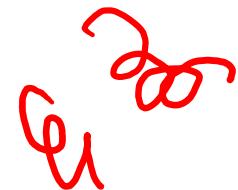


Chapter 12: Mass-Storage Systems

Prof. Li-Pin Chang

CS@NYCU



Chapter 12: Mass-Storage Systems

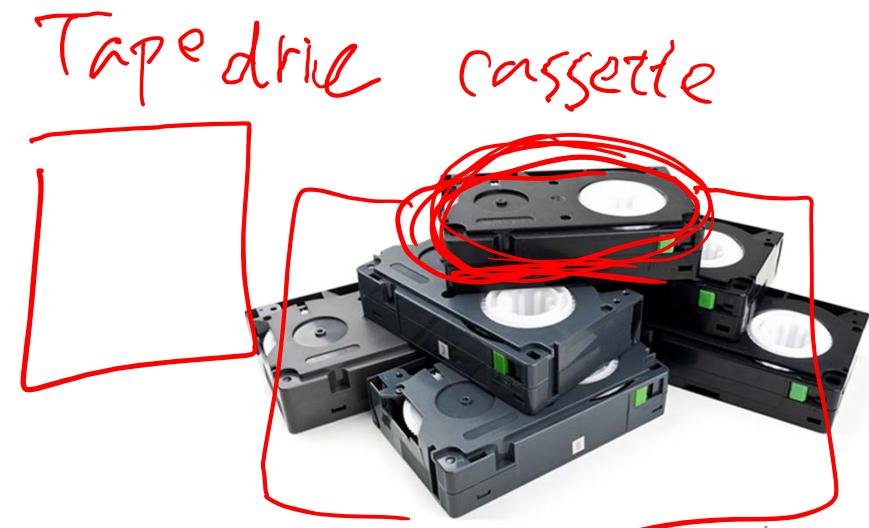
- Magnetic Tape and Disk
- Disk Scheduling
- RAID Structure
- Solid State Disks

Objectives

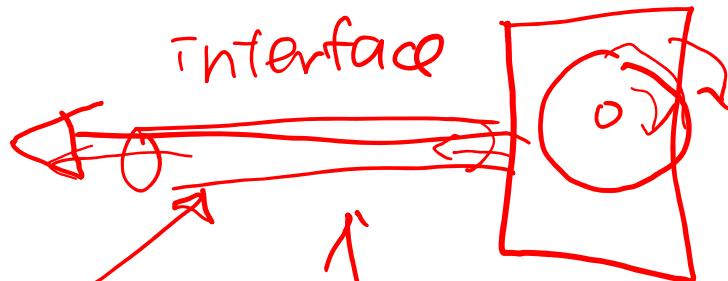
- Describe the physical structure of secondary and tertiary storage devices and the resulting effects on the uses of the devices
- Explain the performance characteristics of mass-storage devices and performance optimization techniques
- Discuss operating-system services provided for mass storage, such as RAID

Magnetic Tape

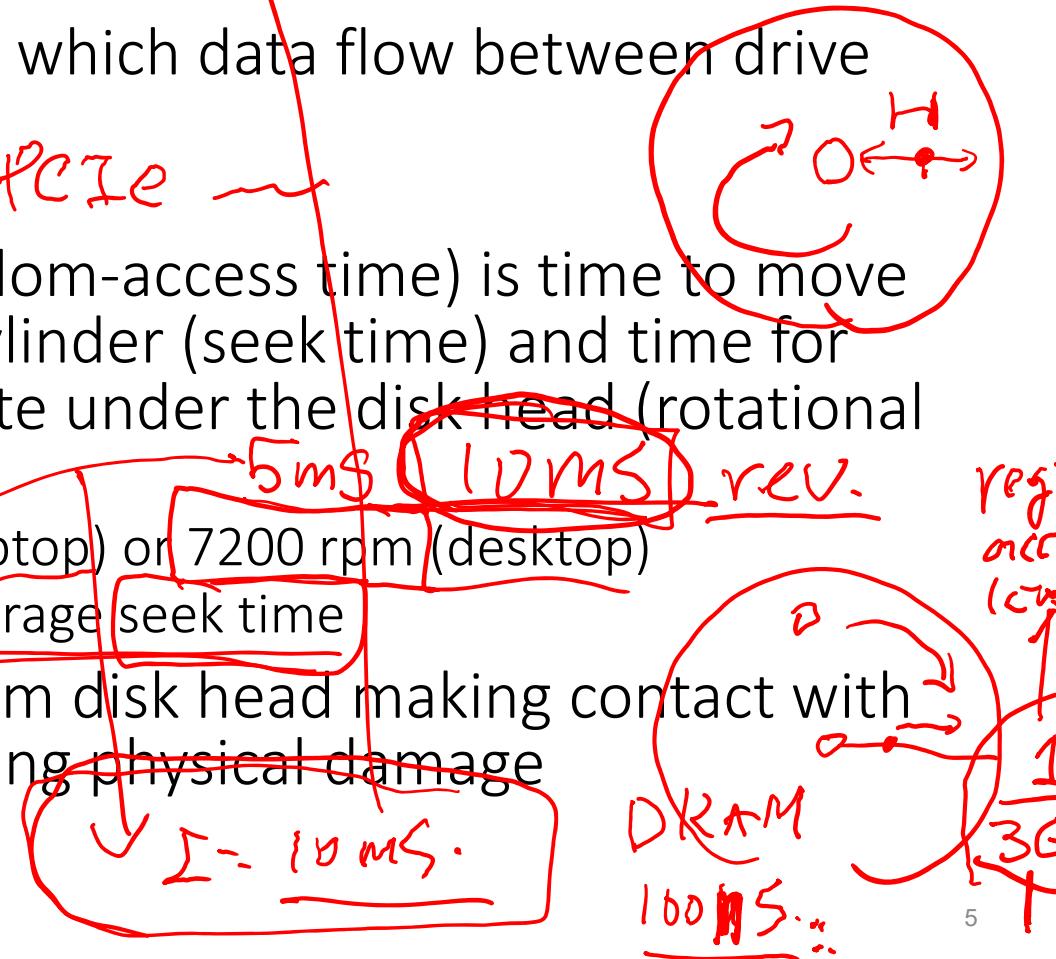
- Very fast sequential read-write, achieving 400 MB/s (LTO-9); extremely slow on random access
 - Kept in spool and wound or rewound past read-write head
 - Once data under head, transfer rates comparable to disk
- Relatively permanent and holds large quantities of data
 - 18 TB per cassette (LTO-9)
 - For data backup



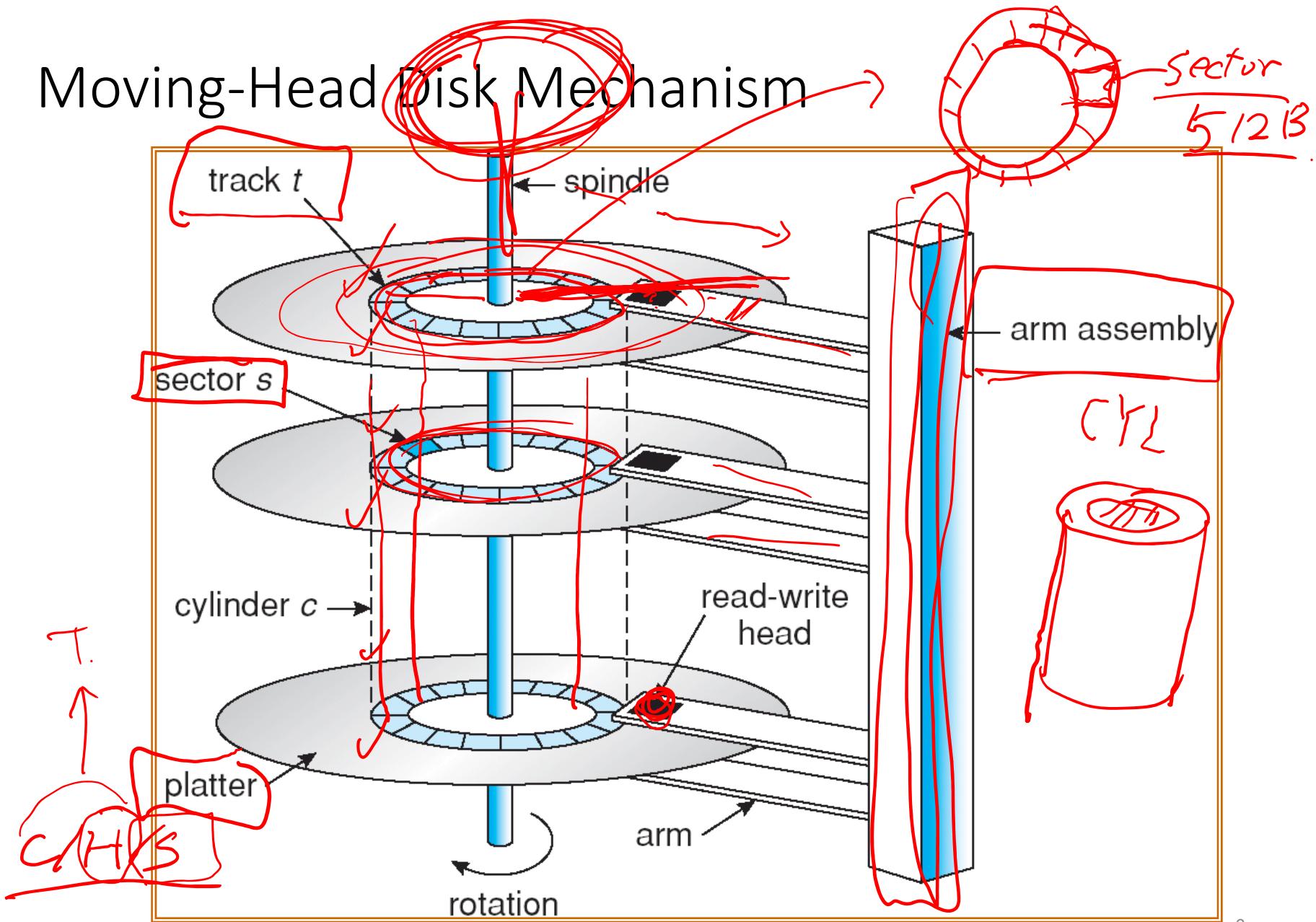
Magnetic Disks



- Provide bulk of secondary storage of modern computers
- **Transfer rate** is rate at which data flow between drive and computer
 - SATA3: 600 MB/s
- **Positioning time** (random-access time) is time to move disk arm to desired cylinder (seek time) and time for desired sector to rotate under the disk head (rotational latency)
 - Typically 5400 rpm (laptop) or 7200 rpm (desktop)
 - Typically 5ms~7ms average seek time
- Head crash results from disk head making contact with the disk surface, causing physical damage



Moving-Head Disk Mechanism

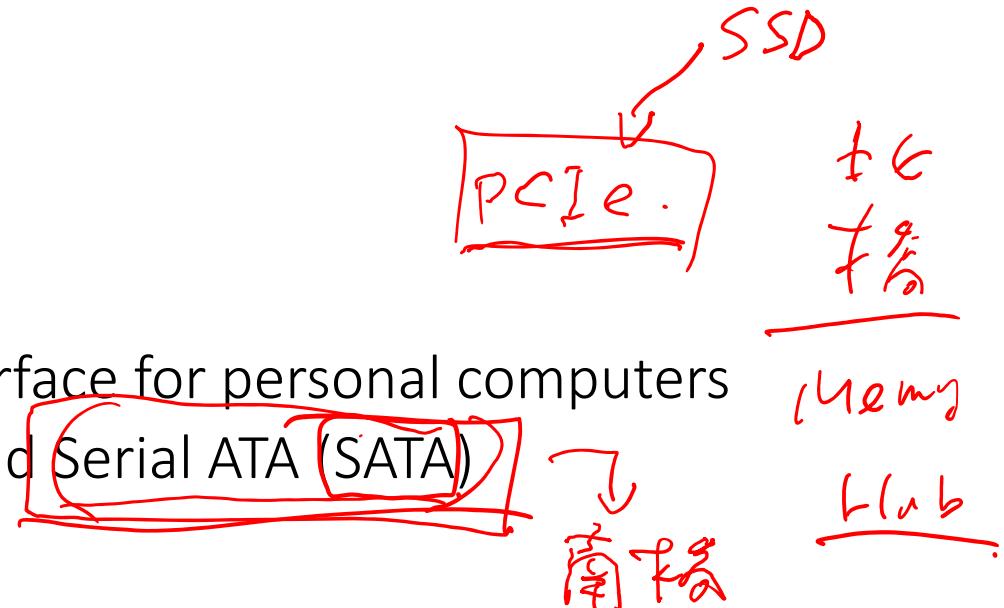


Disk Structure

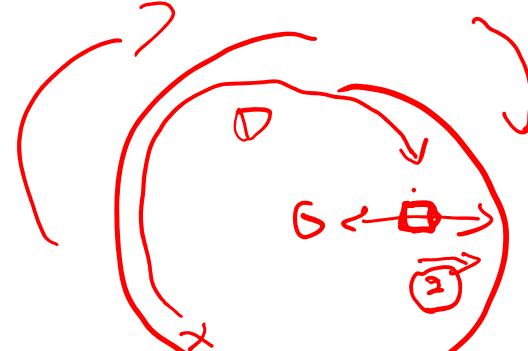
- Disk drives are addressed as large 1-dimensional arrays of logical blocks, where the logical block is the smallest unit of transfer.
- The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially.
 - Sector 0 is the first sector of the first track on the outermost cylinder.
 - Mapping proceeds in order through that track, then the rest of the tracks in that cylinder, and then through the rest of the cylinders from outermost to innermost.

