    (2 levels: internal to rest of OS, and to user via FS)
  **Tension:** too general ⇒ poor performance,
                lack of feature exploitation

**Resource Manager:** OS as scheduler / multiplexor
  ⇒ Arbiter of system resources
  **Tension:** fairness vs. performance of I/O requests

**Protection:** OS as secure VM
  ⇒ must share between legal users, disallow illegal uses

**summary:** [who] gets [what] [when]
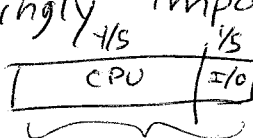
## Today's Focus: I/O

: why care about I/O?
(after all, processors are the coolest
          part of computer systems)

1) without I: programs would produce same result each time
            O: no point in running code

2) I/O performance: increasingly important
   Amdahl's Law
   ⇒ what you do here matters

   | CPU | I/O |

   Processors getting fast much faster!

   CPU gets 8x in 3 years
   ⇒ % of time in I/O?

3) Storage + Networking:
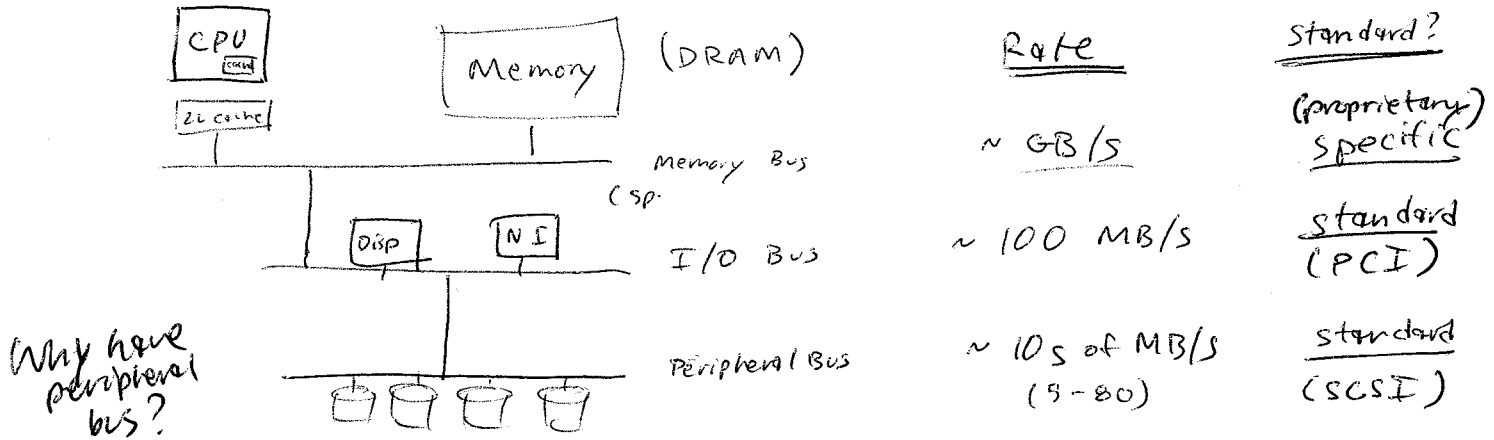   storing your data
   access to remote data  } what Internet is all about!
                            (not processing)

⇒ I/O is where the action is!!

# Basic System Structure

| | Rate | Standard? |
|---|---|---|
| Memory Bus (sp. | ~ GB/s | (proprietary) specific |
| I/O Bus | ~ 100 MB/s | standard (PCI) |
| Peripheral Bus | ~ 10s of MB/s (5-80) | standard (SCSI) |

Why have peripheral bus?

