

CS310 Operating Systems

Lecture 40 : File System Design – 1 Storage Devices and FAT

Ravi Mittal

IIT Goa

References

- CS162, Operating Systems and Systems Programming, University of California, Berkeley
- Various websites on the Internet

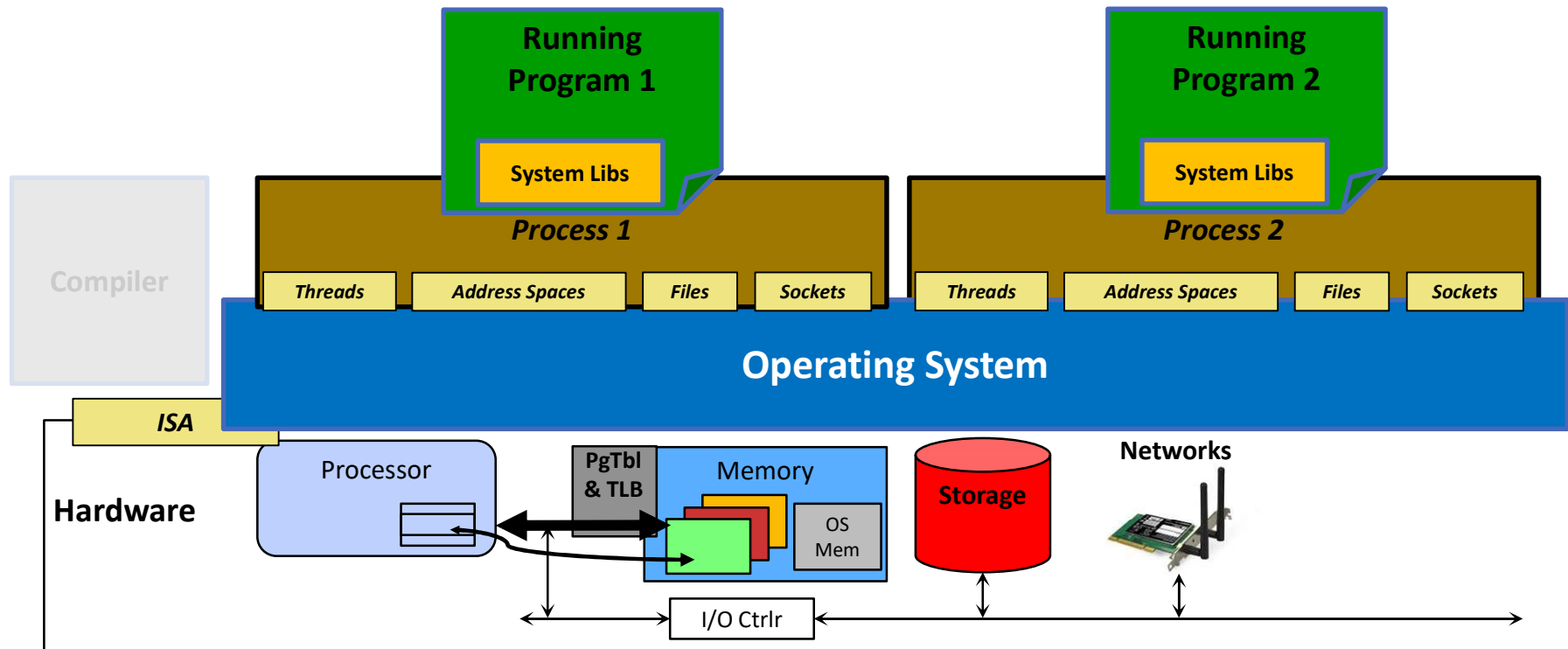
Reading

- CS162, Operating Systems and Systems Programming, University of California, Berkeley
- Book: Operating System Concepts, 10th Edition, by Silberschatz, Galvin, and Gagne

