

Operating Systems

Lecture 13

disk and fs abstraction

Prof. Mengwei Xu

Goals for Today

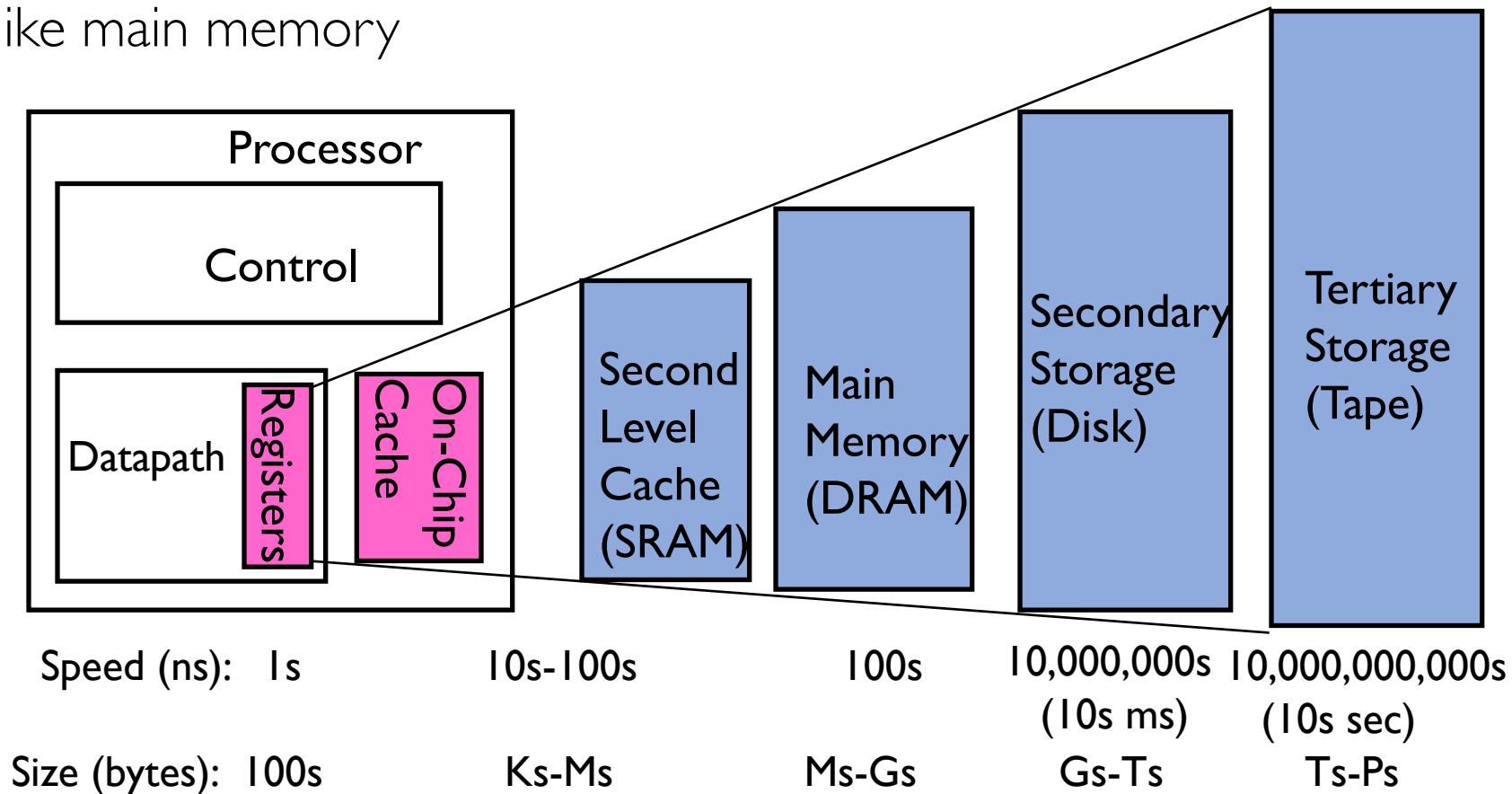
- Storage Devices
- File System Abstraction

Goals for Today

- Storage Devices
- File System Abstraction

Storage Devices

- Why we learn the hardware characteristics? Because they help us build better OSes and applications!
- As secondary storage to computers, storage devices are persistent.
 - Unlike main memory



Storage Devices

1. Magnetic disks (磁盘)

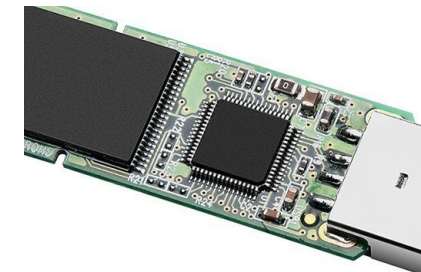
- Storage that rarely becomes corrupted
- Large capacity at low cost
- Block level random access
- Slow performance for random access
- Better performance for sequential access



Servers, workstations,
and laptops

2. Flash memory (闪存)

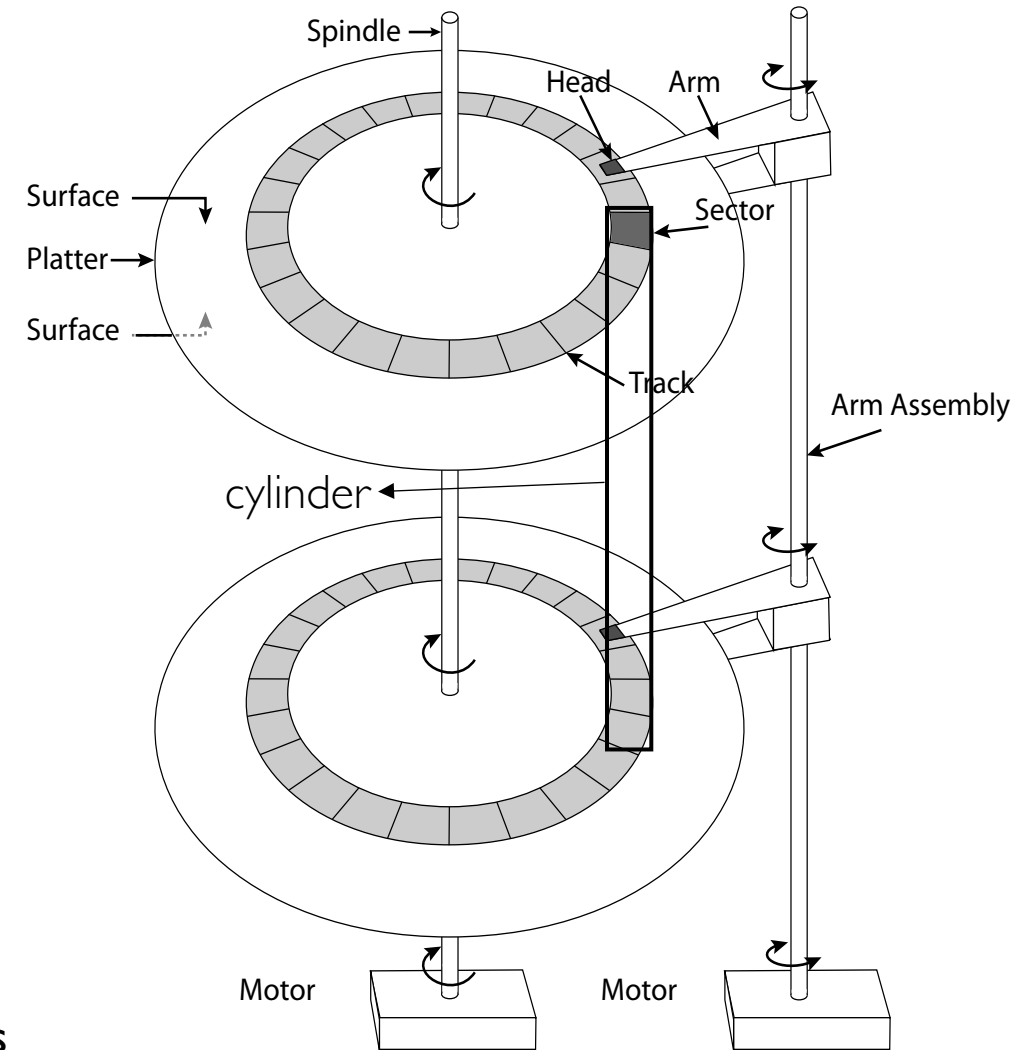
- A type of solid state storage (SSD, 固态硬盘)
- Storage that rarely becomes corrupted
- Capacity at intermediate cost (5-20x disk)
- Block level random access
- Good performance for reads; worse for random writes
- Erasure requirement in large blocks
- Wear patterns issue



