CS162 Operating Systems and Systems Programming Lecture 18

Storage Devices & File Systems

Professor Natacha Crooks & Matei Zaharia https://cs162.org/

Slides based on prior slide decks from David Culler, Ion Stoica, John Kubiatowicz, Alison Norman and Lorenzo Alvisi

Recall : Simplified IO architecture



Follows a hierarchical structure because of cost:

the faster the bus, the more expensive

Recall: How Does Processor Talk to Devices?



status

Registers

(port 0x20)

and/or

Queues

Memory Mapped

Region: 0x8f008020

- 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 in Intel assembly: out 0x21,AL
 - Memory-mapped I/O: load/store instructions
 - » Registers/memory appear in physical address space

- Device Driver: Device-specific code in the kernel for each supported device
 - Implements a standard, internal interface (e.g. block or character device)
 - Device-specific configuration with the ioctl() system call
- Drivers typically have 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()
 - » 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
- Your body is 90% water, your OS is 70% device drivers

Latency - time to complete a task Measured in units of time (s, ms, us, ..., hours, years)

Throughput or *Bandwidth* – rate at which tasks are performed Measured in units of things per unit time (ops/s, GFLOP/s)

Start up or Overhead – time to initiate an operation

Most I/O operations are roughly linear in *b* bytes - Latency(b) = Overhead + b/TransferCapacity

Magnetic disks (HDDs)

- 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 (SSDs)

- Storage that rarely becomes corrupted
- Capacity at intermediate cost (3-20x disk)
 - Block level random access
- Good performance for reads; worse for random writes
 - Wear patterns issue

Hard Disk Drives (HDDs)





Read/Write Head Side View

IBM Personal Computer 1986

30MB Hard Disk for 500 dollars

The Amazing Magnetic Disk

Store data magnetically on thin metallic film bonded to rotating disk of glass, ceramic, or aluminum



The Amazing Magnetic Disk



Track: concentric circle on surface

Sectors: slice of a track Smallest addressable unit Are units of transfers

Cylinder all the tracks under the head at a given point on all surfaces

The Amazing Magnetic Disk



Track lengths vary across disk: outside tracks have more sectors per track, higher bandwidth

Disk is organized into regions of tracks with the same number of sector/tracks

Usually, only outer half of radius is used

Reading/Writing Data

Seek time: position the head/arm over the proper track

Rotational latency: wait for desired sector to rotate under r/w head

Transfer time: transfer a block of bits (sector) under r/w head



Reading/Writing Data

Request Time =

Queueing Time + Controller Time + Seek + Rotational + Transfer



Typical Numbers for Magnetic Disk

Parameter	Info/Range
Space/Density	Space: 18TB (Seagate), 9 platters, in 3½ inch form factor! Areal Density: ≥ 1 Terabit/square inch! (PMR, Helium,)
Average Seek Time	Typically 4-6 milliseconds
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 4-8 milliseconds
Controller Time	Depends on controller hardware
Transfer Time	 Typically 50 to 250 MB/s. Depends on: Transfer size (usually a sector): 512B – 1KB per sector Rotation speed: 3600 RPM to 15000 RPM Recording density: bits per inch on a track Diameter: ranges from 1 in to 5.25 in
Cost	Used to drop by a factor of two every 1.5 years (or faster), now slowing down

Disk Performance Example

Key to using disk effectively (especially for file systems) is to minimize seek and rotational delays

Do access patterns influence how fast can read/write to disk?