Lecture Contents

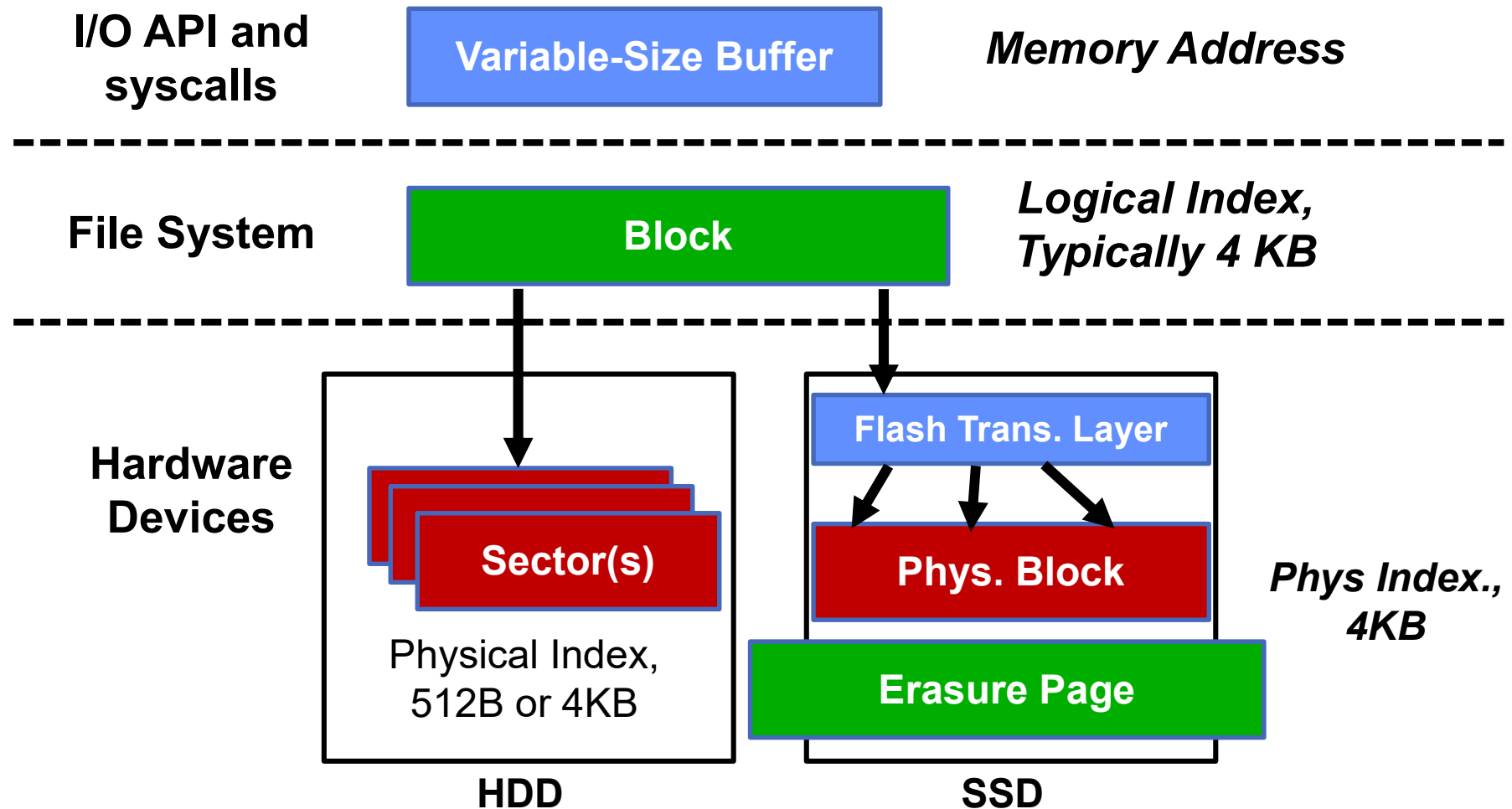
- Persistent Storage Technologies
- Hard Disk Drives
- Solid State Drive (SDD)

Where are we?



Persistent Storage Technologies

Persistent Storage for File System



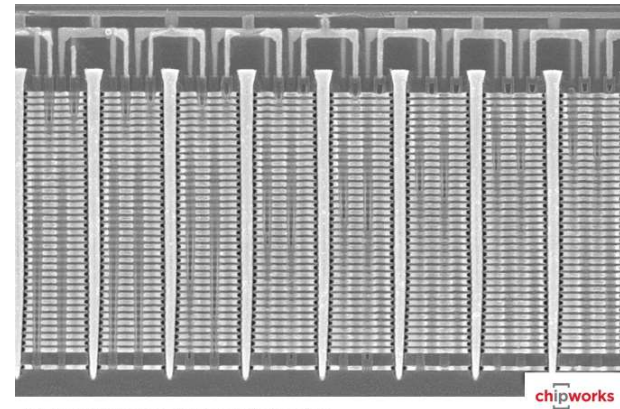
Storage Technologies

Magnetic Disks



- Store on magnetic medium
- Electromechanical access

Nonvolatile (Flash) Memory



- Store as persistent charge
- Implemented with 3-D structure
 - 100+ levels of cells
 - 3 bits data per cell