Smartphones and tablets

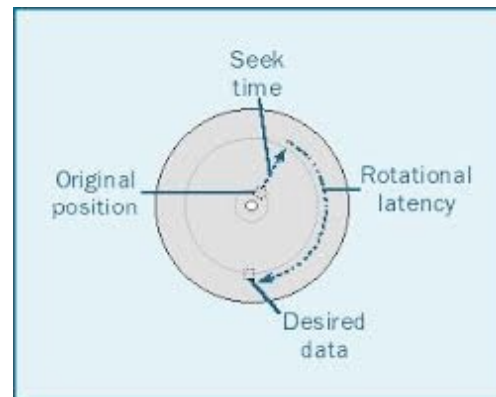
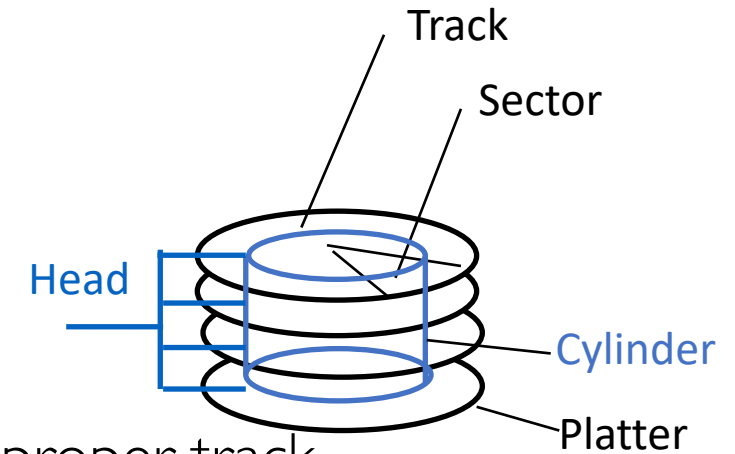
The Magnetic Disk

- Sector (扇区): the unit of transfer
- Track (磁道): ring of sectors
 - $\sim 1\mu\text{m}$ (10^{-6}m) wide
 - Resolution of human eye: $50\mu\text{m}$
 - Wavelength of light is $\sim 0.5\mu\text{m}$
- Cylinder (柱面): stacked tracks
- Head (磁头): attached to movable arms to read data
 - 2 per each platter (磁片) for each surfaces
- Storage capacity =
 $(\text{head \#}) * (\text{cylinder \#}) * (\text{sector \#}) * (\text{sector size})$
Often 512 bytes



The Magnetic Disk

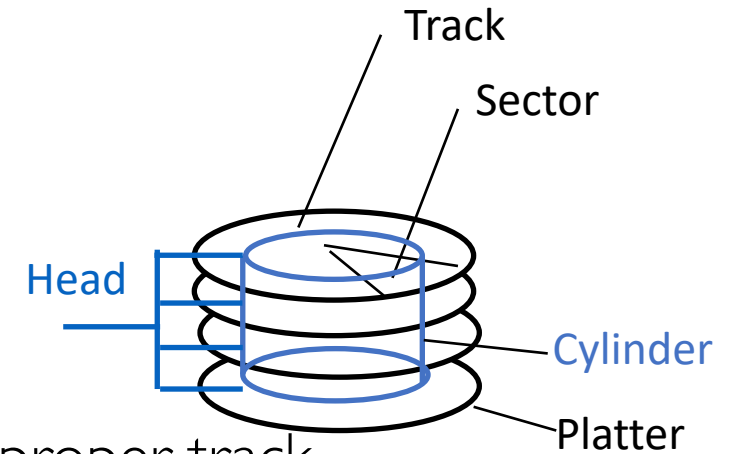
- **Cylinders**: all the tracks under the head at a given point on all surface
- Read/write data is a three-stage process:
 - **Seek time (寻道时间)**: position the head/arm over the proper track
 - **Rotational latency (延迟时间)**: wait for desired sector to rotate under r/w head
 - **Transfer time (传输时间)**: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



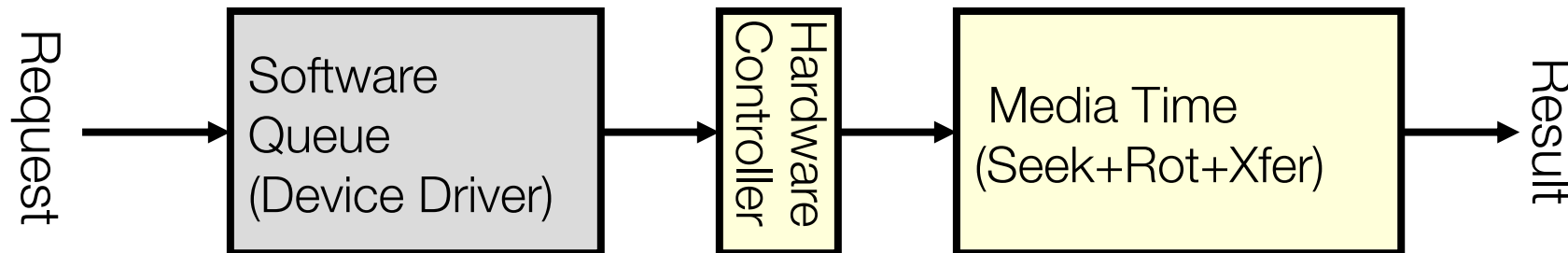
Seek time = 4-8ms
One rotation = 1-2ms
(3600-7200 RPM)