Disk Attachment

- ATA/IDE interface
 - the primary disk interface for personal computers
 - Parallel ATA (PATA) and Serial ATA (SATA)
- SCSI bus
 - Up to 16 devices (disks, printers, etc) on one cable
- FC (Fiber Channel) is high-speed serial architecture
 - The basis of Storage Area Networks (SANs) in which many hosts attach to many storage units



Disk Scheduling



- The operating system is responsible for using hardware efficiently — for the disk drives, this means having a fast access time and disk bandwidth.

- Access time has two major components

OS • Seek time is the time for the disk are to move the heads to the cylinder containing the desired sector.

• Rotational latency is the additional time waiting for the disk to rotate the desired sector to the disk head.

- Disk scheduling optimizes seek time, which is proportional to the total seek distance

seek t ~~at~~ seek dist.

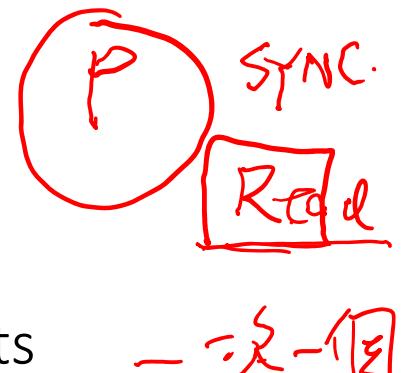


The Need for Disk Scheduling

- Because of

- 1) multiprogramming and
- 2) write buffering, ~~ASOC~~ ~~LISTFC~~

there might be a number of pending disk requests



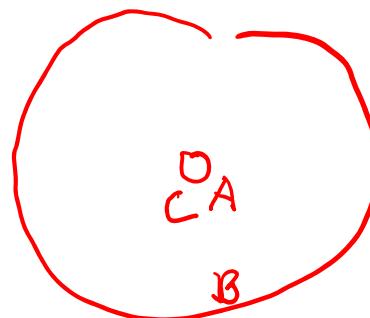
- - - - -



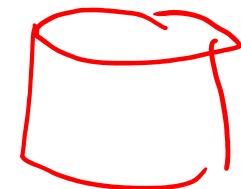
- How to select the next request to serve?

- has impacts on response and throughput

TSP

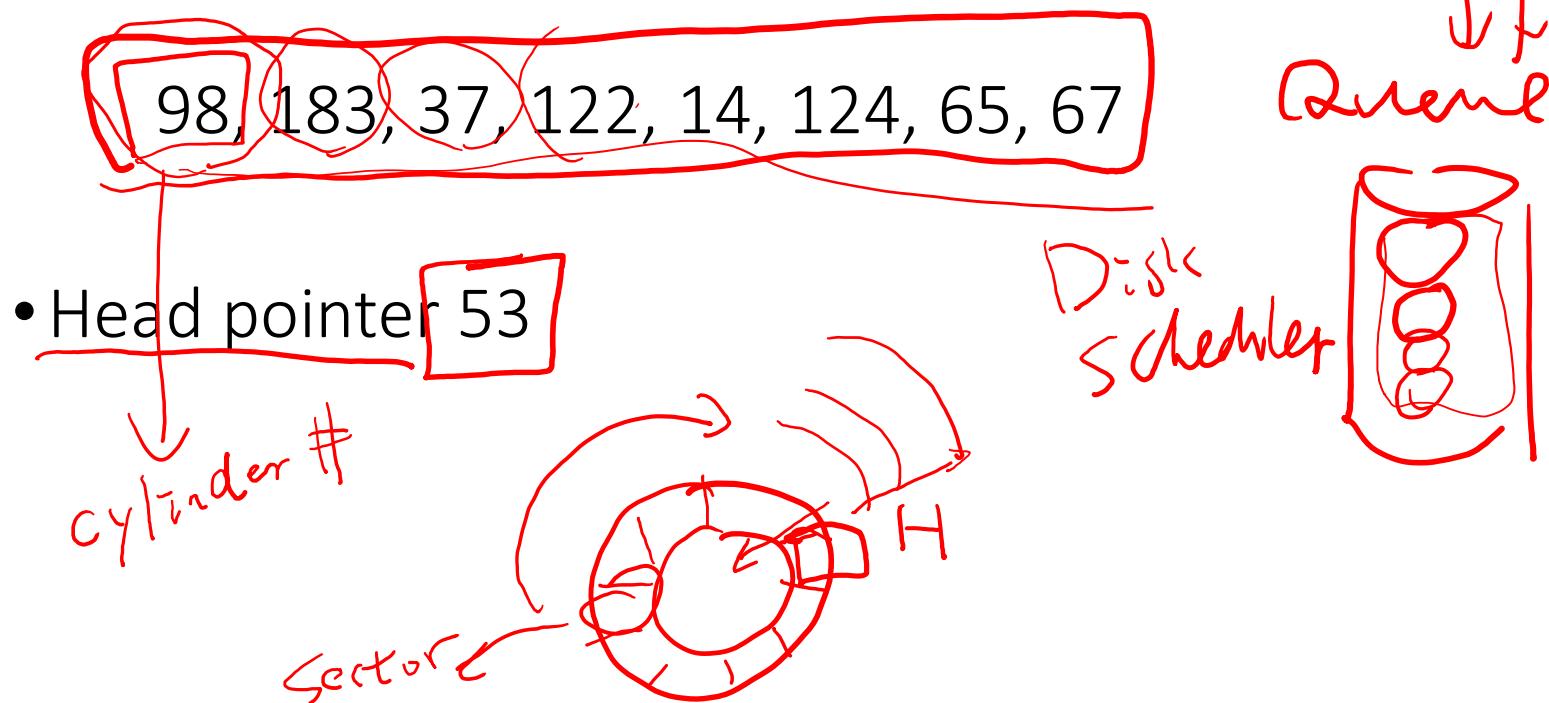
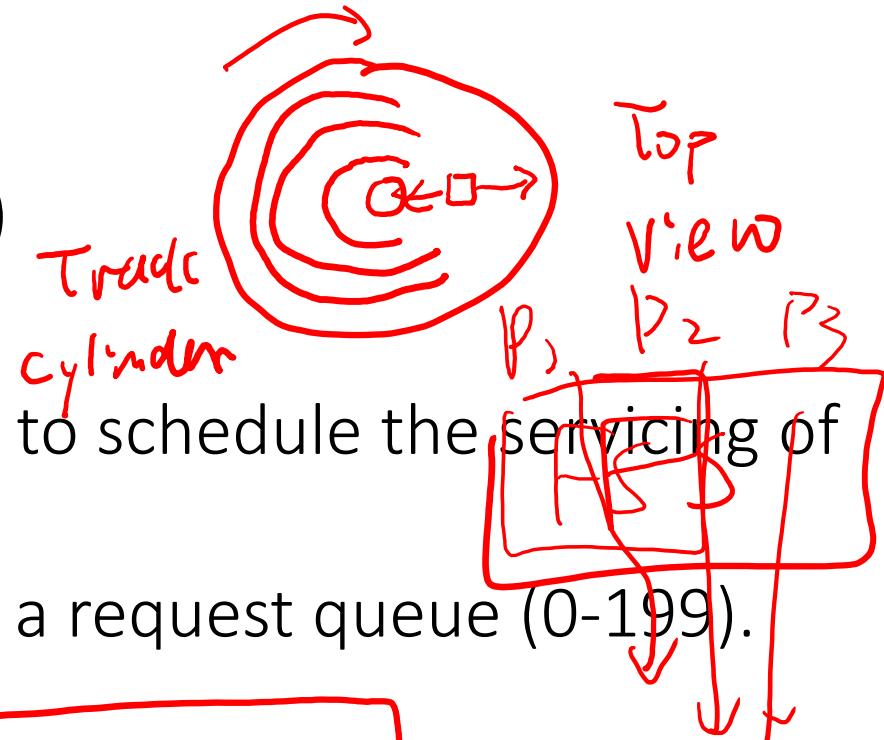


$A \rightarrow B \rightarrow C$
 $A \rightarrow C \rightarrow B$.



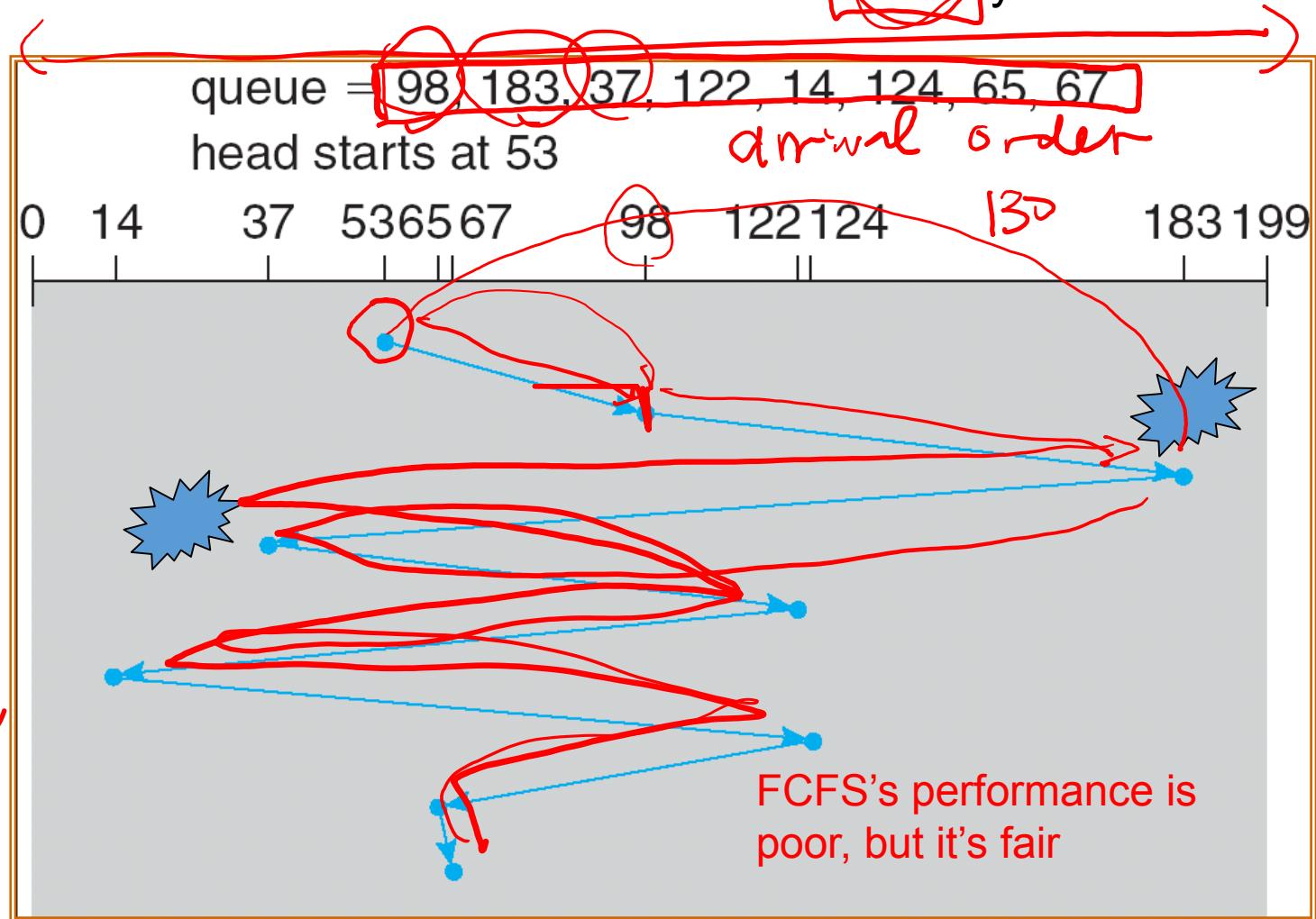
Disk Scheduling (Cont.)