RAM vs Hard Disk vs SSD - 2018

	RAM	HDD	SSD
Typical Size	8 GB	1 TB	256 GB
Cost	\$10 per GB	\$0.05 per GB	\$0.32 per GB
Power	3 W	2.5 W	1.5 W
Read Latency	15 ns	15 ms	30 μ s
Read Speed (Seq.)	8000 MB/s	175 MB/s	550 MB/s
Read/Write Granularity	word	sector	page*
Power Reliance	volatile	non-volatile	non-volatile

In SSD Each cell has limited program/erase lifetime (thousands, for modern devices)
 – Cells become slowly less reliable

Popular Storage Devices

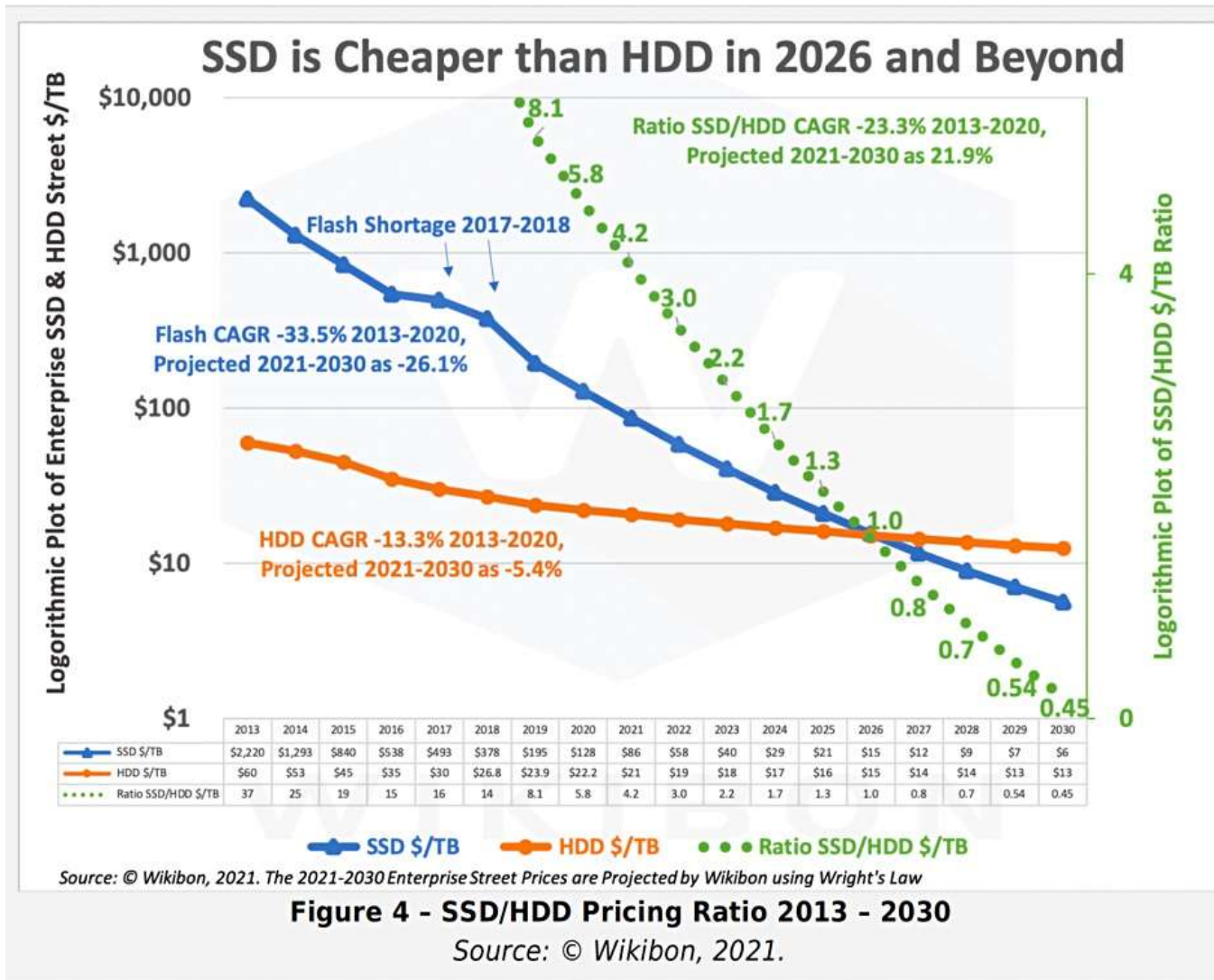