The Magnetic Disk

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



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



Disk Performance Example

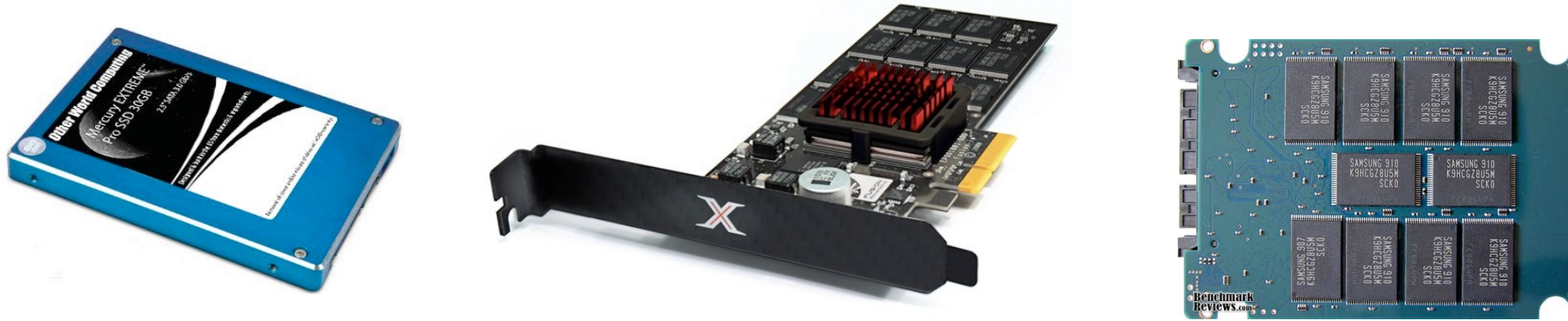
- Assumptions:
 - Ignoring queuing and controller times for now
 - Avg seek time of 5ms,
 - 7200RPM \Rightarrow Time for rotation: $60000 \text{ (ms/minute)} / 7200 \text{ (rev/min)} \approx 8\text{ms}$
 - Transfer rate of 4MByte/s, sector size of 1 Kbyte \Rightarrow
 $1024 \text{ bytes} / 4 \times 10^6 \text{ (bytes/s)} = 256 \times 10^{-6} \text{ sec} \approx .26 \text{ ms}$
- Read sector from random place on disk:
 - Seek (5ms) + Rot. Delay (4ms) + Transfer (0.26ms) = 9.26ms
 - Approx 10ms to fetch/put data: **100 KByte/sec**
- Read sector from random place in same cylinder:
 - Rot. Delay (4ms) + Transfer (0.26ms) = 4.26ms
 - 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

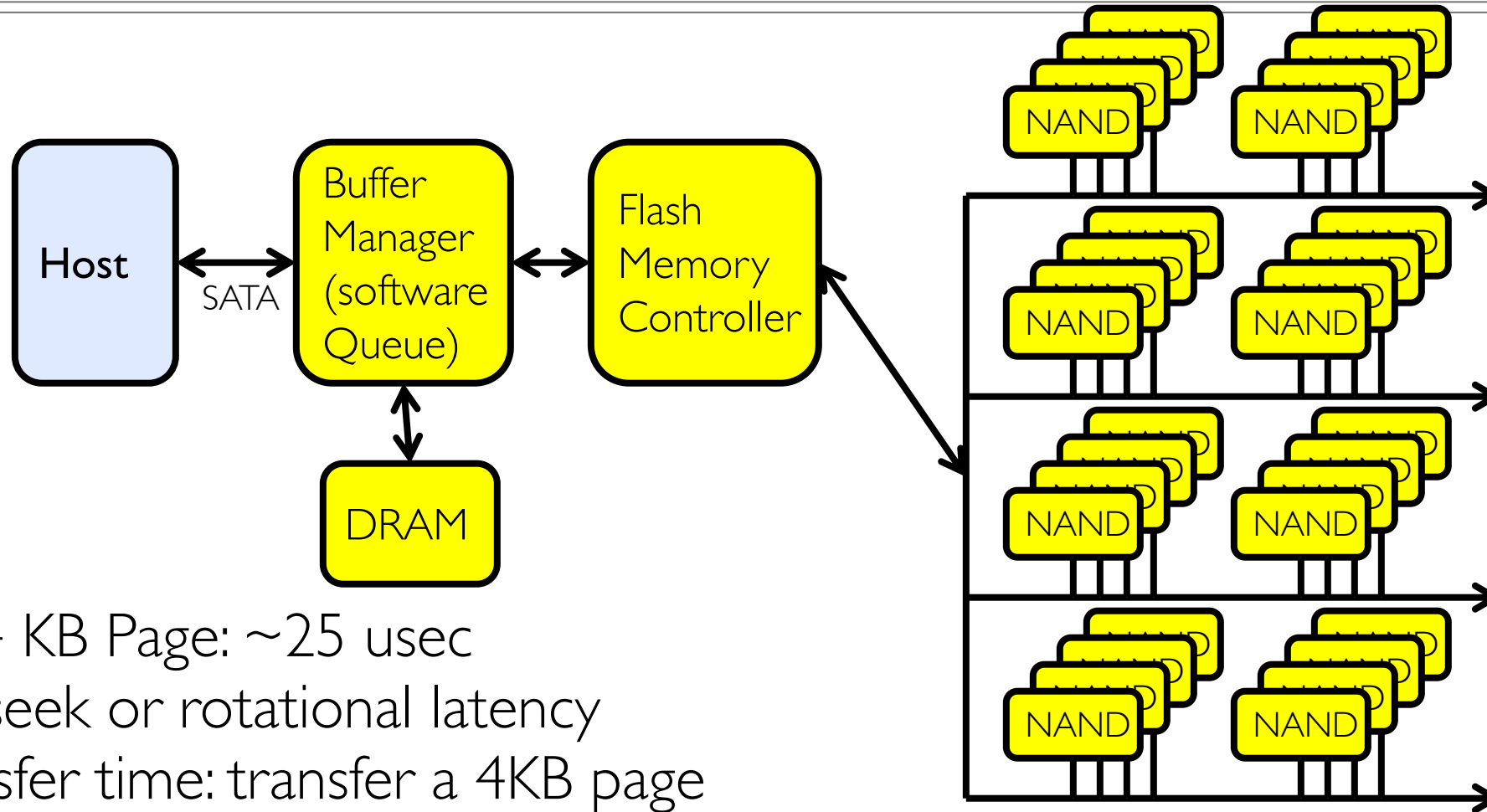
- Sectors contain sophisticated error correcting codes
 - Disk head magnet has a field wider than track
 - Hide corruptions due to neighboring track writes
- Sector sparing
 - Remap bad sectors transparently to spare sectors on the same surface
- Slip sparing
 - Remap all sectors (when there is a bad 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
- ...

Solid State Disks (SSDs)



- 1995 – Replace 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!
- A 5-min video on SSD: <https://www.bilibili.com/video/BV1644yI57mB>

SSD Architecture – Reads



Read 4 KB Page: ~25 usec

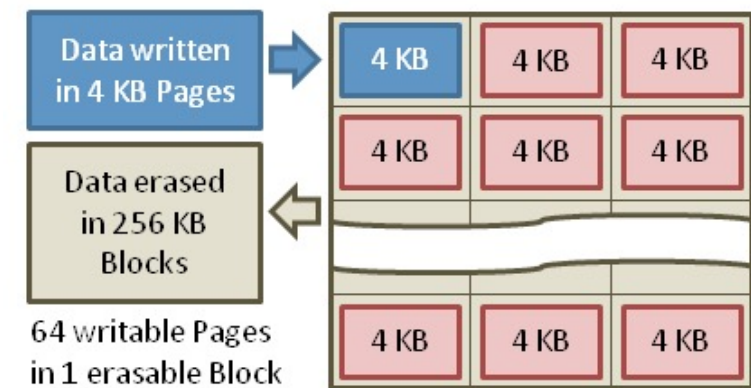
- No seek or rotational latency
- Transfer time: transfer a 4KB page

□ SATA: 300-600MB/s => $\sim 4 \times 10^3 \text{ b} / 400 \times 10^6 \text{ bps} \Rightarrow 10 \text{ us}$

- Latency = Queuing Time + Controller Time + Xfer Time
- Highest Bandwidth: Sequential OR Random reads

SSD Architecture – Writes

- Writing data is complex! ($\sim 200\mu\text{s}$ – 1.7ms)
 - 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

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,
- According to John Kubiawicz (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 (0.1–0.5× disk), expensive (3–20× disk)~~
 - ❑ 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!