## Observations:

→ Closer to processor, lower the latency

→ Faster the bus, shorter the bus (electrical properties)
   ⟹ fewer devices / bus

▷

## Device Types

| | I/O? | B/C? | S/C? | Rate? |
|---|---|---|---|---|
| Display | O | C | C (w/ human) | 100s MB/s |
| Network Interface (NI) | I/O | C | C (w/ computer) | 1/10s/100s |
| Disks, Tapes | I/O | B | S | 10s MB/s |
| Keyboard, Mouse | I | C | C (w/ human) | bytes/sec |

## How to categorize?

Input / Output / Both

Block / Character
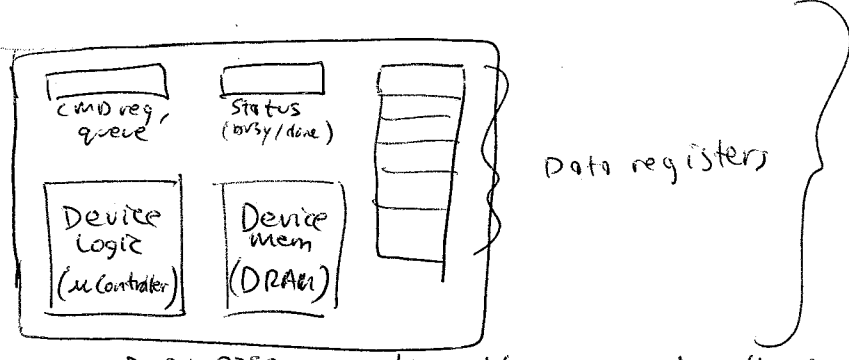
Storage / Communication

Rate of operation

# Devices: The H/W side   "care + feeding"

Device may require lots of attention during operation

Observe low-level progress, provide detailed cmds, correct minor errors
(move disk arm to this location)          (disk read bit)

⇒ put some of this into <u>H/W</u> itself ⇒ <u>device controller</u>

```
┌─────────────────────────────────┐
│ ┌─────────┐  ┌─────────┐  ┌───┐ │
│ │cmd reg/  │ │status    │ │   │ │
│ │queue     │ │(busy/done)│ │   │ │
│ └─────────┘  └─────────┘ │   │ │
│ ┌─────────┐  ┌─────────┐ │   │ │
│ │Device   │  │Device   │ │   │ │
│ │logic    │  │mem      │ │   │ │
│ │(µ Controler)│ │(DRAM)│ │   │ │
│ └─────────┘  └─────────┘ └───┘ │
└─────────────────────────────────┘
```
Data registers

CPU, mem, etc ⇒

<u>"Computer w/ in the computer"</u>!

⇒ expose cmds, H/W executes them (a little program)

## Basic operation

⇒ check status, wait until <u>free</u>
⇒ Put data in data register(s)
⇒ Put cmd in cmd reg
    ( device executes cmd)

⇒ check status
    if busy, keep checking
    if not, all done, do another I/O?

} exactly <u>how</u> we will discuss later

⇒ <u>All done?</u>  <u>no</u>

~~Problem #~~ ~~Lots of devices out there (all w/ different interfaces, registers, etc)~~

## Some <u>problems</u>!

> Lots of different devices out there, how to simplify use?

> How to get processor to talk to devices?

> How to do so efficiently? (overlap)
    (the above)

Problem: Lots of devices out there, how to simplify use?

⸱ not just for apps, but for rest of OS itself!    [FS]

(characteristic of systems in general)    [device drive]
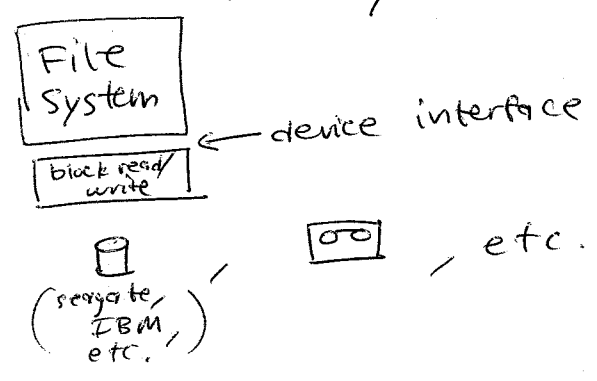
⟹ internal Abstraction Layer