Magnetic Disks

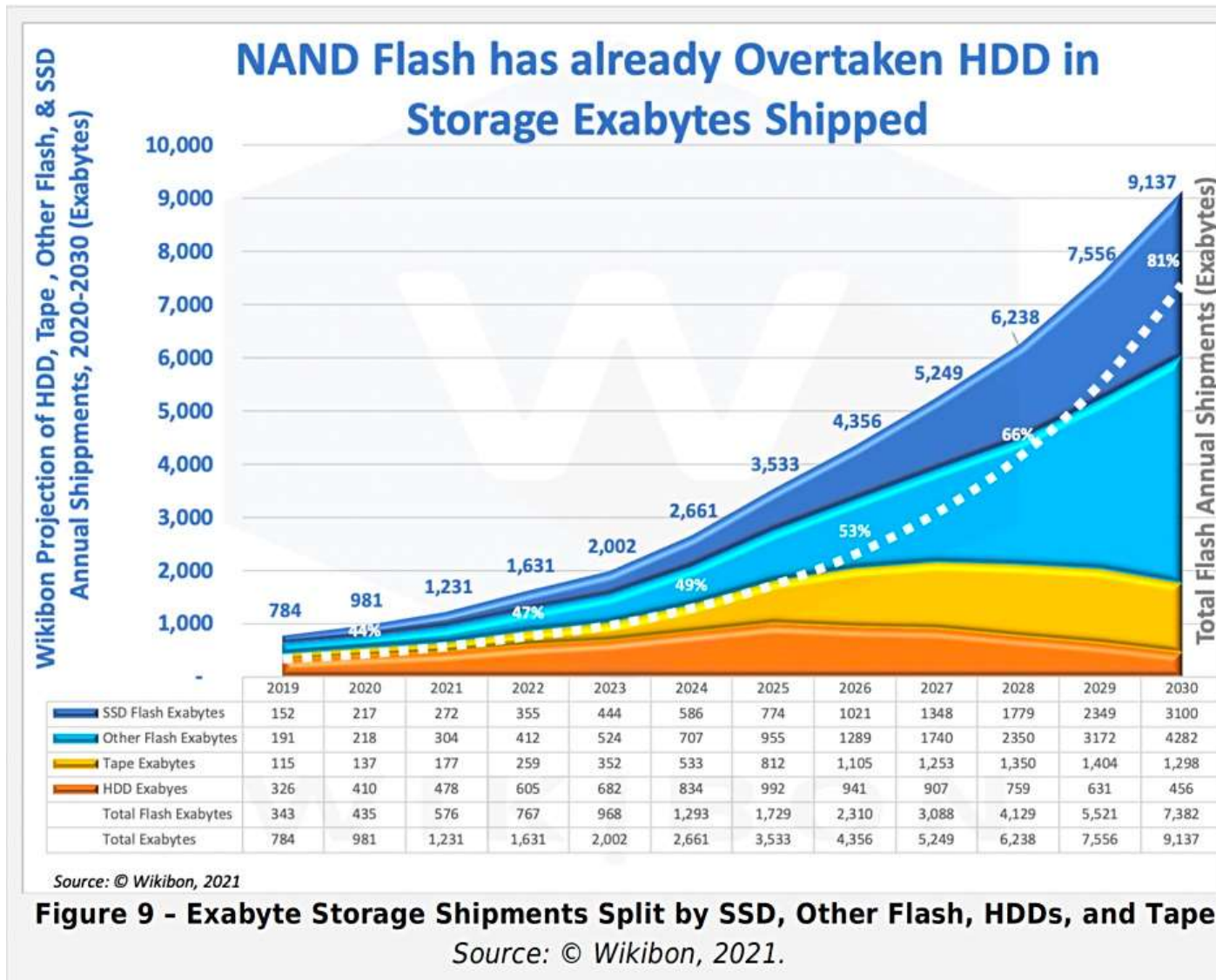
- Rarely becomes corrupted
- Traditionally: large capacity at low cost
- Block level random access
- Slow performance for random access
- Better performance for sequential access

Flash Memory

- Rarely becomes corrupted
- Increasingly larger and cheaper
- Block level random access
- Good performance for reads, worse for random writes
- Have to erase data in large blocks
- Challenge: Wear Levelling