No longer
true!

Enterprise

10TB (2016)

- 7 platters, 14 heads
- 7200 RPMs
- 6 Gbps SATA / 12Gbps SAS interface
- 220MB/s transfer rate, cache size: 256MB
- Helium filled: reduce friction and power usage
- Price: \$500 (\$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 !!)

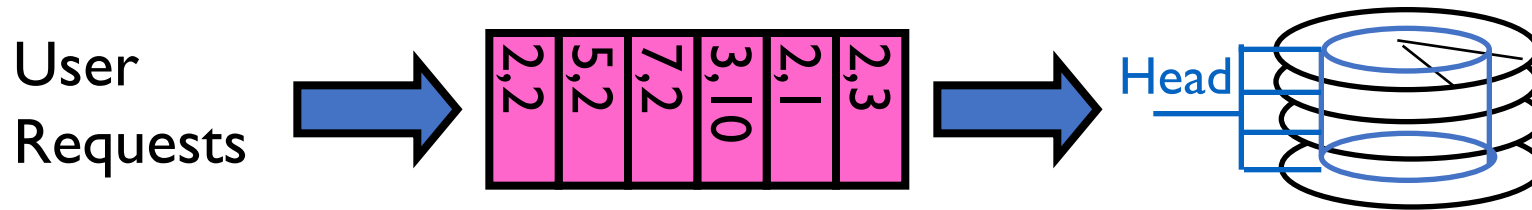
Largest SSDs

- 60TB (2016)
- Dual port: 16Gbs
- Seq reads: 1.5GB/s
- Seq writes: 1GB/s
- Random Read Ops (IOPS): 150K
- Price: ~ \$20K (\$0.33/GB)

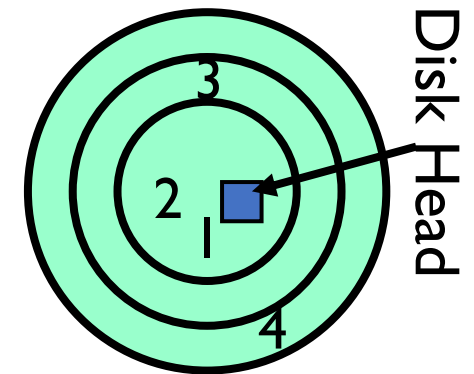


Disk Scheduling

- Disk can do only one request at a time; What order do you choose to do queued requests?
 - The scheduling can be done in OS, firmware, or both.

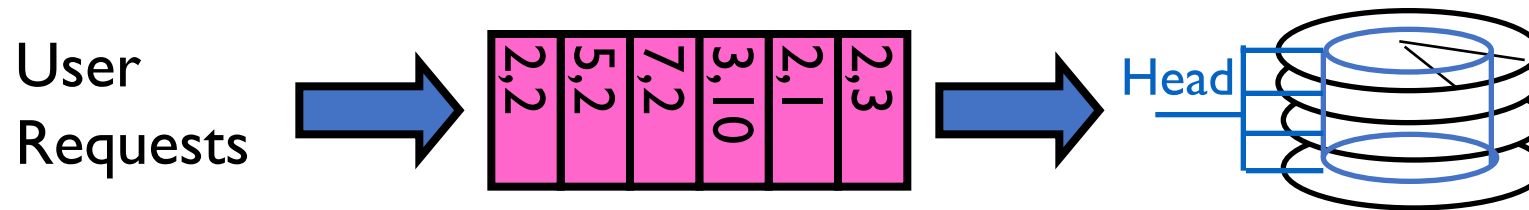


- FIFO Order
 - Fair among requesters, but order of arrival may be to random spots on the disk \Rightarrow Very long seeks

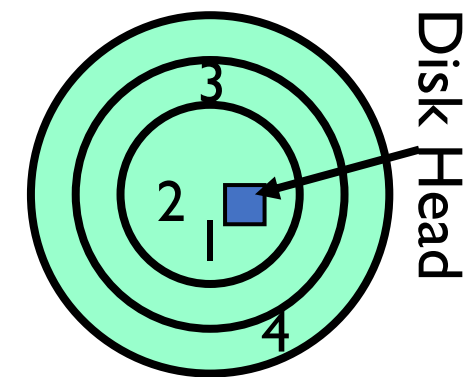


Disk Scheduling

- Disk can do only one request at a time; What order do you choose to do queued requests?
 - The scheduling can be done in OS, firmware, or both.

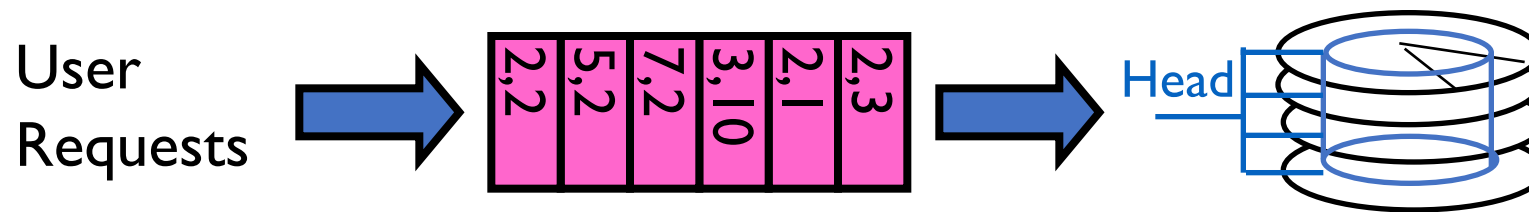


- SSTF: Shortest seek time first
 - Pick the request that's closest on the disk
 - Although called SSTF, today must include rotational delay in calculation, since rotation can be as long as seek
 - Con: SSTF good at reducing seeks, but may lead to starvation

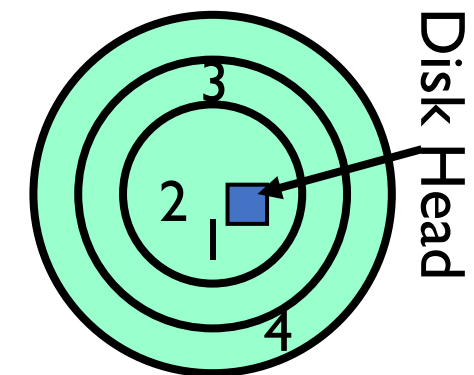
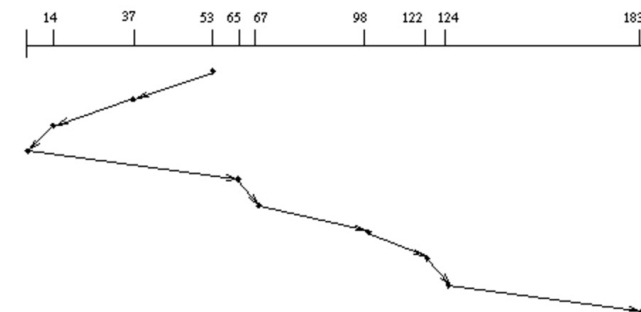


Disk Scheduling

- Disk can do only one request at a time; What order do you choose to do queued requests?
 - The scheduling can be done in OS, firmware, or both.