Avg seek time of 5ms,

7200RPM \Rightarrow Time for rotation: 60000 (ms/min)/7200(rev/min) ~= 8ms

Transfer rate of 50MByte/s, block size of 4Kbyte \Rightarrow 4096 bytes/50×10⁶ (bytes/s) = 81.92 × 10⁻⁶ sec \cong 0.082 ms for 1 sector

Read block from random place on disk (random reads):

- Seek (5ms) + Rot. Delay (4ms) + Transfer (0.082ms) = 9.082ms
- Approx 9ms to fetch/put data: 4096 bytes/9.082×10⁻³ s \cong 451KB/s

Read block from random place in same cylinder:

- Rot. Delay (4ms) + Transfer (0.082ms) = 4.082ms
- Approx 4ms to fetch/put data: 4096 bytes/4.082×10⁻³ s \cong 1.03MB/s

Read next block on same track (sequential reads):

- Transfer (0.082ms): 4096 bytes/ 0.082×10^{-3} s \cong 50MB/sec

When there are big sequential reads, or

When there is so much work to do that they can be piggy backed (reordering queues—one moment)

OK to be inefficient when things are mostly idle Bursts are both a threat and an opportunity <your idea for optimization goes here> Waste space for speed? Disk can do only one request at a time; What order do you choose to do queued requests?





Disk Scheduling (1/3)



FIFO Order Fair among requesters, but order of arrival may be to random spots on the disk

> SSTF: Shortest seek time first Pick the request that's closest on the disk Con: SSTF good at reducing seeks, but may lead to starvation

Crooks & Zaharia CS162 © UCB Spring 2025

Disk Scheduling (2/3)



SCAN: Implements an Elevator Algorithm: take the closest request in the direction of travel

– No starvation, but retains flavor of SSTF



Disk Scheduling (3/3)

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

- Fairer than SCAN, not biased towards pages in middle



Lots of Intelligence in the Device Controller

Shingled Magnetic Recording (SMR): pack more data with error correcting codes Disk write head has wider field than read head, so overlap tracks Hide corruptions due to neighboring track writes (but must fix in background)

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

Example of Current HDD

- Seagate Exos X24 (2023)
 - 24 TB hard disk
 - » 10 platters, 20 heads
 - » 1.26 TB/in²
 - » Helium filled: reduce friction and power
 - 4.16 ms average seek time
 - 4096 byte physical sectors
 - 7200 RPMs
 - Dual 6 Gbps SATA /12Gbps SAS interface
 - » 285MB/s MAX transfer rate
 - » Cache size: 512MB
 - Price: \$ 479 (~ \$0.02/GB) 385 x
- Origin BM Personal Computer/AT (1986)
 - 30 MB hard disk 850k x
 - 30-40 ms average see ime
 - 0.7-1 MB/s (est.)
 - Price: \$500 (\$17K/GB)

10 x

 $\overline{\mathbf{s}}$

Solid State Drives (SSDs)

1995 – Replace rotating magnetic media with non-volatile memory (battery backed DRAM)

- 2009 Use 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"

The Flash Cell

Encode bit by trapping electrons into a cell

Single-level cell (SLC) Single bit is stored within a transistor Faster, more lasting (50k to 100k writes before wear out)

Multi-level cell (MLC)

Two/three bits are encoded into different levels of charge Wear out much faster (1k to 10k writes)

Of Banks, Blocks, and Cells

Distinction between blocks and pages important in operations!

How do you read?

- Chip supports reading pages
- 10s of microseconds, independently of the previously read page

What about writing? More complicated!

– Must first erase the block

» Erase quite expensive (milliseconds)

- Once block has been erased, can then program a page

» Change 1s to 0s within a page.

» 100s of microseconds.

- Blocks can only be erased a limited number of times!

Low-level Flash Operations

Erase() Program(0) Program(0) Program(1) Erase()

 \rightarrow

 \rightarrow

- iiii Initial: pages in block are invalid (i)
- EEEE State of pages in block set to erased (E)
- \rightarrow VEEE *Program page 0; state set to valid* (V)
 - **error** *Cannot re-program page after programming*
 - VVEE Program page 1
 - EEEE Contents erased; all pages programmable

Low-level Flash Operations

Assume block of 4 pages. All valid. Want to write Page 0

Page 0	Page 1	Page 2	Page 3		
00011000	11001110	00000001	00111111		
VALID	VALID	VALID	VALID		

Step 1: erase full block

Page 0	Page 1	Page 2	Page 3
11111111	11111111	11111111	11111111
ERASED	ERASED	ERASED	ERASED

Step 2: program page 0

 Page 0	Page 1	Page 2	Page 3		
00000011 1111111		11111111	11111111		
VALID	ERASED	ERASED	ERASED		

SSD Architecture

Recall that SSDs uses low-level Flash operations to provide same interface as HDD – read and write chunk (4KB) at a time

Reads are easy, but for writes, can only overwrite data one block (256KB) at a time!