Emergence of SSDs





The Emergence of SSDs

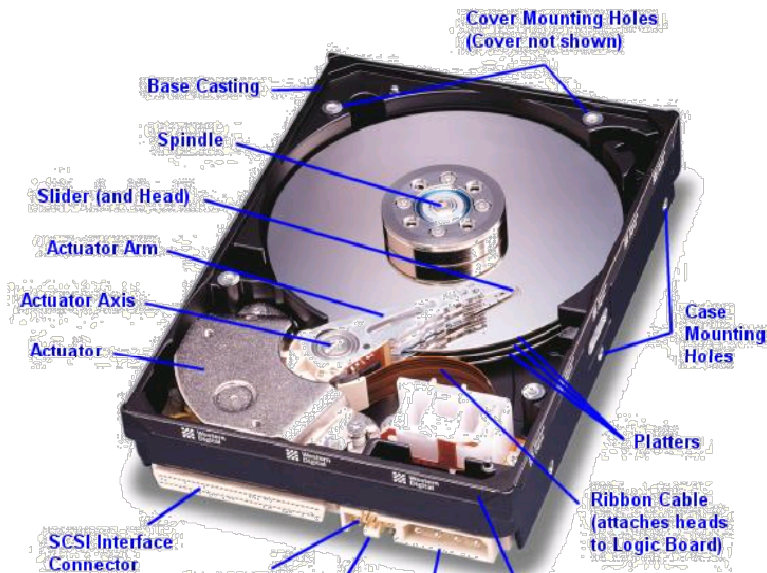
- Faster
- Lower power
- No moving parts
- But HDDs have their place
 - Cheapest online storage per byte
 - Application: Archival storage



SSD vs HDD		
Usually 10 000 or 15 000 rpm SAS drives		
0.1 ms	Access times SSDs exhibit virtually no access time	5.5 ~ 8.0 ms
SSDs deliver at least 6000 io/s	Random I/O Performance SSDs are at least 15 times faster than HDDs	HDDs reach up to 400 io/s
SSDs have a failure rate of less than 0.5 %	Reliability This makes SSDs 4 - 10 times more reliable	HDD's failure rate fluctuates between 2 ~ 5 %
SSDs consume between 2 & 5 watts	Energy savings This means that on a large server like ours, approximately 100 watts are saved	HDDs consume between 6 & 15 watts
SSDs have an average I/O wait of 1 %	CPU Power You will have an extra 6% of CPU power for other operations	HDDs' average I/O wait is about 7 %
the average service time for an I/O request while running a backup remains below 20 ms	Input/Output request times SSDs allow for much faster data access	the I/O request time with HDDs during backup rises up to 400 ~ 500 ms
SSD backups take about 6 hours	Backup Rates SSDs allows for 3 - 5 times faster backups for your data	HDD backups take up to 20 ~ 24 hours

Hard Disk Drives (HDD)

Hard Disk Drivers (HDDs)