- Several algorithms exist to schedule the servicing of disk I/O requests.
- We illustrate them with a request queue (0-199).



FCFS Disk Scheduling

cylinder (track)
Illustration shows total head movement of 640 cylinders.



SSTF Scheduling

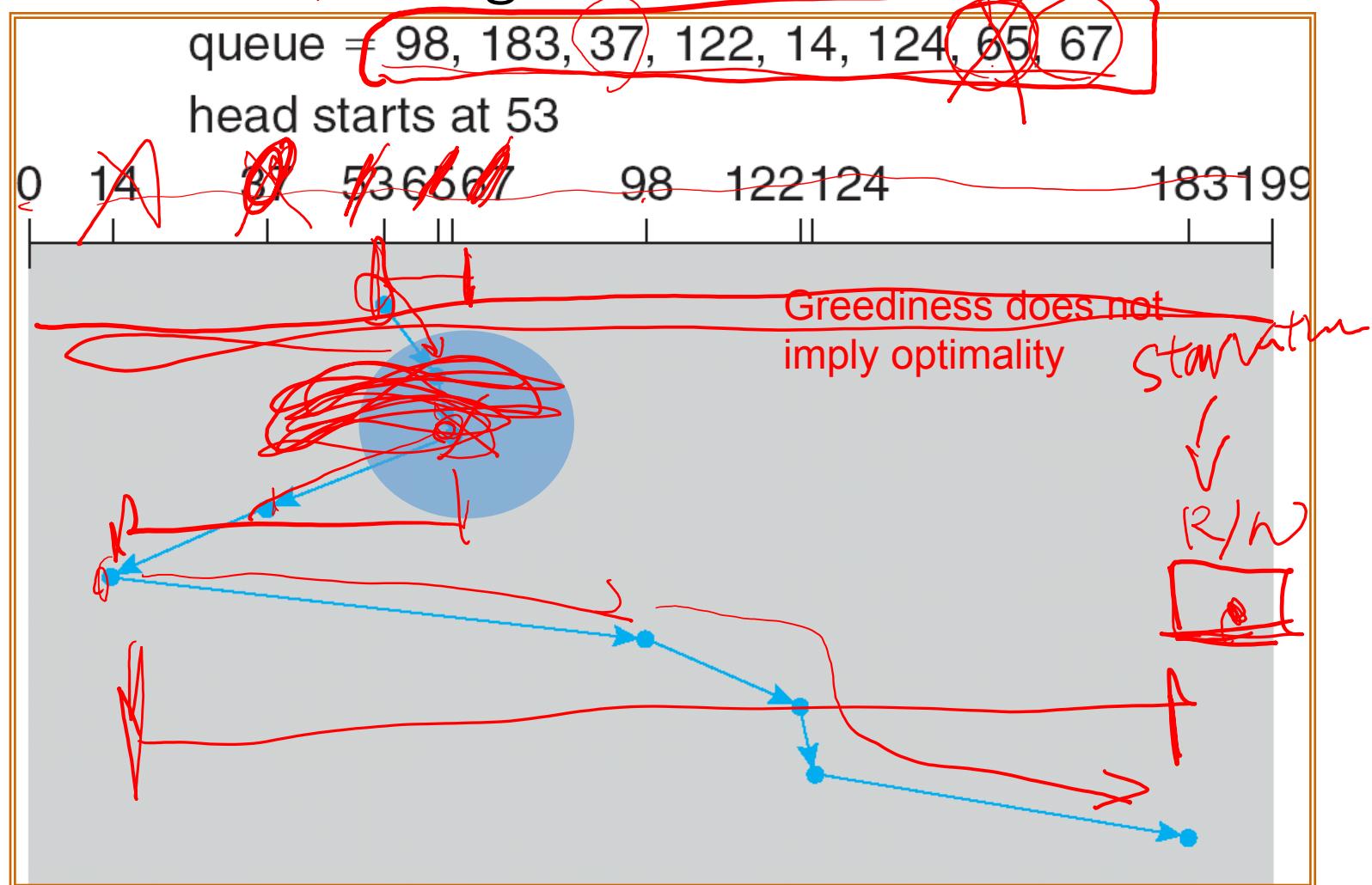
(Greedy)

Shortest Seek Time First

- Selects the request with the minimum seek time from the current head position
- SSTF scheduling is a form of SJF scheduling; may cause **starvation** of some requests
 - How to avoid starvation?
- Illustration shows total head movement of cylinders



SSTF Disk Scheduling



SCAN Scheduling

- The disk arm starts at one end of the disk, and moves toward the other end, servicing requests until it gets to the other end of the disk, where the head movement is reversed and servicing continues
- Sometimes it is called the **elevator** algorithm
- Illustration shows total head movement of cylinders

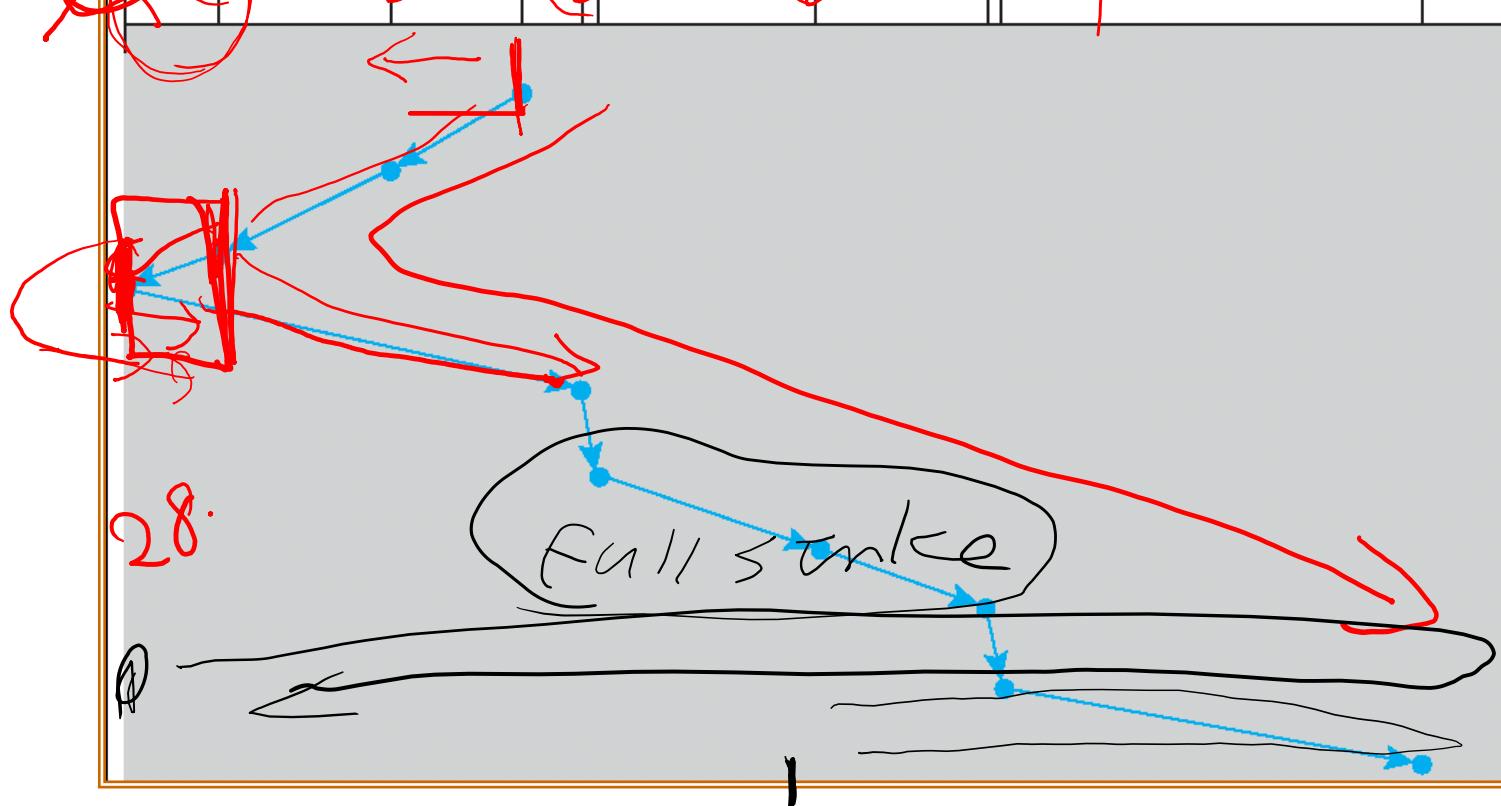
Look → 236
208

SCAN Disk Scheduling

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

0 14 37 53 65 67 98 122 124 183 199



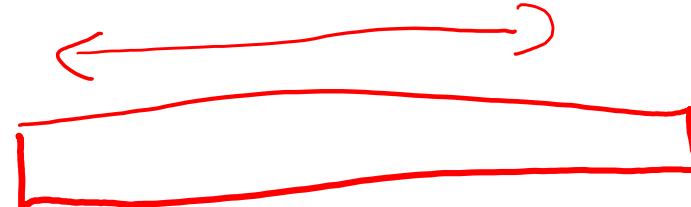
SCAN



- Much less prone to starvation: bounded waiting if all requests are known
- The waiting time of each cylinder is not uniform
 - At the outermost or the innermost cylinder:
 - Max 2 full strokes and 1 reverse
 - At the middle of the disk:
 - Max 2 half full strokes and 1 reverse

C-SCAN Scheduling

↳ circular



- Provides a more uniform wait time than SCAN
- The head moves from one end of the disk to the other, servicing requests as it goes. When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
- Treats the cylinders as a circular list that wraps around from the last cylinder to the first one