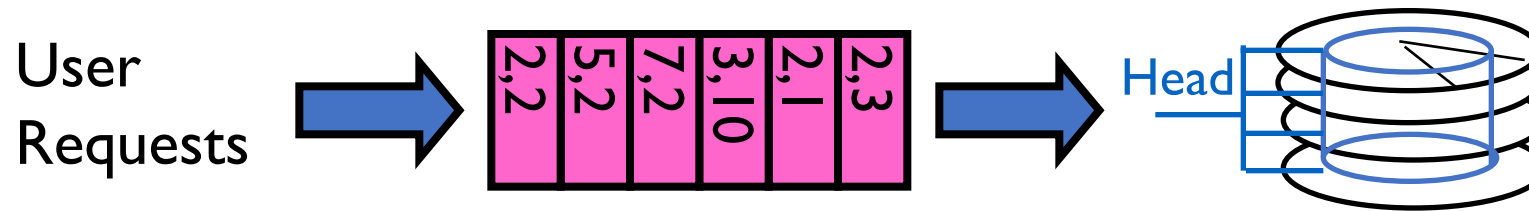


- SCAN: Implements an Elevator Algorithm (电梯算法): take the closest request in a fixed direction of travel (reversed at the end)
 - No starvation, but retains flavor of SSTF

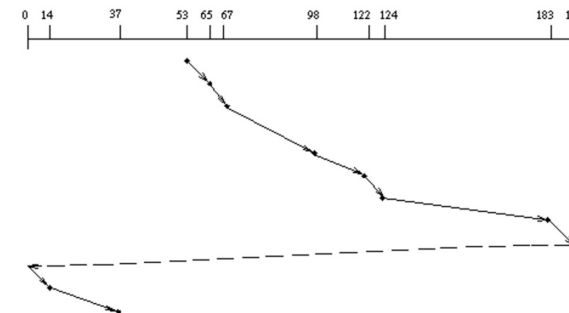


Disk Scheduling

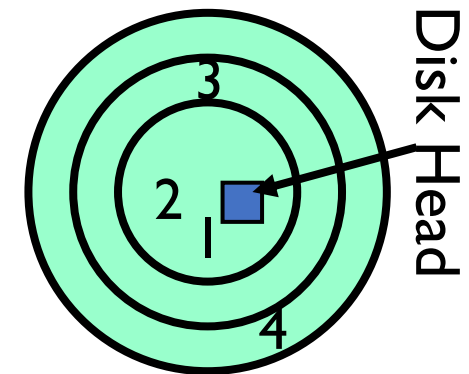
- Disk can do only one request at a time; What order do you choose to do queued requests?
 - The scheduling can be done in OS, firmware, or both.



- C-SCAN: Circular-Scan: only goes in one direction
 - Skips any requests on the way back
 - Fairer than SCAN, not biased towards pages in middle



This diagram shows a C-SCAN disk scheduling path. At the top, a horizontal axis represents the disk cylinders with tick marks at 0, 14, 37, 53, 65, 67, 98, 122, 124, 183, and 199. A solid line with arrows shows the head moving from cylinder 53 towards cylinder 199, visiting requests at 65, 67, 98, 122, and 124. A dashed line with an arrow shows the head jumping back to cylinder 0 and then moving towards cylinder 53, illustrating the circular nature of the scan.



A Simple Read() Lifecycle

- A process issues a syscall `read()`
- OS moves the calling thread to a wait queue (state=WAITING)
- OS uses memory-mapped I/O to tell the disk to read the requested data and set up DMA so the disk can place the data in kernel's memory
- Disk reads the data and DMA's it into main memory
- Disk triggers an interrupt
- OS's interrupt handler copies the data from the kernel's buffer into the process's address space
- OS moves the thread to the ready list
- The thread is scheduled on CPU, and returns from the `read()`

Goals for Today

- Storage Devices
- File System Abstraction

I/O & Storage Layers

Operations, Entities and Interface

Application / Service

High Level I/O

streams

Low Level I/O

handles

Syscall

registers

`file_open, file_read, ... on struct file * & void *`

File System

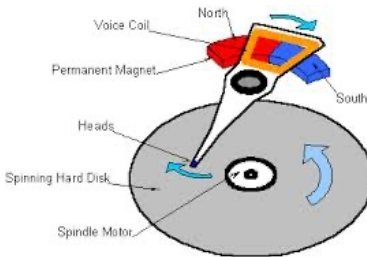
descriptors

we are here ...