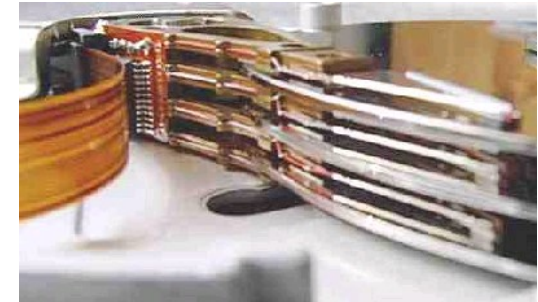


Western Digital Drive

<http://www.storagereview.com/guide/>



IBM/Hitachi Microdrive

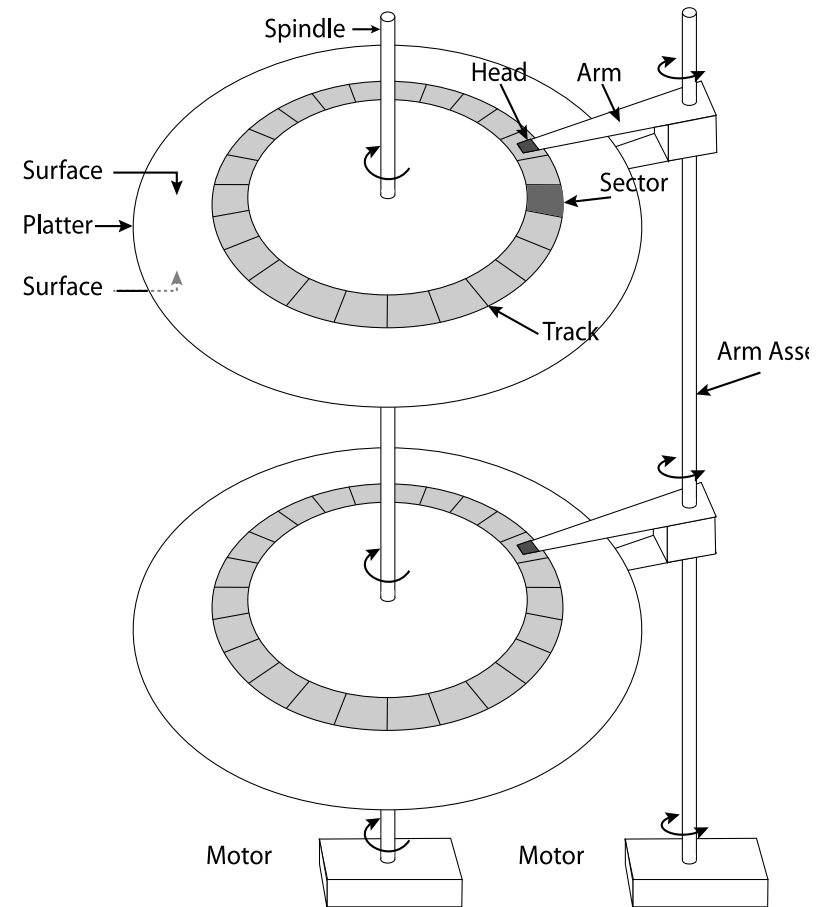


Read/Write Head
Side View

IBM Personal Computer/AT (1986)
30 MB hard disk - \$500
30-40ms seek time
0.7-1 MB/s (est.)

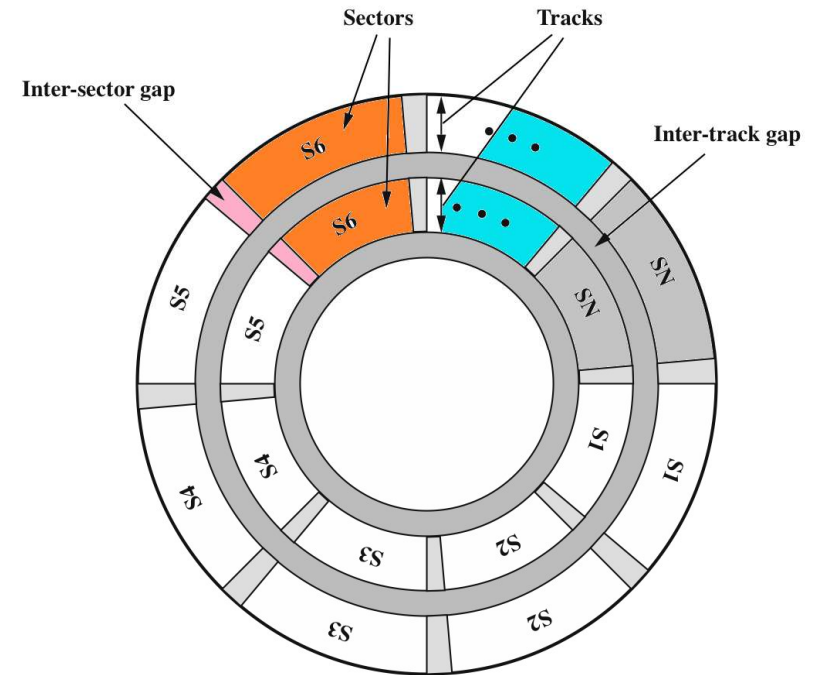
The Amazing Magnetic Disk

- Unit of Transfer: **Sector**
 - Ring of sectors form a track
 - Stack of tracks form a cylinder
 - Heads position on cylinders
- **Disk Tracks ~ $1\mu\text{m}$ (micron) wide**
 - Wavelength of light is ~ $0.5\mu\text{m}$
 - Resolution of human eye: $50\mu\text{m}$
 - 100K tracks on a typical 2.5" disk
- **Separated by unused guard regions**
 - Reduces likelihood neighboring tracks are corrupted during writes (still a small non-zero chance)