Why not just erase and rewrite new version of entire 256KB block?

- Erasure is very slow (milliseconds)
- Each block has a finite lifetime, can only be erased and rewritten about 10K times
- Heavily used blocks likely to wear out quickly

SSD Architecture (Simplified)

Add a layer of indirection: the flash translation layer

Translates request for logical blocks (device interface) to low-level Flash blocks and pages

Reduce write amplification

Ratio of the total write traffic in bytes issues by the flash chip by the FTL divided by the total write traffic issued by the OS to the device

Avoid wear out

A single block should not be erased too often

FTL – Two Systems Principles

FTL uses indirection and copy-on-write

Maintains mapping tables in DRAM

 Map virtual block numbers (which OS uses) to physical page numbers (which flash mem. controller uses)

Can now freely relocate data w/o OS knowing

Copy on Write/Log-structured FTL

- Don't overwrite a page when OS updates its data
 - Instead, write new version in a free page
 - Update FTL mapping to point to new location

FTL Example

Some Recent SSDs

- Silicon Power 4TB SATA Internal SSD (2023)
 - Seq reads 540 MB/s
 - Seq writes 500 MB/s
 - Price (Amazon): \$184 (\$0.46/GB)

- Micron 9300 Pro 15.36TB NVMe U.2 Enterprise SSD (2019)
 - Seq reads/writes: 3500 MB/s
 - Random Read Ops (IOPS): 100K+
 - Price: \$1750 (\$1.13/GB)

HDD vs. SSD Comparison

SSD prices falling faster than 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 high 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

- Pros (vs. hard disk drives):
 - Low latency, high throughput (eliminate seek/rotational delay)
 - No moving parts:

» Very light weight, low power, silent, very shock in No

- Read at high speeds (limited by controller and I/O bus longer
- Cons

- Small storage (0.1-0.5x disk), expensive (3-20x 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

true!

- 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!

Recall: I/O and Storage Layers

From Storage to File Systems

Layer of OS that transforms block interface of disks (or other block devices) into Files, Directories, etc.

Building a File System

OS as an illusionist: Take limited hardware interface (array of blocks) and provide a more convenient/useful interface with:

Naming: Find file by name, not block numbers

Organize file names with directories

Organization: Map files to blocks

Protection: Enforce access restrictions

Reliability: Keep files intact despite crashes, failures, etc.

Crooks & Zaharia CS162 © UCB Spring 2025

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

Translation 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 writing bytes 2 - 12?

– Fetch block, modify relevant portion, write out block

Everything inside file system is in terms of whole-size blocks

Basic entities on a disk: File: user-visible group of blocks arranged sequentially in logical space Directory: user-visible index mapping names to files

The disk is accessed as linear array of sectors

Old: Physical Position [cylinder, surface, sector]
 New: Logical Block Addressing (LBA)
 Every sector has integer address
Controller translates from address ⇒ physical position
 Shields OS from structure of disk

What Does the File System Need?

Track free disk blocks

- Need to know where to put newly written data

Track which blocks contain data for which files – Need to know where to read a file from

> Track files in a directory – Find list of file's blocks given its name

Where do we maintain all of this?

- Somewhere on disk

Recall: FD & File Descriptors

Description Table

Critical Factors in File System Design

(Hard) Disks Performance !!!

Open before Read/Write

Size is determined as they are used !!!