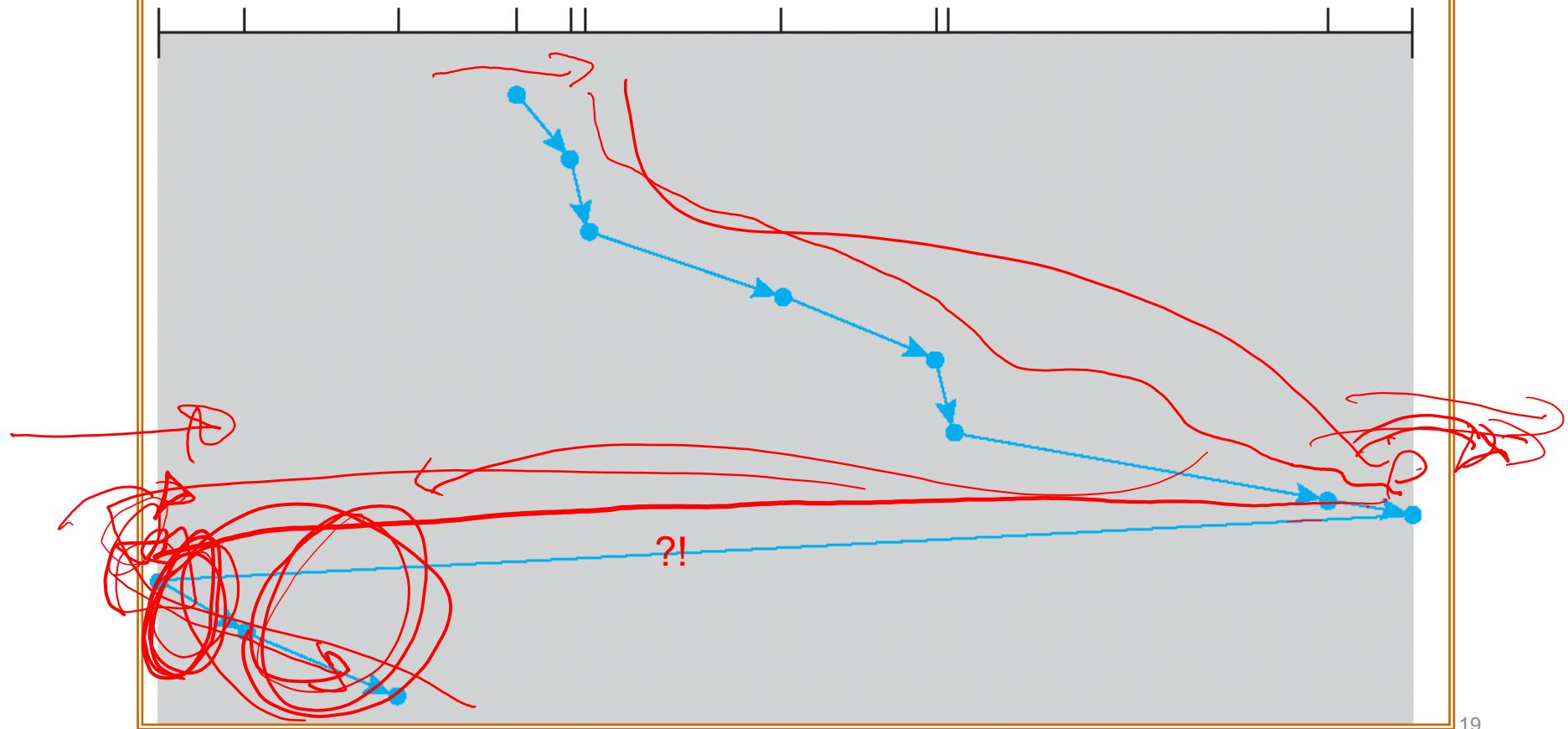


C-SCAN Disk Scheduling

queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

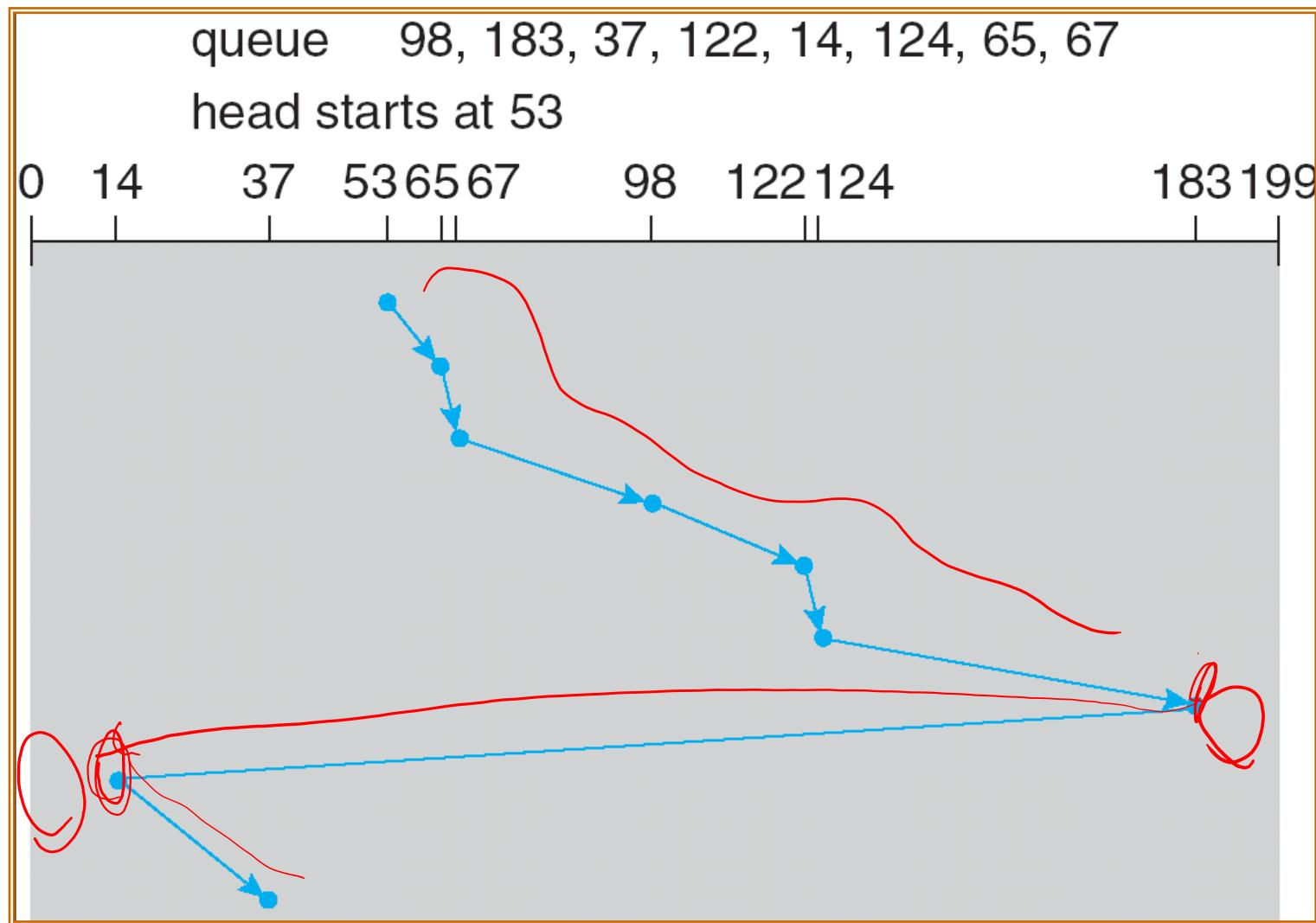
0 14 37 53 65 67 98 122 124 183 199



C-LOOK

- Version of C-SCAN
- Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk.

C-LOOK DISK Scheduling

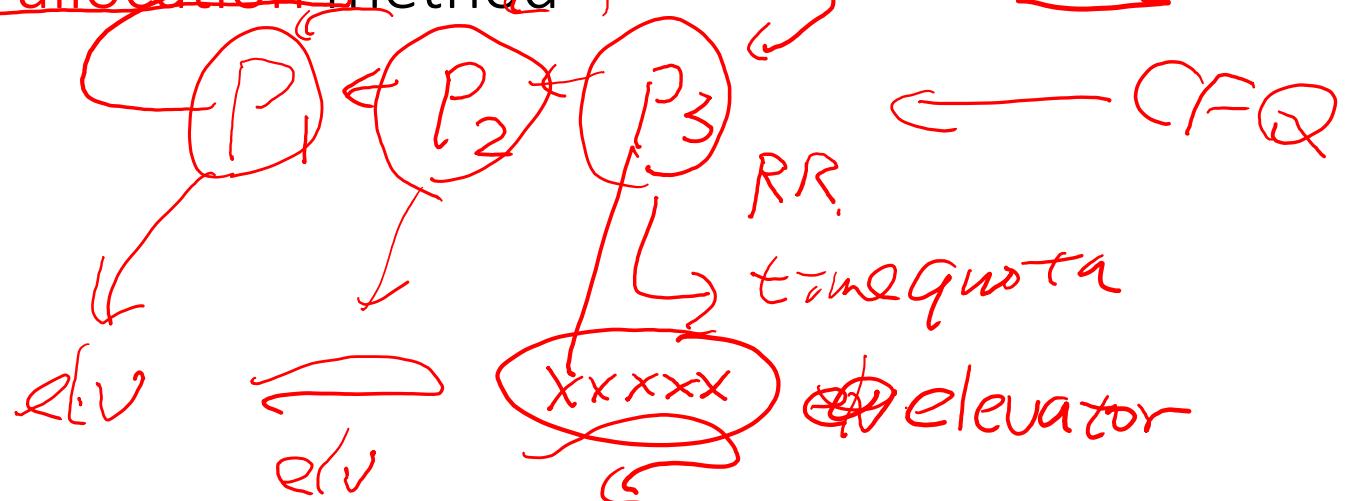


M. U

ext

Selection of a Disk-Scheduling Algorithm

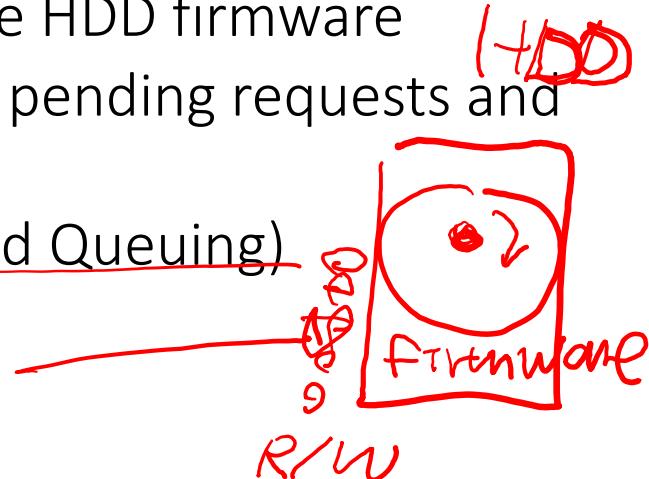
- SSTF has a natural appeal, but it risks starvation
- SCAN and C-SCAN perform better for systems that place a heavy load on the disk; either one is a reasonable default choice
- Requests for disk service can be influenced by the file-allocation method



Disk Seek Optimization

- Disk scheduling problem is NP-hard
 - All the methods mentioned above are not optimal (in terms of the total seek distance)
- Hard to optimize rotational delay from operating systems
 - The HDD firmware knows the current rotation angle better
 - Delegating disk scheduling to the HDD firmware
 - New HDDs accept a number of pending requests and then reorder them internally
 - E.g., SATA NCQ / Native Command Queuing

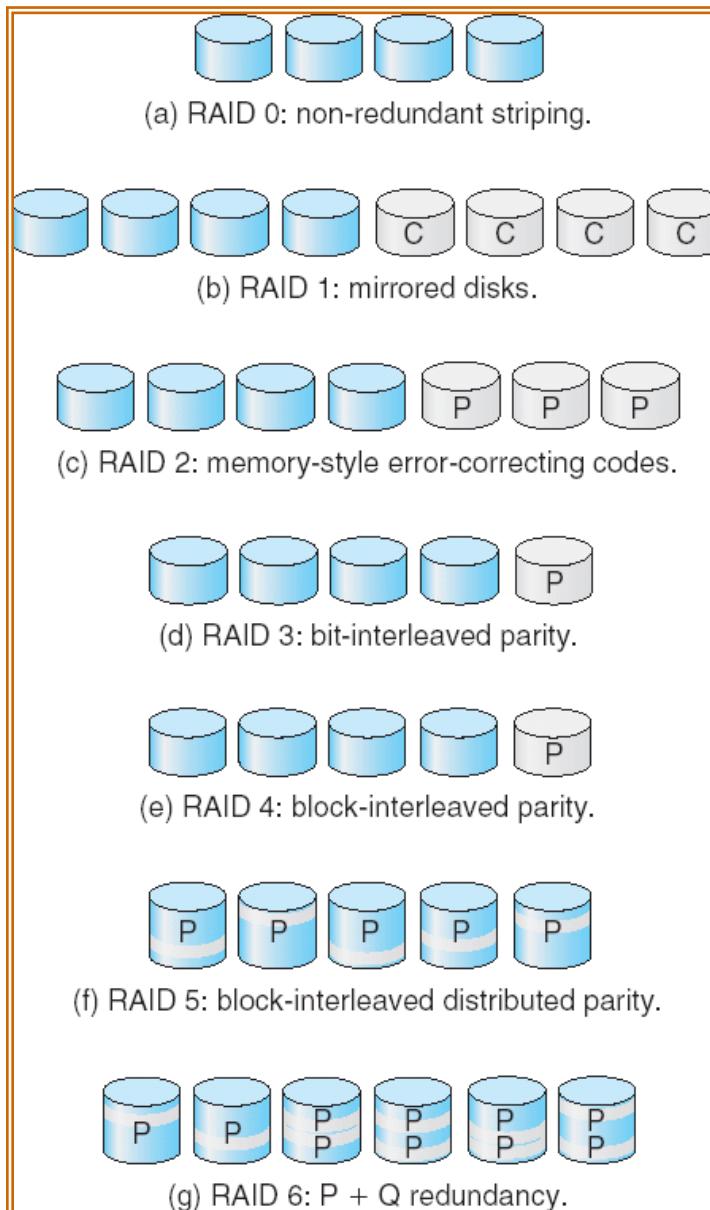
SATA

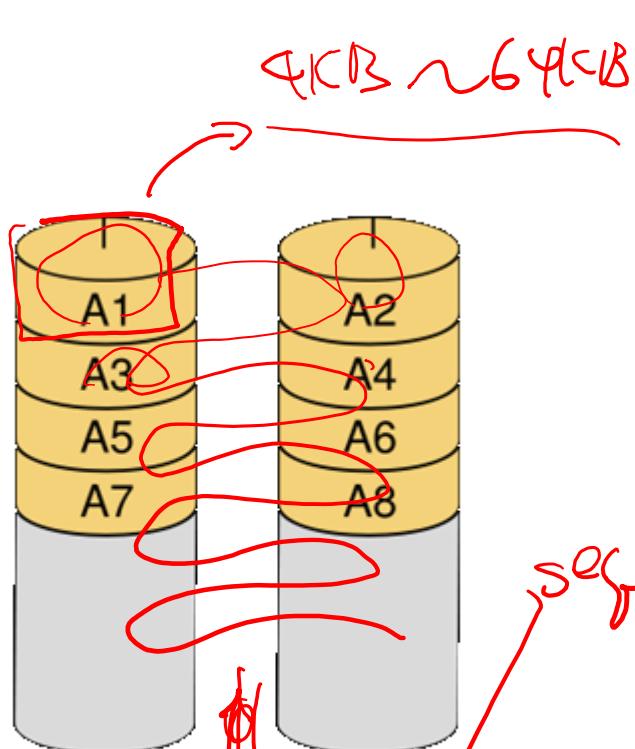


RAID Structure

- RAID – Redundant Array of Inexpensive Disks
 - Performance improvement through parallelism
 - Reliability improvement through redundancy
- RAID schemes improve performance and improve the reliability of the storage system by storing redundant data.
 - Mirroring or shadowing keeps duplicate of each disk.
 - Block interleaved parity uses much less redundancy.