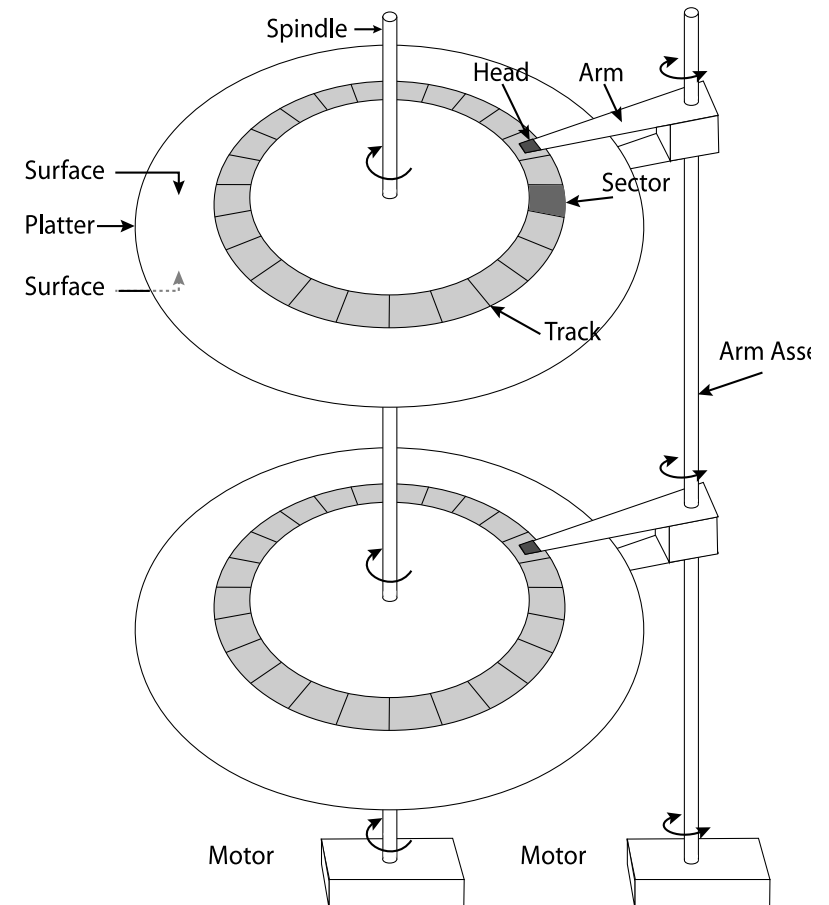
Disk platter – top view

- Typically sector = 512 bytes
- Minimum reading on disk: 1 sector
 - Can't read individual byte or word



The Amazing Magnetic Disk

- Track length varies across disk
 - Outside: More sectors per track, higher bandwidth
 - Disk is organized into regions of tracks with same # of sectors/track
 - Only outer half of radius is used
 - Most of the disk area in the outer regions of the disk



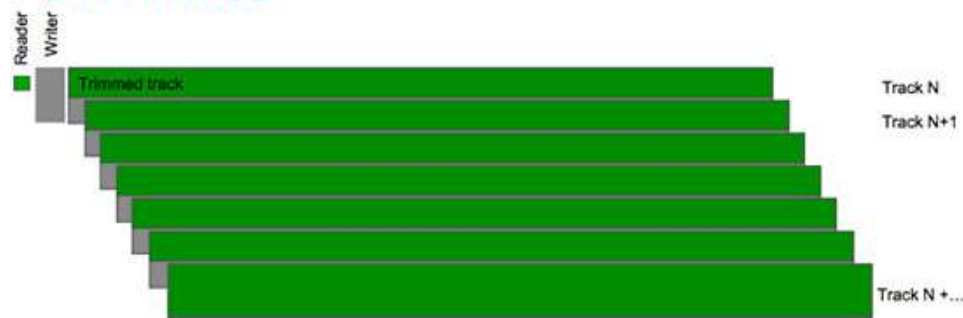
Shingled Magnetic Recording (SMR)

- Shingled magnetic recording is a magnetic storage data recording technology used in hard disk drives to increase storage density and overall per-drive storage capacity
- Overlapping tracks yields greater density, capacity
- Restrictions on writing, complex DSP for reading