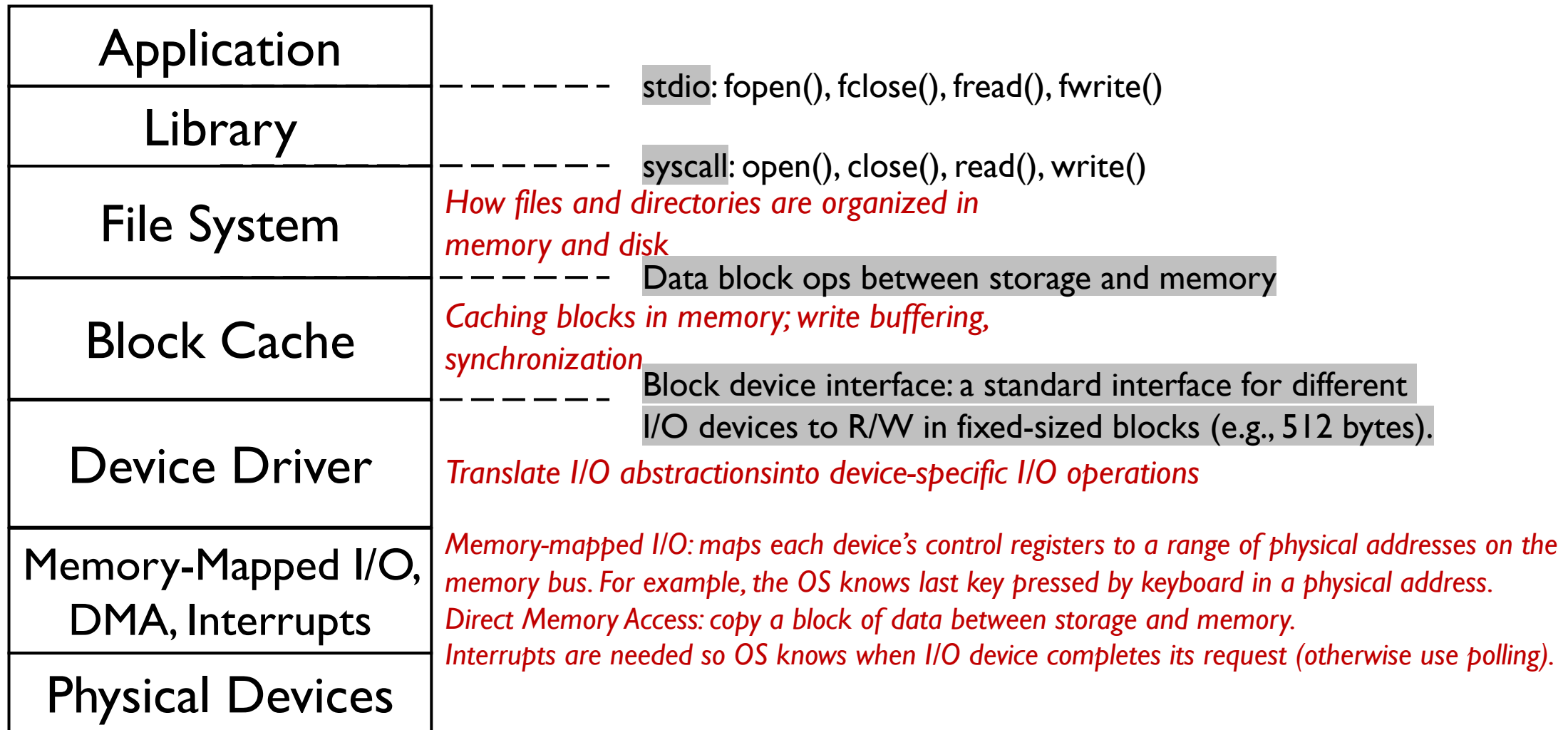
I/O Driver

Commands and Data Transfers

Disks, Flash, Controllers, DMA



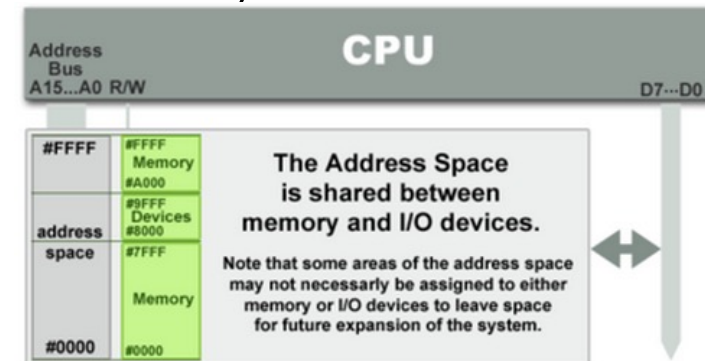
Layered abstractions of I/O and storage



Memory-mapped I/O vs. Port-mapped I/O

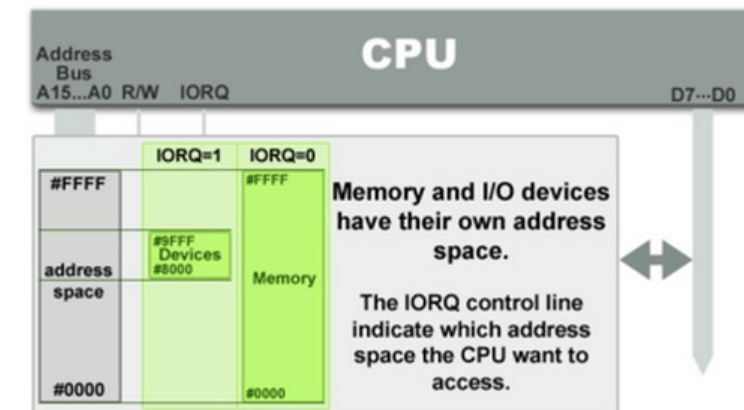
- Two complementary ways for CPU to access I/O devices
 - I/O devices have their own registers (or memory)
- Memory-mapped I/O (MMIO): let memory and devices share the physical address space.

- Most widely adopted
- Shared address bus

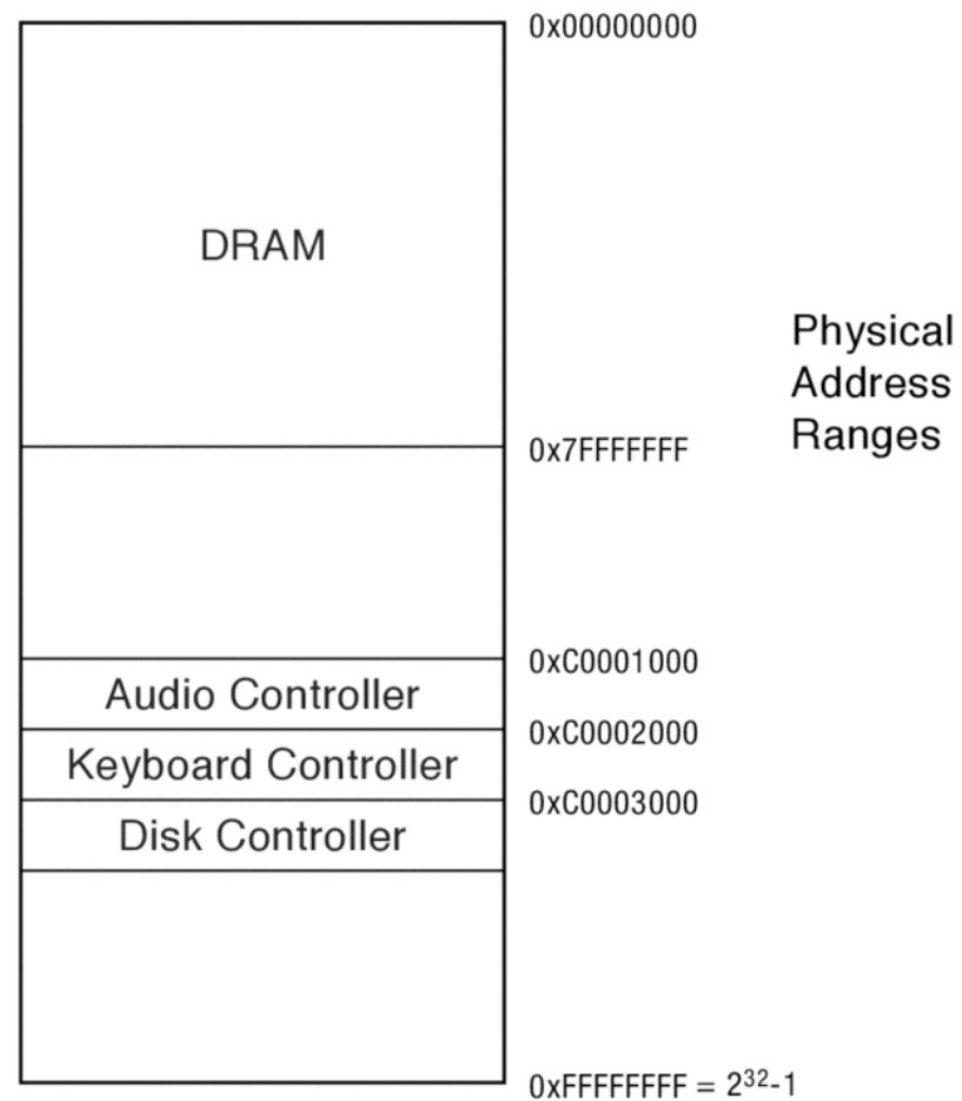
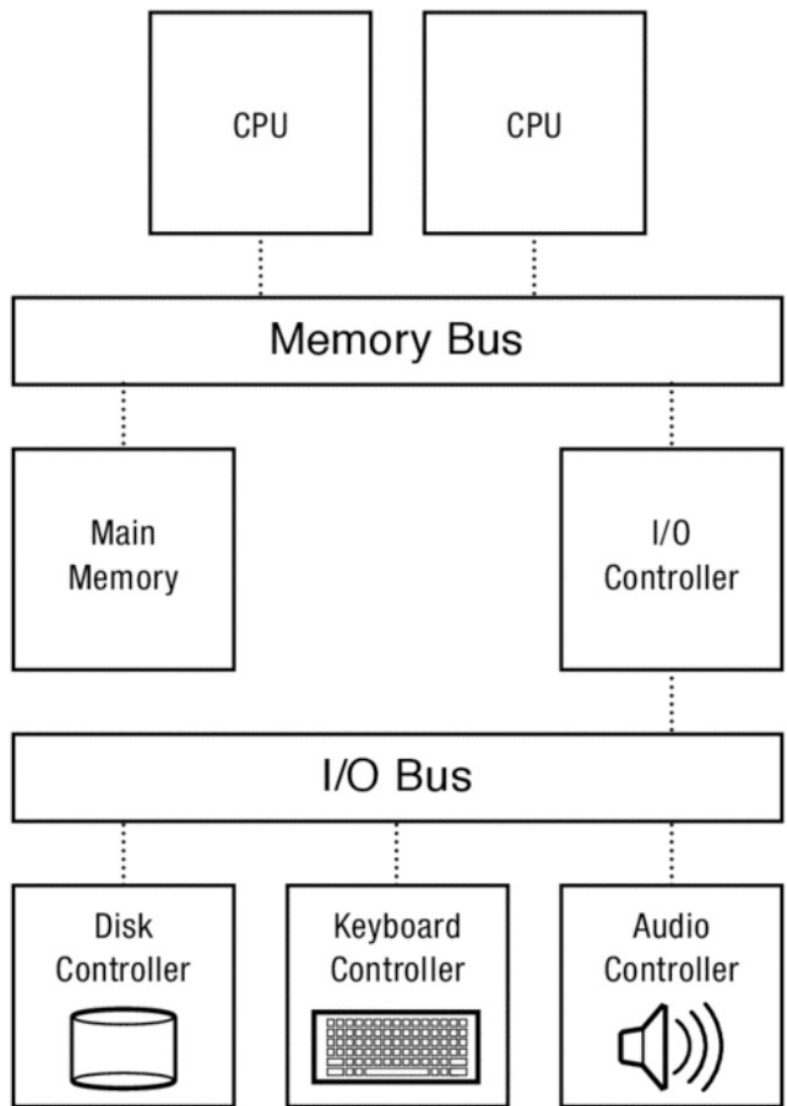