RAID Levels





①

A5 & A4

seg read

②

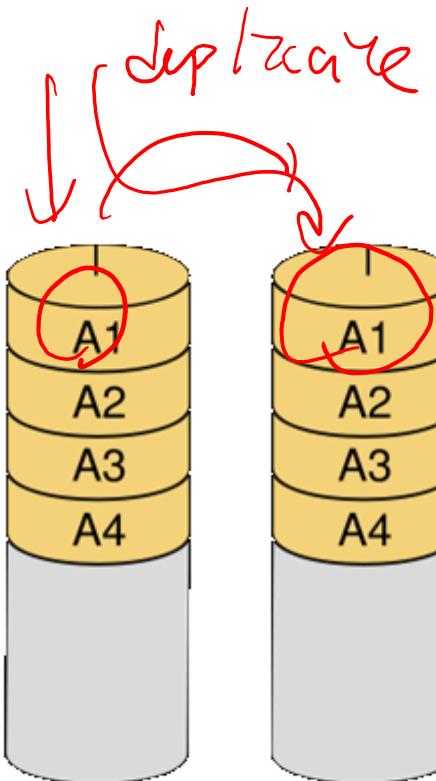
A7 RAID0

striping block

A2 ↗ RWD

- Block striping
- Best performance (high parallelism)
- Less reliable

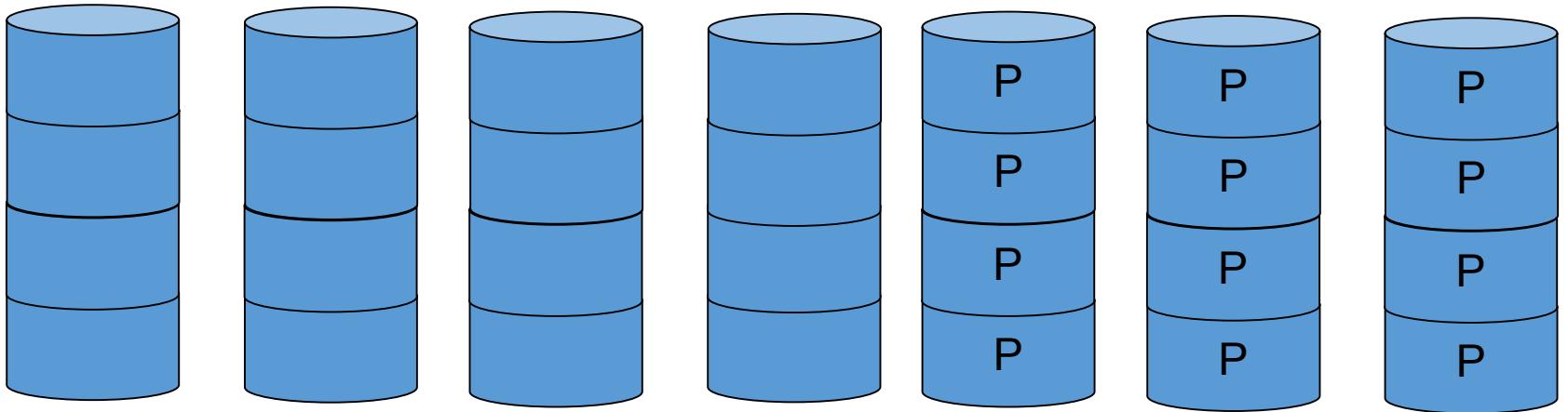
striping



- Mirroring, 100% redundancy
- Less efficient
- More reliable

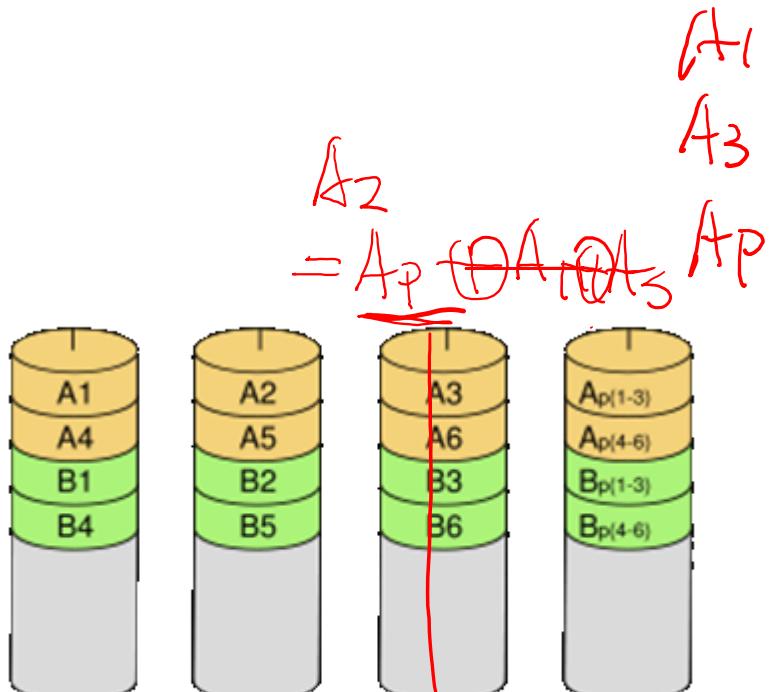
RAID1

Mirror



RAID-2

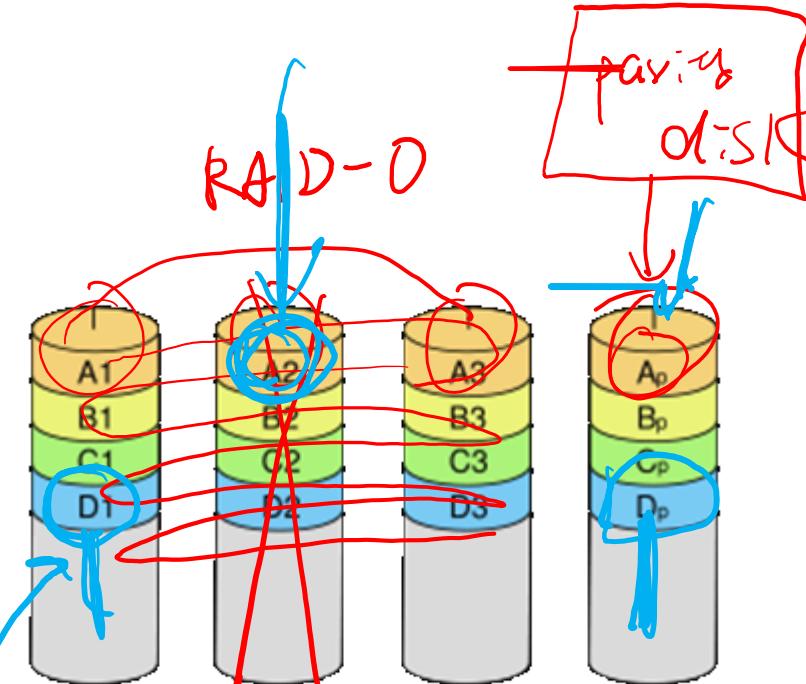
- Memory-style ECC, such as Hamming code
- n parity bit covers up to $2^n - 1$ bits (including parity)
- e.g., 3 parity disks are needed for a 7-disk array



$A_1 \oplus A_2 \oplus A_3 \oplus A_p = A_p = A_1 \oplus A_2 \oplus A_3$

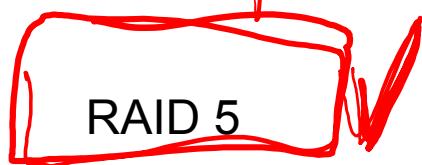
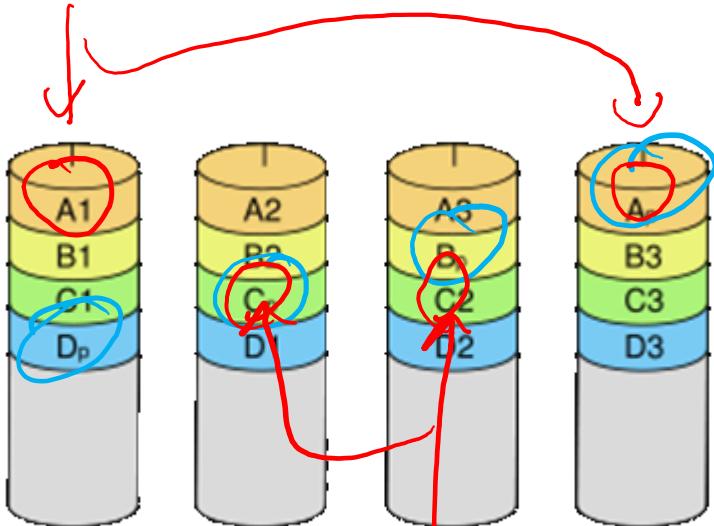
- Sub-block-interleaved
- Fully interleaved, one R/W involves all disks

$$D_p' = D_p \oplus D_1 \oplus D_2 \oplus D_3$$

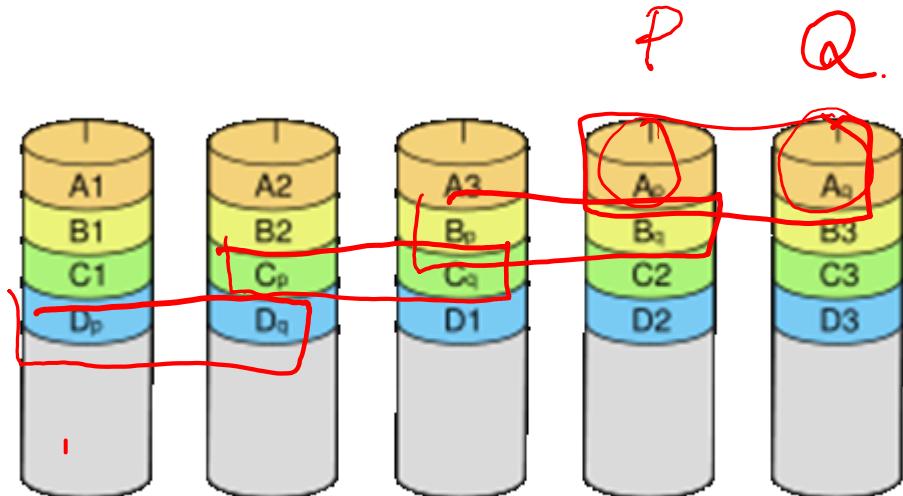


RAID4

- Block-interleaved with parity blocks in the parity disk
- XOR-based parity
- One R involves one disk
- One W involves two disks
- Parity disk → bottleneck