Conventional Writes

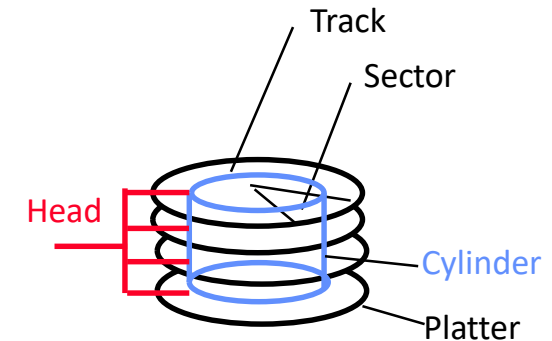


SMR Writes

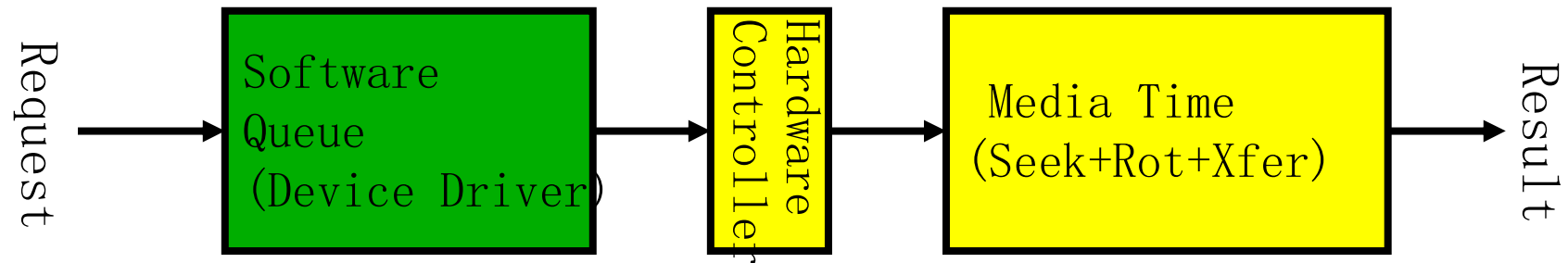


Review: Magnetic Disks

- **Cylinders:** all the tracks under the head at a given point on all surfaces
- Read/write data is a three-stage process:
 - **Seek time:** position the head/arm over the proper track –Average of 5-10 ms
 - **Rotational latency:** wait for desired sector to rotate under r/w head
 - 4-8 ms (3600-7200 rpm), 2-4 ms (15000ms)
 - **Transfer time:** Time to actually read sectors
 - 50-100 MB/sec



Bigger Picture



Latency = Queue Time + Controller Time +
Seek Time + Rotational Latency +
Transfer Time

To Achieve Best Bandwidth: Large Transfers of Physically Adjacent Sectors from one track

Disk Performance Example – Self reading

- Assumptions:
 - Ignoring queuing and controller times for now
 - Avg seek time of 5ms,
 - 7200RPM \Rightarrow Time for rotation: $60000 \text{ (ms/min)} / 7200 \text{ (rev/min)} \approx 8\text{ms}$
 - Transfer rate of 50MByte/s, block size of 4Kbyte \Rightarrow
 $4096 \text{ bytes} / 50 \times 10^6 \text{ (bytes/s)} = 81.92 \times 10^{-6} \text{ sec} \cong 0.082 \text{ ms per block}$
- Read block from random place on disk:
 - Seek (5ms) + Rot. Delay (4ms) + Transfer (0.082ms) = 9.082ms
 - Approx 9ms to fetch/put data: $4096 \text{ bytes} / 9.082 \times 10^{-3} \text{ s} \cong 451\text{KB/s}$
- Read block from random place in same cylinder:
 - Rot. Delay (4ms) + Transfer (0.082ms) = 4.082ms
 - Approx 4ms to fetch/put data: $4096 \text{ bytes} / 4.082 \times 10^{-3} \text{ s} \cong 1.03\text{MB/s}$
- Read next block on same track:
 - Transfer (0.082ms): $4096 \text{ bytes} / 0.082 \times 10^{-3} \text{ s} \cong 50\text{MB/sec}$
- Key to using disk effectively (especially for file systems) is to minimize seek and rotational delays

HDD Controller

HDD Controllers

- Old Days: Device driver would address block of data by cylinder number, head (platter) number, and sector number
- Now: Hard drive is just an array of sectors
 - Sector number mapped internally to physical location
 - Numerically close sectors are probably physically close
- Lots of other intelligent features
 - Error Correcting Codes
 - Sector Sparing
 - Slip Sparing
 - Track Skewing

Intelligence in the Controller

- Error Correcting Codes

- Disk encodes each sector with additional error correcting code data
 - allowing it to fix imperfectly read or written data

- Sector sparing

- Disks are made with tracks and sectors as thin and small as possible
- If there are imperfections in a sector, then the sector can not be used for reliably store data
- Manufacturers include spare sectors distributed across each surface
- Disk firmware or formatting can 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
- Helps in retaining good sequential performance by remapping all sectors from the bad sector to the next spare

Intelligence in the Controller

- Track Skewing
 - Sector numbers offset from one track to the next, to allow for disk head movement for sequential operations
 - Logical sector zero on each track is staggered from sector zero on the previous track
 - By an amount corresponding to time it takes the disk to move the head from one track to another

Typical Numbers for Magnetic Disk (self reading)

Parameter	Info/Range
Space/Density	Space: 14TB (Seagate), 8 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: <ul style="list-style-type: none">• 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

Example of Current HDDs (Self reading)

- Seagate Exos X14 (2018)
 - 14 TB hard disk
 - 8 platters, 16 heads
 - Helium filled: reduce friction and power
 - 4.16ms average seek time
 - 4096 byte physical sectors
 - 7200 RPMs
 - 6 Gbps SATA /12Gbps SAS interface
 - 261MB/s MAX transfer rate
 - Cache size: 256MB
 - Price: \$615 (< \$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 !!)



Lecture Summary

- Hard Disk Drives – had magical run
- Tremendous improvement in HDD technology
 - Smaller sizes
 - More information
 - Less power
 - Low cost
 - Improved Performance
 - Higher reliability
 - A lot of intelligence
- SSD are now becoming main stream