- Port-mapped I/O (PMIO), or isolated I/O: use specialized instructions to R/W I/O devices

- In Intel: outb, outw, etc.



Storage Stack



Recall: C Low level I/O

- File Descriptors – as OS object representing the state of a file
 - User has a “handle” on the descriptor

```
#include <fcntl.h>
#include <unistd.h>
#include <sys/types.h>

int open (const char *filename, int flags [, mode_t mode])
int create (const char *filename, mode_t mode)
int close (int filedes)
```

Bit vector of:

- Access modes (Rd,Wr, ...)
- Open Flags (Create, ...)
- Operating modes (Appends, ...)

Bit vector of Permission Bits:

- User|Group|Other X R|W|X

http://www.gnu.org/software/libc/manual/html_node/Opening-and-Closing-Files.html

Recall: C Low level I/O

- File Descriptors – as OS object representing the state of a file
 - User has a “handle” on the descriptor

`ssize_t read (int fildes, void *buffer, size_t maxsize)`

- returns bytes read, 0 => EOF, -1 => error

`ssize_t write (int fildes, const void *buffer, size_t size)`

- returns bytes written

`off_t lseek (int fildes, off_t offset, int whence)`

- set the file offset

- * if whence == SEEK_SET: set file offset to “offset”

- * if whence == SEEK_CUR: set file offset to crt location + “offset”

- * if whence == SEEK_END: set file offset to file size + “offset”

`int fsync (int fildes)`

- wait for i/o of fildes to finish and commit to disk

`void sync (void)` - wait for ALL to finish and commit to disk

- When write returns, data is on its way to disk and can be read, but it may not actually be permanent!

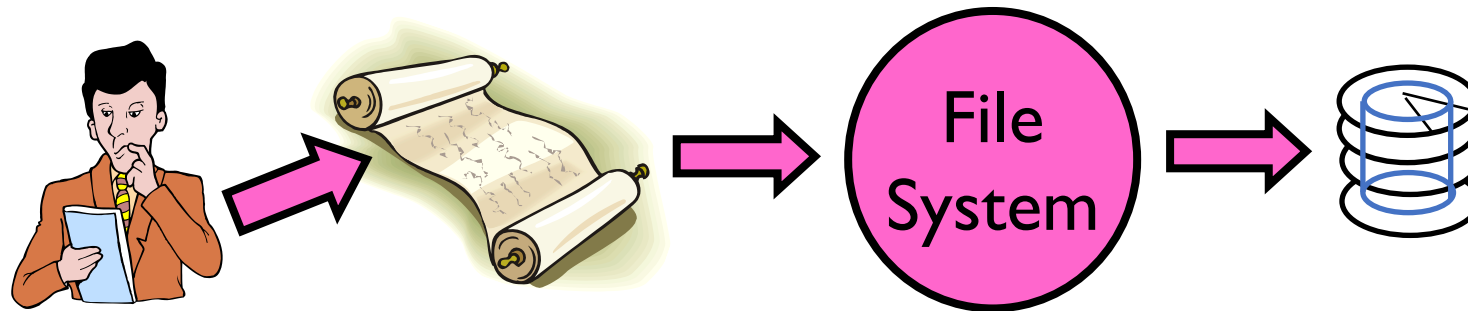
Building a File System

- **File System:** Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.
- File System Components
 - **Naming:** Interface to find files by name, not by blocks
 - **Disk Management:** collecting disk blocks into files
 - **Protection:** Layers to keep data secure
 - **Reliability/Durability:** Keeping of files durable despite crashes, media failures, attacks, etc.

User vs. System View of a File

- User's view:
 - Durable Data Structures
- System's view (system call interface):
 - Collection of Bytes (UNIX)
 - Doesn't matter to system what kind of data structures you want to store on disk!
- System's view (inside OS):
 - Collection of blocks (a block is a logical transfer unit, while a sector is the physical transfer unit)
 - Block size \geq sector size; in UNIX, block size is 4KB

Translating from User to System View



- What happens if user says: give me bytes 2—12?
 - Fetch block corresponding to those bytes
 - Return just the correct portion of the block
- What about: write bytes 2—12?
 - Fetch block
 - Modify portion
 - Write out Block
- Everything inside File System is in whole size blocks
 - For example, **getc()**, **putc()** \Rightarrow buffers something like 4096 bytes, even if interface is one byte at a time
- From now on, file is a collection of blocks

Disk Management Policies (1/2)

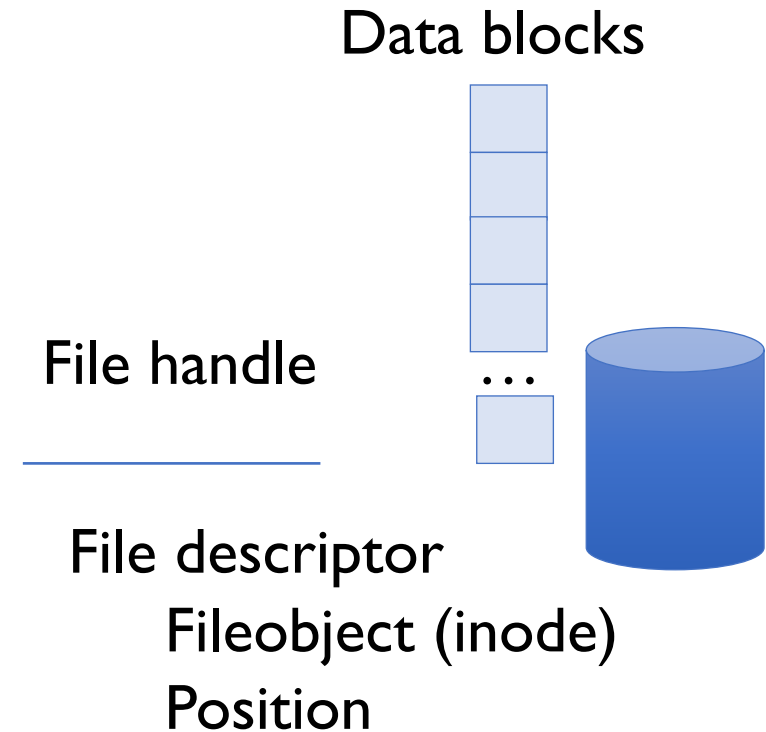
- Basic entities on a disk:
 - **File**: user-visible group of blocks arranged sequentially in logical space
 - **Directory**: user-visible index mapping names to files
- Access disk as linear array of sectors. Two Options:
 - Identify sectors as vectors [cylinder, surface, sector], sort in cylinder-major order
 - ❑ Used in BIOS, but not in OSes anymore
 - **Logical Block Addressing (LBA, 逻辑块寻址)**: Every sector has integer address from zero up to max number of sectors
 - Controller translates from address \Rightarrow physical position
 - ❑ First case: OS/BIOS must deal with bad sectors
 - ❑ Second case: hardware shields OS from structure of disk