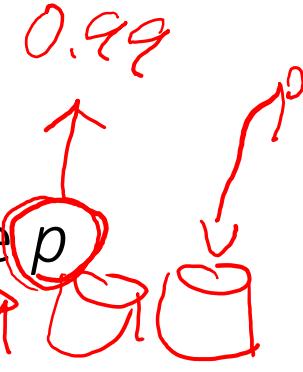
- One R involves one disk
- One W involves two disks
- Parity blocks are distributed among disks
- Simple XOR-based parity
- Software RAID-5 is possible



RAID 6

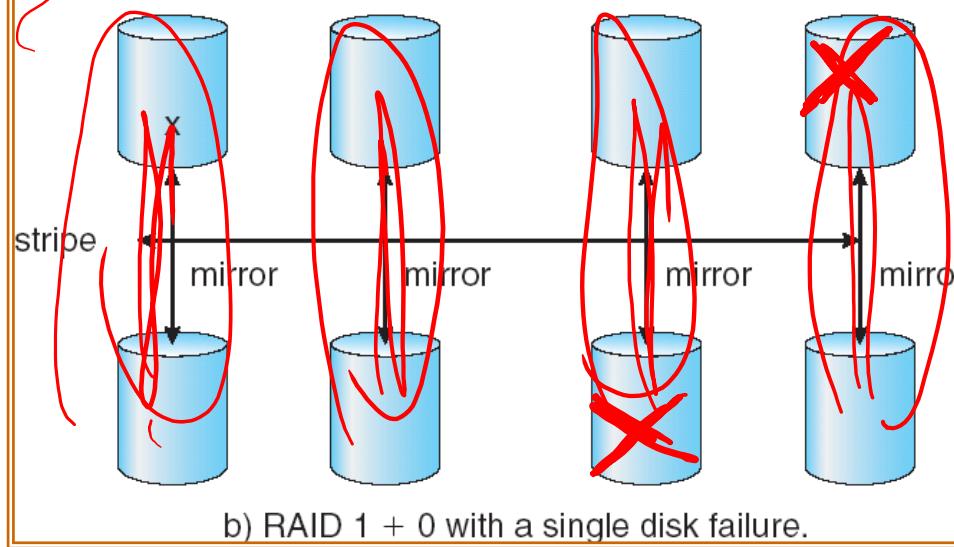
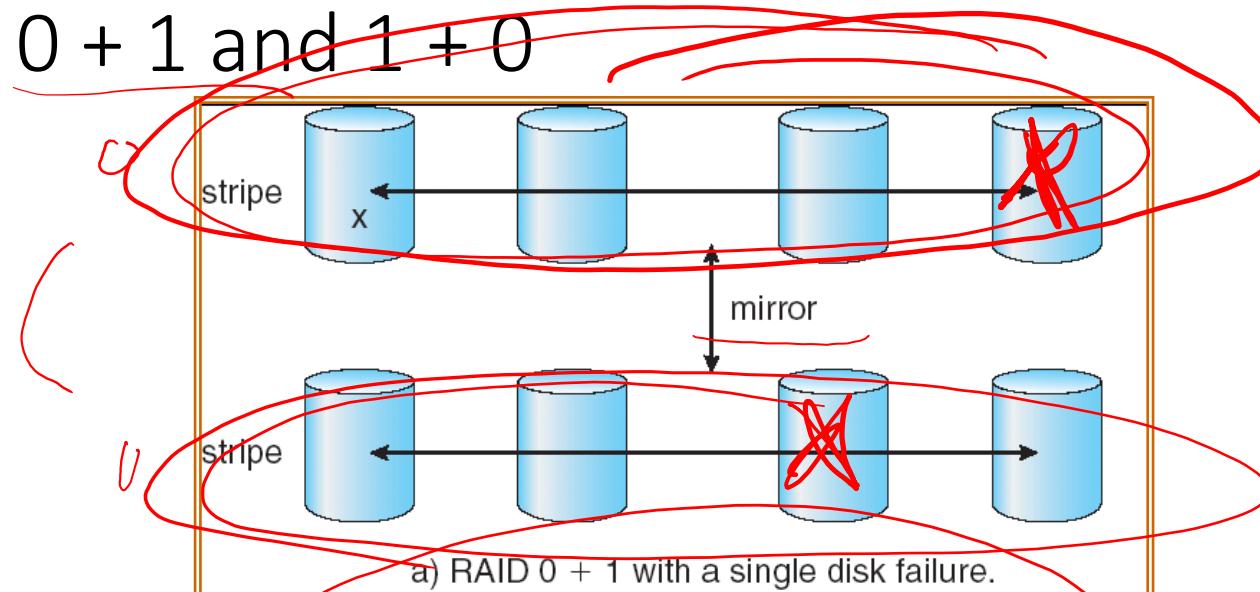
- Choosing 4 blocks out of $\{1,2,3,4,p,q\}$ sufficiently reconstruct $\{1,2,3,4\}$
- Reed-Solomon code or EVENODD parity
- Need hardware support

RAID-5 Reliability



- Let the probability of 1-year up of a disk be p
- The probability of 1-year up of a RAID-0 of 4 disks:
 - p^4 ; ~ 0.96 if $p=0.99$
- The probability of 1-year up of a RAID-5 of 5 disks:
 - $p^5 - (5,1)(1-p)*p^4$; ~ 0.999 if $p=0.99$

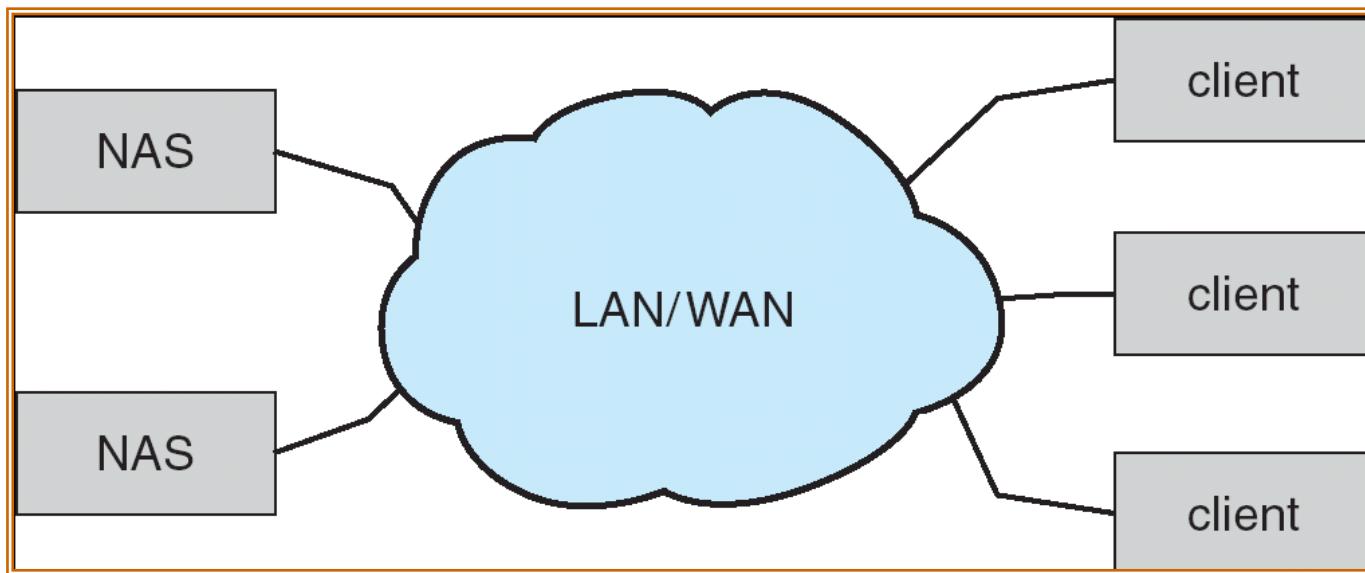
RAID 0 + 1 and 1 + 0



← Better survivability

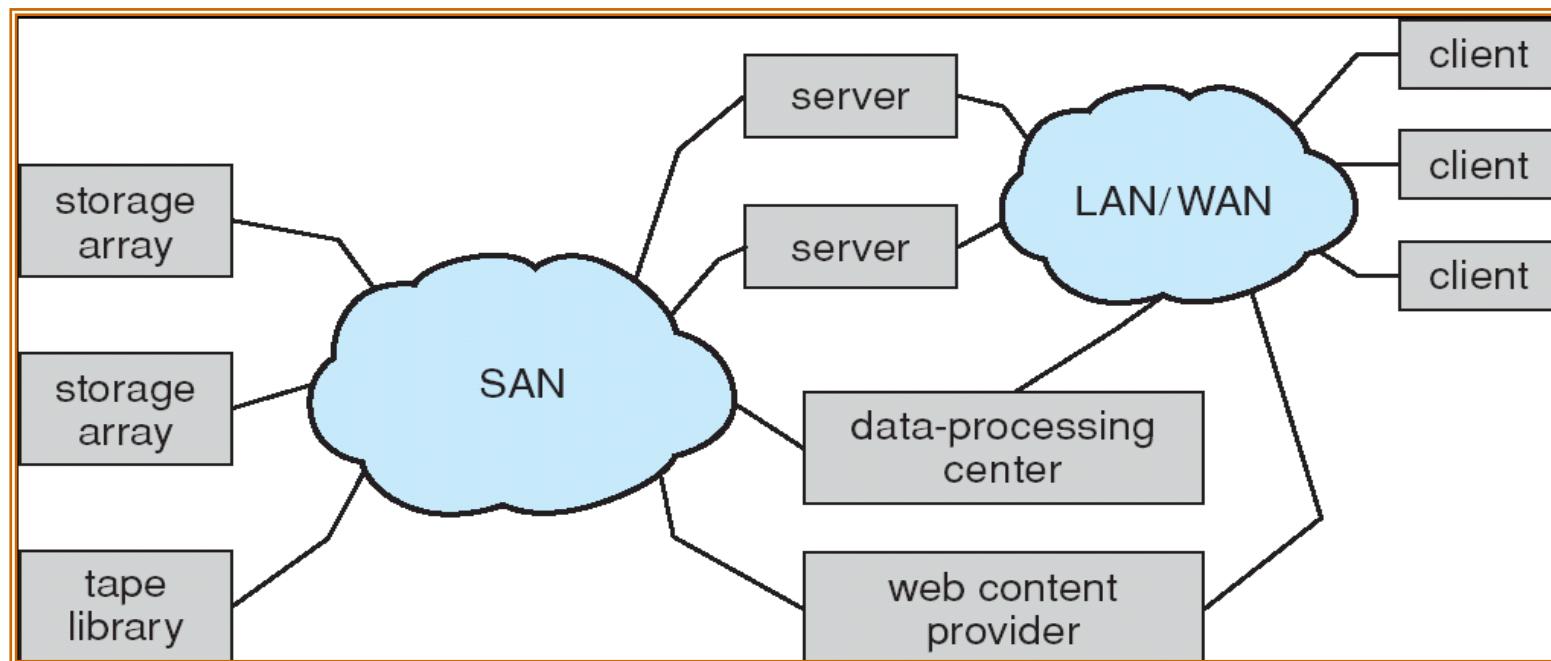
Network-Attached Storage

- Distributed storage operating over **legacy** network
- NFS, CIFS, SAMBA are common file-level protocols
- iSCSI is block-level protocol



Storage-Area Network

- Common in large storage environments
- Involving dedicated network infrastructures
- Multiple hosts attached to multiple storage arrays – flexible



SAN vs. NAS

- NAS may operate over legacy network
 - Interoperability is much more important
- SAN is a network dedicated for storage
 - Performance and resource allocation is the primary concern

Solid-State Disks (SSDs)

- Storage devices that emulate standard block devices using non-volatile memory

- Flash memory or battery-backed RAM
 - The OS use the legacy I/O stack on top of SSDs