Uniform interface to similar devices (e.g. disk ⟹ block_read / block_write)

S/W takes abstract requests

⟹ specific H/W register manipulations

What's this s/w called?  Device driver

How to build file system then?

| File System |
| block read/write |  ⟵ device interface

🛢 , 📼 , etc.
( seagate, IBM, etc. )

⟹ Simplifies use, but at what cost

⟹ Implementation Notes

> Old days: new driver ⟹ recompile, reboot entire OS
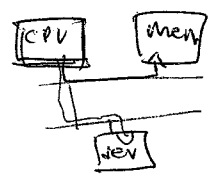
now: loadable on the fly  "plug & play"

> Linux: over half of source code is in drivers

Problem: How to comm. w/ devices?

⟹ Special I/O instructions
   only valid in kernel mode
   was popular in mainframes (no longer)

⟹ store reg$_x$, [M$_y$]

[CPU] [mem]
[dev]

⟹ Memory-mapped I/O
   read + write special memory addresses
   protect by placing in KVM or PM
   simple, general, widely used

⟹ H/W s/w interface

Problem: How to interact efficiently?

Control: How to know when event is complete?

> Polling:

Handshake by setting, clearing flags

{ Plus: simple
{ Minus: CPU cycles are wasted busy-waiting
    (if not attentive, may lose data)

```
        action
start-event()
while (!done)
  check-status
```

> Interrupts

Handle events asynchronously
Device asserts interrupt request line
   when device is ready

CPU: jump to correct Int. Service Routine (ISR)
Interrupt vector: Table of handler addresses
Index by interrupt #

Plus: frees CPU from checking
Minus: could be costly — context switch into handler routine

```
                action()
1) start-event
   notify-me-when-done()

2)   handler
```

⇒ when to use polling, when interrupts?
   ⇒ Depends on speed of device
   If slow, interrupts (allows OS to switch to other job)
     ⇒ OVERLAP
   if quick, polling (cost of ctxt switch avoided)

Problem (cont.) : How to do so efficiently?

· Even w/ interrupts, CPU may have to move data to an fro

Data

=> Programmed I/O (PIO)

CPU moves every byte of data to/fro device

e.g. block read from disk, sitting in controller memory

CPU does device => main mem copy

{ Pro : simple
{ Con : CPU overhead!

=> Direct Memory Access (DMA)

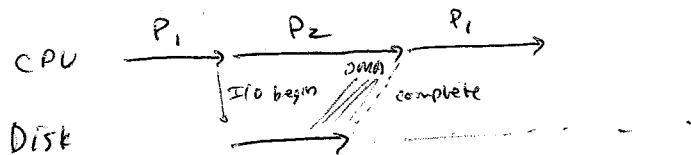Offload work to special-purpose transfer engine