Disk Management Policies (2/2)

- Need way to track free disk blocks
 - Link free blocks together \Rightarrow too slow today
 - Use bitmap to represent free space on disk
- Need way to structure files: File Header
 - Track which blocks belong at which offsets within the logical file structure
 - Optimize placement of files' disk blocks to match access and usage patterns

File

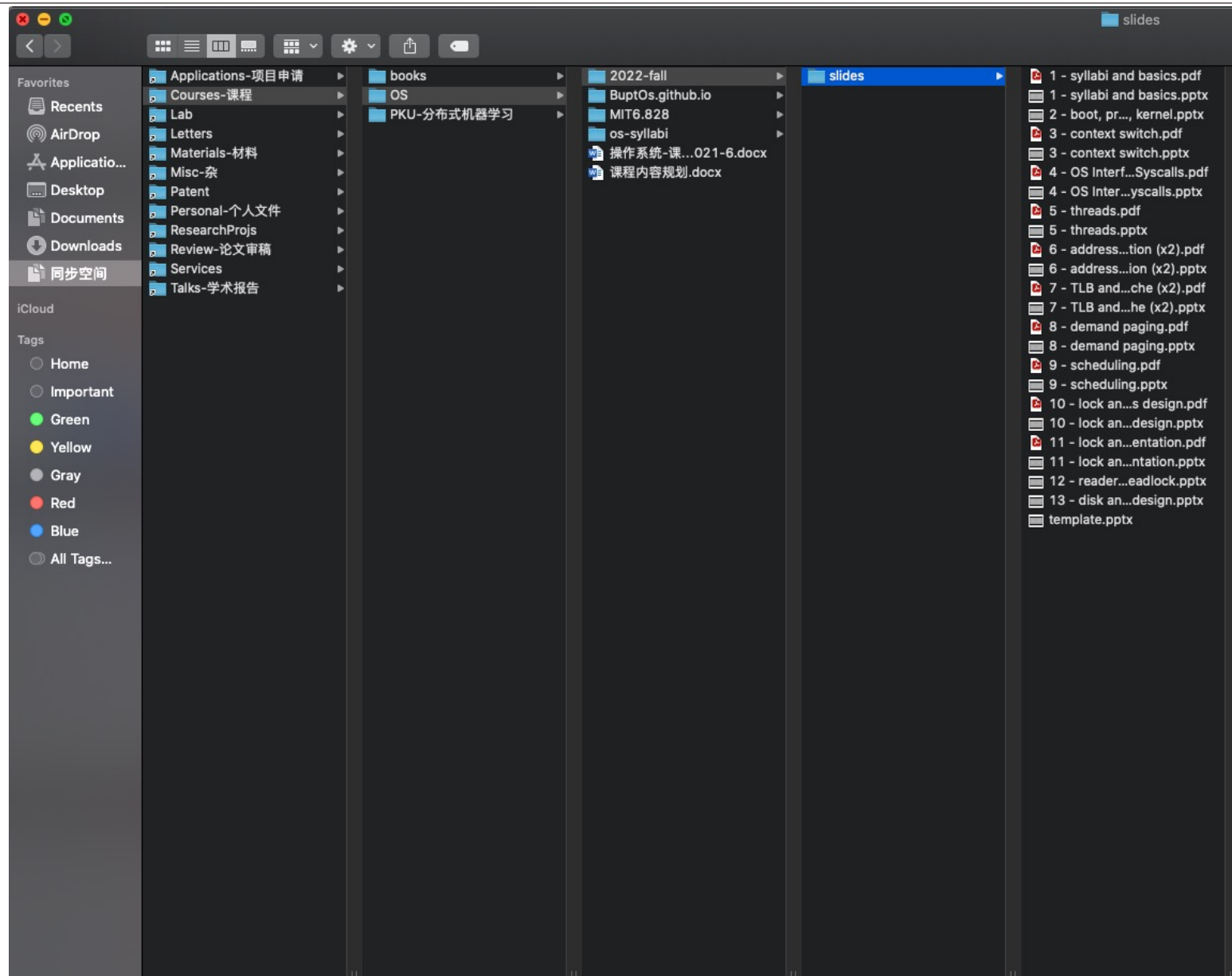
- Named permanent storage
- Contains
 - Data
 - ☐ Blocks on disk somewhere
 - Metadata (Attributes)
 - ☐ Owner, size, last opened, ...
 - ☐ Access rights
 - R, W, X
 - Owner, Group, Other (in Unix systems)
 - Access control list in Windows system



Directory

- Basically a hierarchical structure
- Each directory entry is a collection of
 - Files
 - Directories
 - A link to another entries
- Each has a name and attributes
 - Files have data
- Links (hard links) make it a DAG, not just a tree
 - Softlinks (aliases) are another name for an entry

Directory



Directory

- Conventions of directory
 - Root directory (根目录): “/”
 - Home directory (主目录): “~/cur_dir/file.txt”
 - Absolute path (绝对路径): “/home/mwx/cur_dir/file.txt”
 - Relative path (相对路径): “file.txt”
- Volume (卷): a collection of physical storage resources that form a logical storage device. Could be a part of or many physical devices.
- Mount (挂载): an operation that creates a mapping from some path in the existing file system to the root directory of the mounted volume's file system

`mount -t type device dir`

Directory

```
mw@Dragon21:~$ findmnt -t ext4
```

TARGET	SOURCE	FSTYPE	OPTIONS
/	/dev/sda6	ext4	rw,relatime,errors=remount-ro
└─/data2	/dev/sdc	ext4	rw,relatime
└─/data	/dev/sdb1	ext4	rw,relatime
└─/var/lib/snapd	/dev/sdc[/zi/snap/snapd]	ext4	rw,relatime
└─/boot	/dev/sda1	ext4	rw,relatime

Designing a File System ...

- What factors are critical to the design choices?
- Durable data store => it's all on disk
- (Hard) Disks Performance !!!
 - Maximize sequential access, minimize seeks
- Open before Read/Write
 - Can perform protection checks and look up where the actual file resource are, in advance
- Size is determined as they are used !!!
 - Can write (or read zeros) to expand the file
 - Start small and grow, need to make room
- Organized into directories
 - What data structure (on disk) for that?
- Need to allocate / free blocks
 - Such that access remains efficient