Organized into directories

Need to carefully allocate / free blocks

Files & Directories

		* ~		Q. Search
avorites	Name	Date Modified	Size	Kind
11 Drophov	V static	Feb 10, 2016, 12:45 PM		Folder
- Diobpox	CSS	Jan 14, 2016, 11:51 AM		Folder
iCloud Drive	exams	Mar 10, 2016, 9:03 PM		Folder
AirDrop	Image:	Jan 14, 2016, 11:51 AM		Folder
E Busians	🔻 🛄 hw	Mar 1, 2016, 7:29 PM		Folder
Ш Desktop	hw0.pdf	Jan 20, 2016, 3:19 PM	175 KB	PDF Document
î adj	hw1.pdf	Feb 11, 2016, 9:42 AM	128 KB	PDF Document
Anniications	hw2.pdf	Feb 16, 2016, 9:00 PM	180 KB	PDF Document
-D -	hw3.pdf	Mar 1, 2016, 7:29 PM	200 KB	PDF Document
[91 Documents	🕨 🚞 js	Jan 14, 2016, 11:51 AM		Folder
O Downloads	lectures	Apr 1, 2016, 5:41 PM		Folder
D Movies	pics	Jan 18, 2016, 6:13 PM		Folder
	profiles	Jan 25, 2016, 3:32 PM		Folder
Box Sync	projects	Mar 26, 2016, 10:07 AM		Folder
Google Drive	🔻 🛅 readings	Jan 14, 2016, 11:51 AM		Folder
	endtoend.pdf	Jan 14, 2016, 11:51 AM	38 KB	PDF Document
evices	FFS84.pdf	Jan 14, 2016, 11:51 AM	1.3 MB	PDF Document
Remote Disc	garman_bug_81.pdf	Jan 14, 2016, 11:51 AM	610 KB	PDF Document
	jacobson-congestion.pdf	Jan 14, 2016, 11:51 AM	1.2 MB	PDF Document
hared	Original_Byzantine.pdf	Jan 14, 2016, 11:51 AM	1.2 MB	PDF Document
💻 adj-MBP	atterson_queue.pdf	Jan 14, 2016, 11:51 AM	1.3 MB	PDF Document
adi-mini	TheracNew.pdf	Jan 14, 2016, 11:51 AM	299 KB	PDF Document
	🔻 🔜 sections	Mar 17, 2016, 10:03 AM		Folder
🖻 tido	section1.pdf	Jan 18, 2016, 6:13 PM	130 KB	PDF Document
Al	section2.pdf	Jan 26, 2016, 7:13 PM	108 KB	PDF Document
	section2sol.pdf	Jan 28, 2016, 10:10 AM	127 KB	PDF Document
ags	section3.pdf	Feb 5, 2016, 10:15 AM	115 KB	PDF Document
	section3sol.pdf	Feb 5, 2016, 10:15 AM	134 KB	PDF Document
	a section4.pdf	Feb 10, 2016, 12:45 PM	114 KB	PDF Document
	section4sol.pdf	Feb 11, 2016, 9:42 AM	134 KB	PDF Document
	a section5.odf	Eah 16 2016 1-55 PM	109.68	PDE Document
	Accintosh HD > Documents > G	tHub F website		
		51 items, 39.01 GB available		

Files & Directories

System calls to access directories

- open / creat / readdir traverse the structure
- mkdir / rmdir add/remove entries
- link / unlink (rm)

libc support

- DIR * opendir (const char *dirname)
- struct dirent * readdir (DIR *dirstream)

Components of a File System

Superblock object: information about file system

Free bitmaps: what is allocated/not allocated

Inode object: represents a specific file

Dentry object: directory entry, single component of a path

File object: open file associated with a process.

Blocks: How files are stored on disk

Components of a File System

The (In)famous Inode

How to get the Inode number?

Look up in directory structure

Directory is a specialised file containing <file_name : inode number> mappings

File number could be a file or another directory

Each <file_name : inode> mapping is called a directory entry

How to read a file from disk

Let's read file /foo/bar.txt (Time goes downwards)

	data	inode	root	foo	bar	root	foo	bar	bar	bar
	bitmap	bitmap	inode	inode	inode	data	data	data	data	data
								[0]	[1]	[2]
			read							
						read				
open(bar)				read						
							read			
					read					
					read					
read()								read		
					write					
					read					
read()									read	
					write					
					read					
read()										read
					write					