- Products

- Embedded flash cards, SD cards, USB thumb drives, SSDs, PCI-e flash cards

- Performance

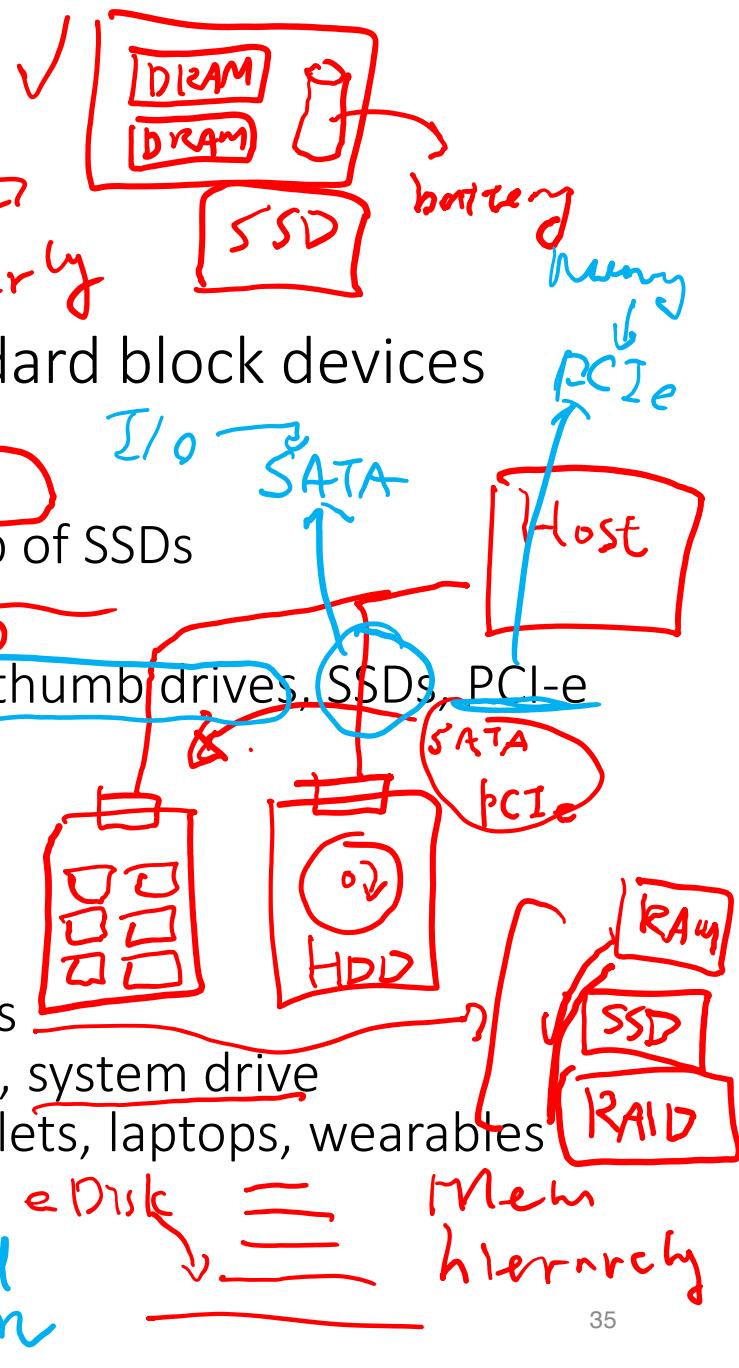
- RAM SSD > flash SSD >> HDD

- Applications

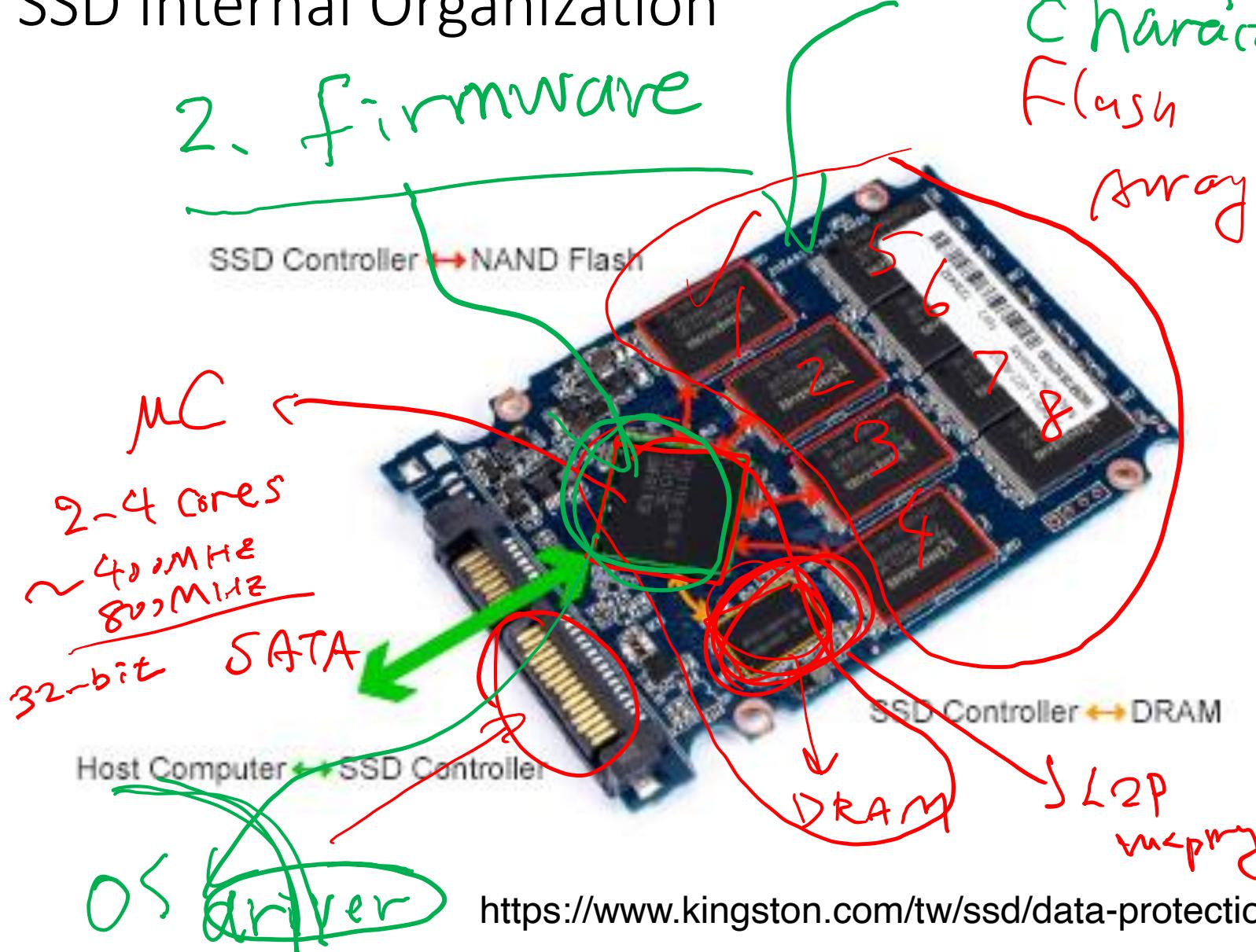
- Cloud storage: tier storage, cache SSDs
- Personal computer: HDD replacement, system drive
- Embedded storage: Smartphones, tablets, laptops, wearables