CPU sets up DMA

set up addresses for src, dest, xfer size

DMA controller handles xfer
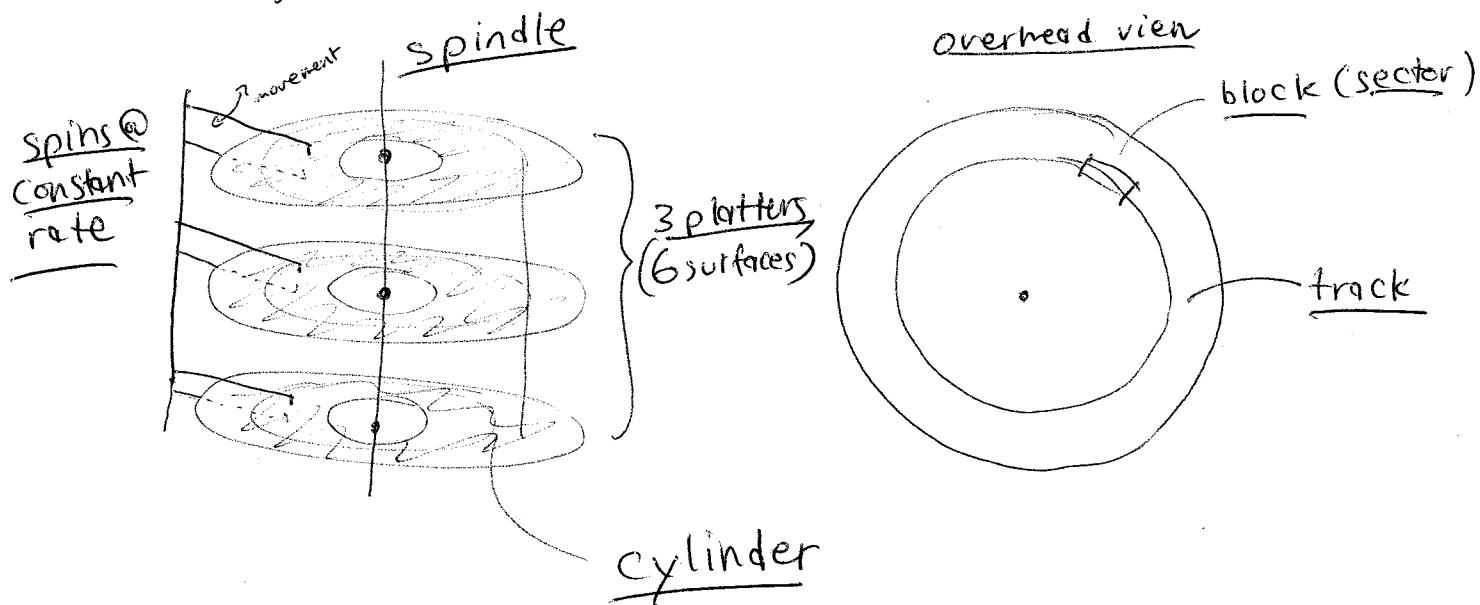
interrupts CPU when finished

=> Now, true overlap is (almost) possible



CPU  ——P₁——→  ——P₂——→  ——P₁——→
                I/o begin    DMA complete
Disk        ——————→

# Case Study : Disks (today and tomorrow)
### (storage)

H/W structure :
Important to understand so we can design better file system!
(general rule : understand H/W to enable you to build S/W)

spins @ constant rate

spindle

↗ movement

3 platters
(6 surfaces)

overhead view

block (sector)

track

cylinder

## To read/write block :

Seek : Position arm/head over correct cylinder
(accelerate, coast, decel., settle)

Rotational Delay : wait for right block/sector
to rotate underneath.

Transfer : when right block is there, transfer to
device memory

$$\Rightarrow T_{I/O} = T_{seek} + T_{rotation} + T_{transfer}$$

only real work

## Typical "modern" disk

IBM 97X

seek : 3 → 18 ms
rot : 0 → 12 ms
transfer : 12 → 20 MB/s

e.g. small        large

why is outer track B/W
higher than inner track?

# Disk Scheduling, etc.

$$T_{I/O} = T_S + T_R + \underbrace{T_{transfer}}_{real\ work}$$

"How to make disks run fast"

FS layout
layout +
request
Scheduling
> Long, sequential transfers
> Avoid seeks, rotations
} treat disk like a tape!

$$T_{seek}, T_{rot} \Rightarrow 0 \qquad OR \qquad T_{transfer} \gg T_{seek}, T_{rot}$$

Problem

Queue of requests: read/write $B_x$
what order to process?

THEME: fairness/perf

Most basic: FCFS (first come, first served)

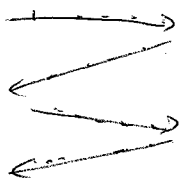service in order
problem? (pathological case?)

SSTF (shortest-seek-time first)
solves previous problem
but, adds a new one: ~~interlock~~ starvation

Scan (Elevator)      Look: variant (don't go to end)

order

starvation: not a problem
new problem:

C-SCAN

(reset)

1) new request pile up
2) 2-scan delay for block (potentially)

oops! all seek centric
new research takes rotational delay into account too!