QLC [eMMC 192GB]
UFS 4K

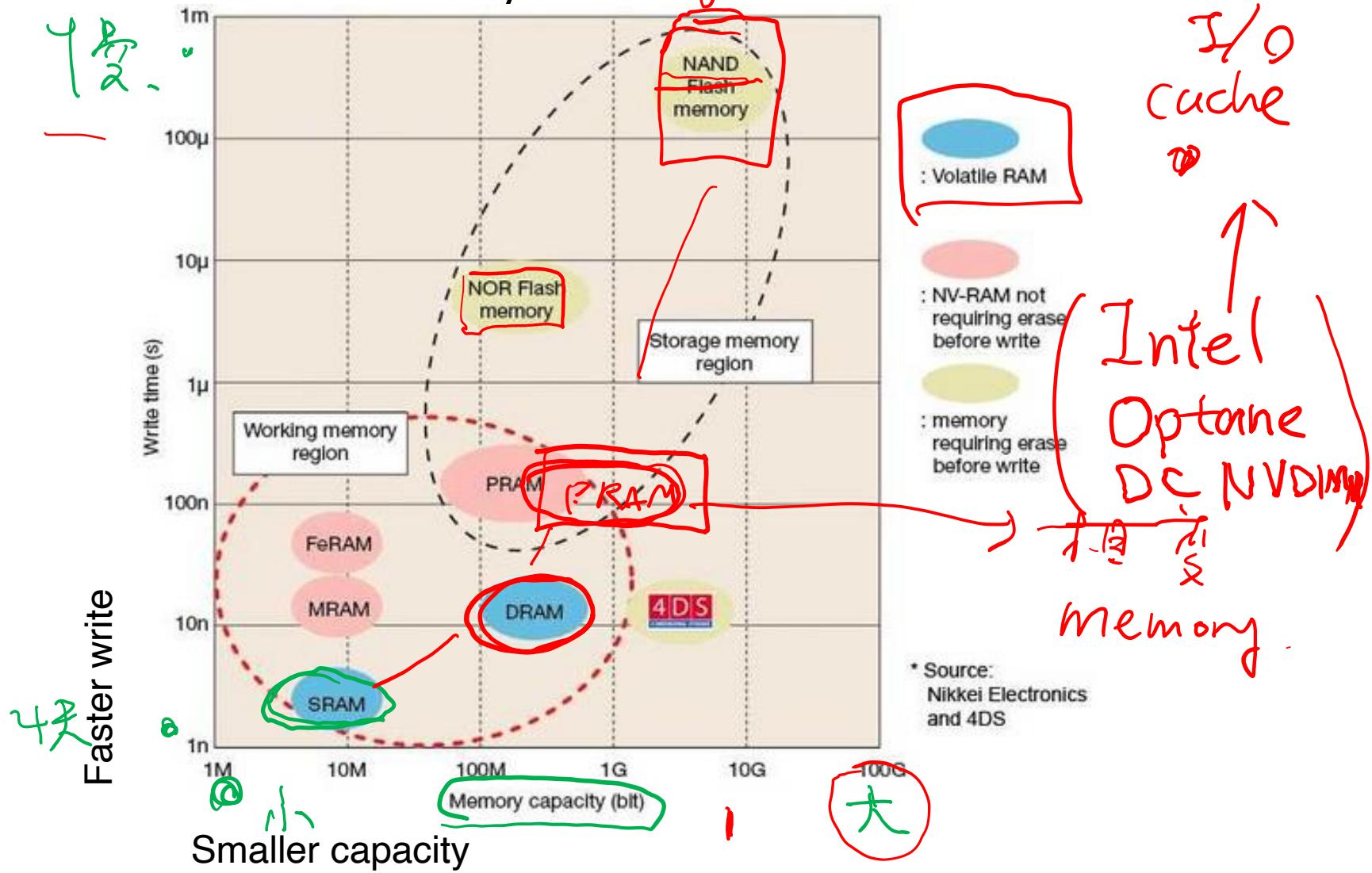
⇒ Android
4Kn



SSD Internal Organization



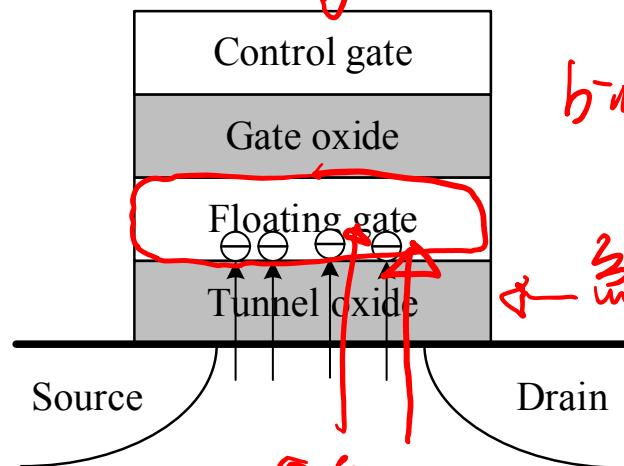
Non-Volatile Memory



Flash Memory

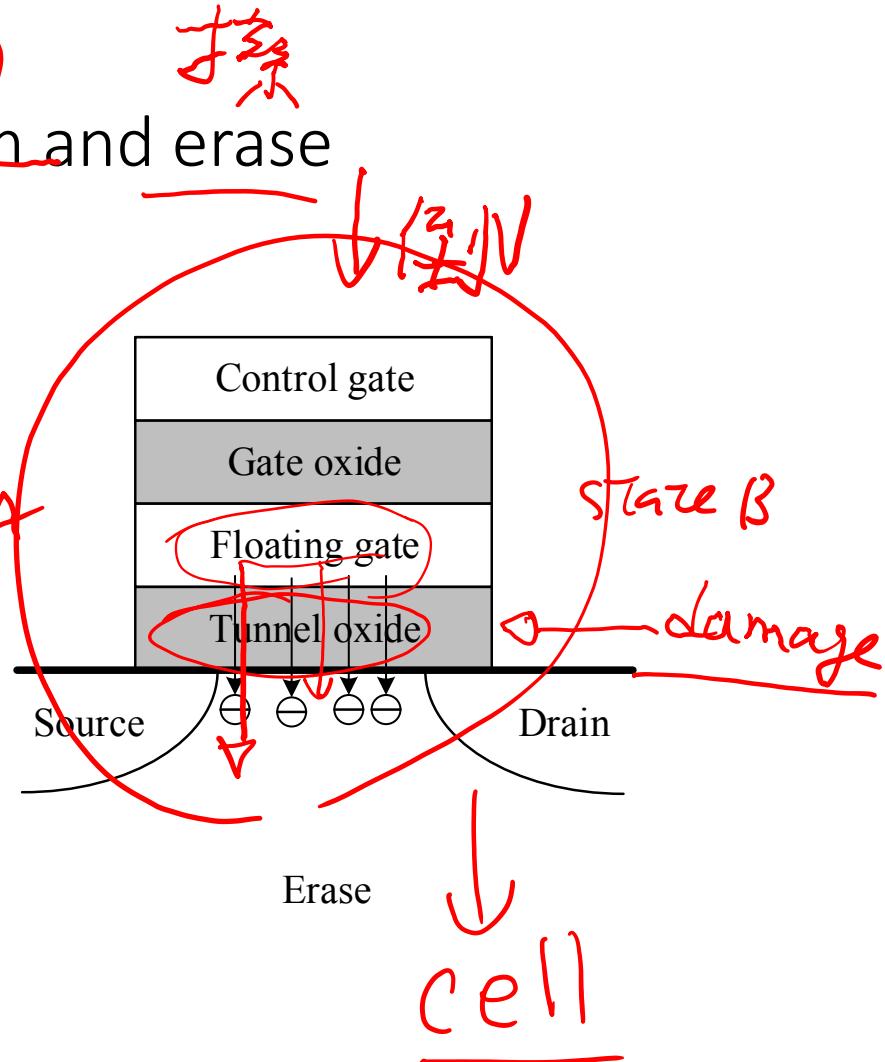
- Cell structure, flash program and erase

- P/E cycle endurance

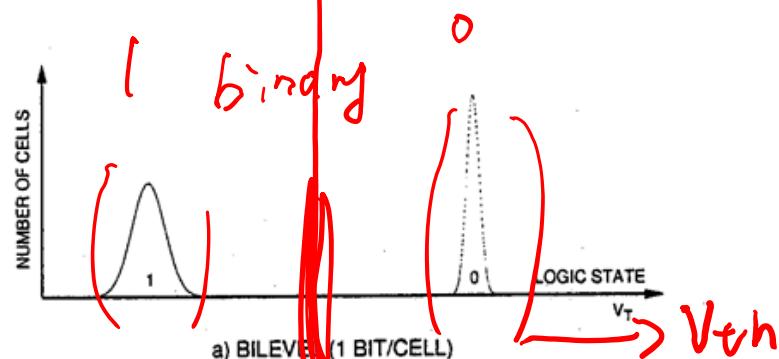


Program (write)

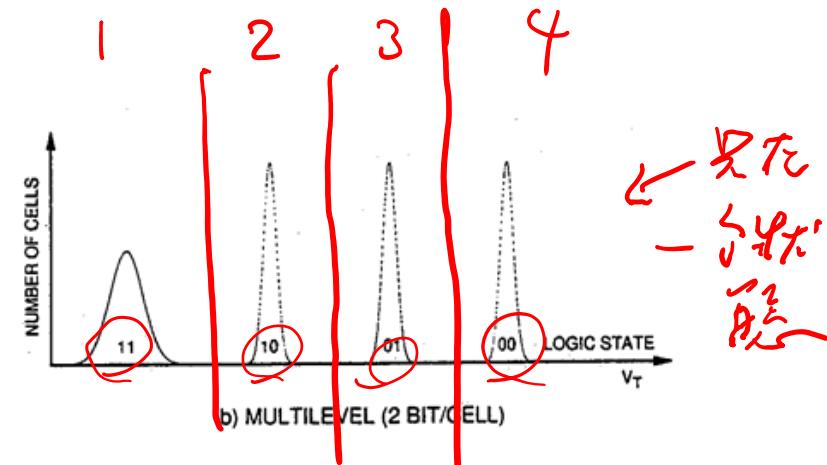
$FG \Rightarrow$ 存储单元



Multilevel Cells



a) BILEVEL (1 BIT/CELL)



b) MULTILEVEL (2 BIT/CELL)

SLC

Single-level cell

單用

• SLC vs. MLC flash

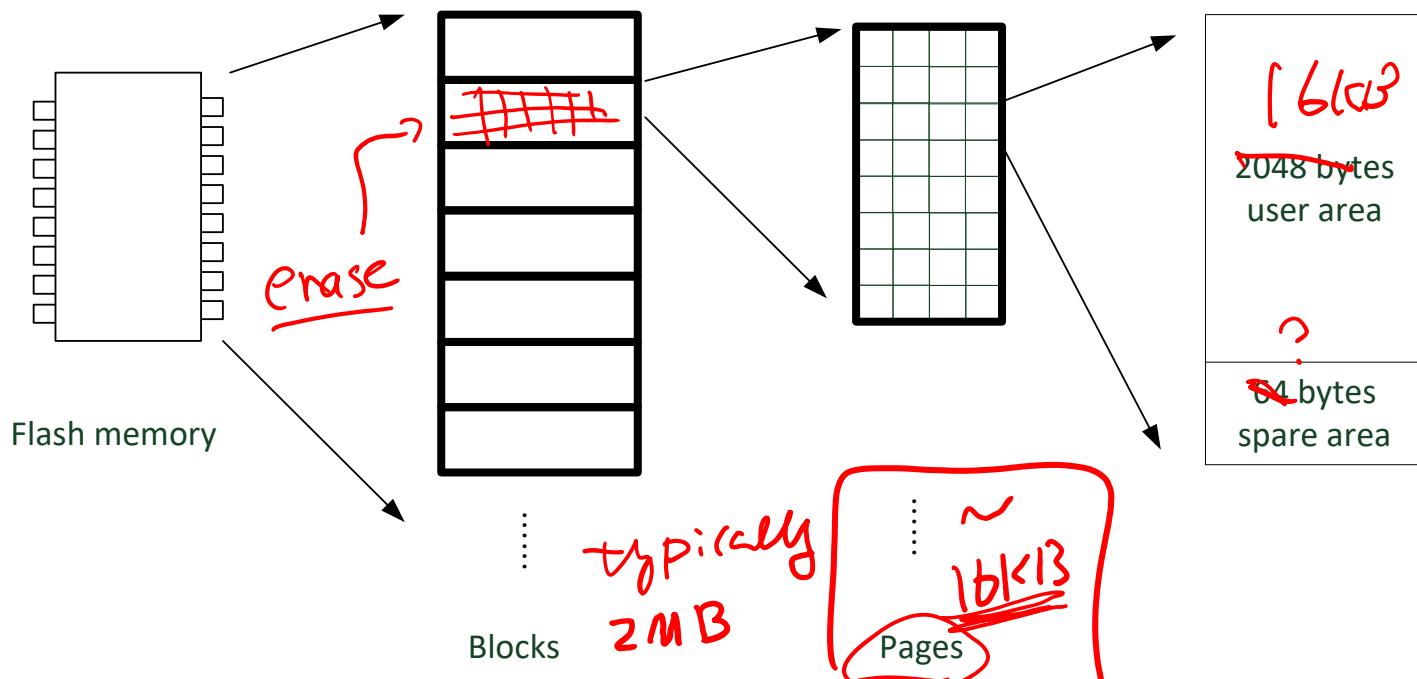
- Read performance: comparable
- Write performance: SLC is 2x~3x faster
- P/E endurance: 100K (SLC), 3K cycles (MLC)
- Cost: SLC is twice or thrice more expensive

• TLC and QLC are currently in mass production

Individual 3bits.

4 bits.
PLC
5 bits

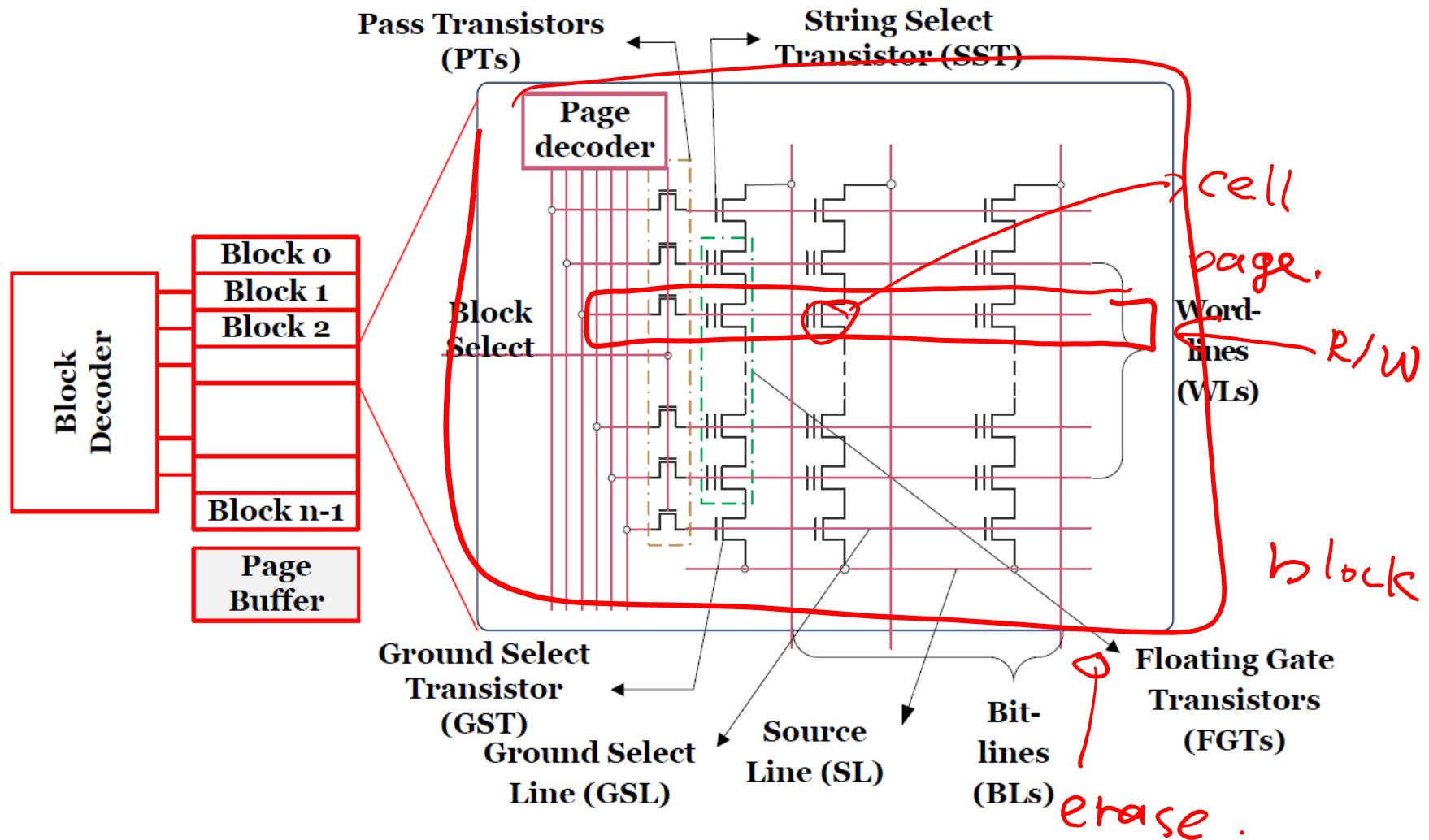
NAND Flash Geometry



- Units of operation
 - Read/write: pages
 - Erase: blocks

erase

Read
write



Flash Translation Layer (FTL)



- A firmware layer inside of SSDs
 - Hiding flash memory physics from the host
- Provide block device emulation to the host
- Manage flash memory inside of SSDs
 - Logical-to-physical address translation ↗^{2P}
 - Garbage collection
 - Wear leveling

@ 11:12 am

Logical-to-Physical Address Translation

why?

- Pages cannot be overwritten unless being erased



block
(τ)

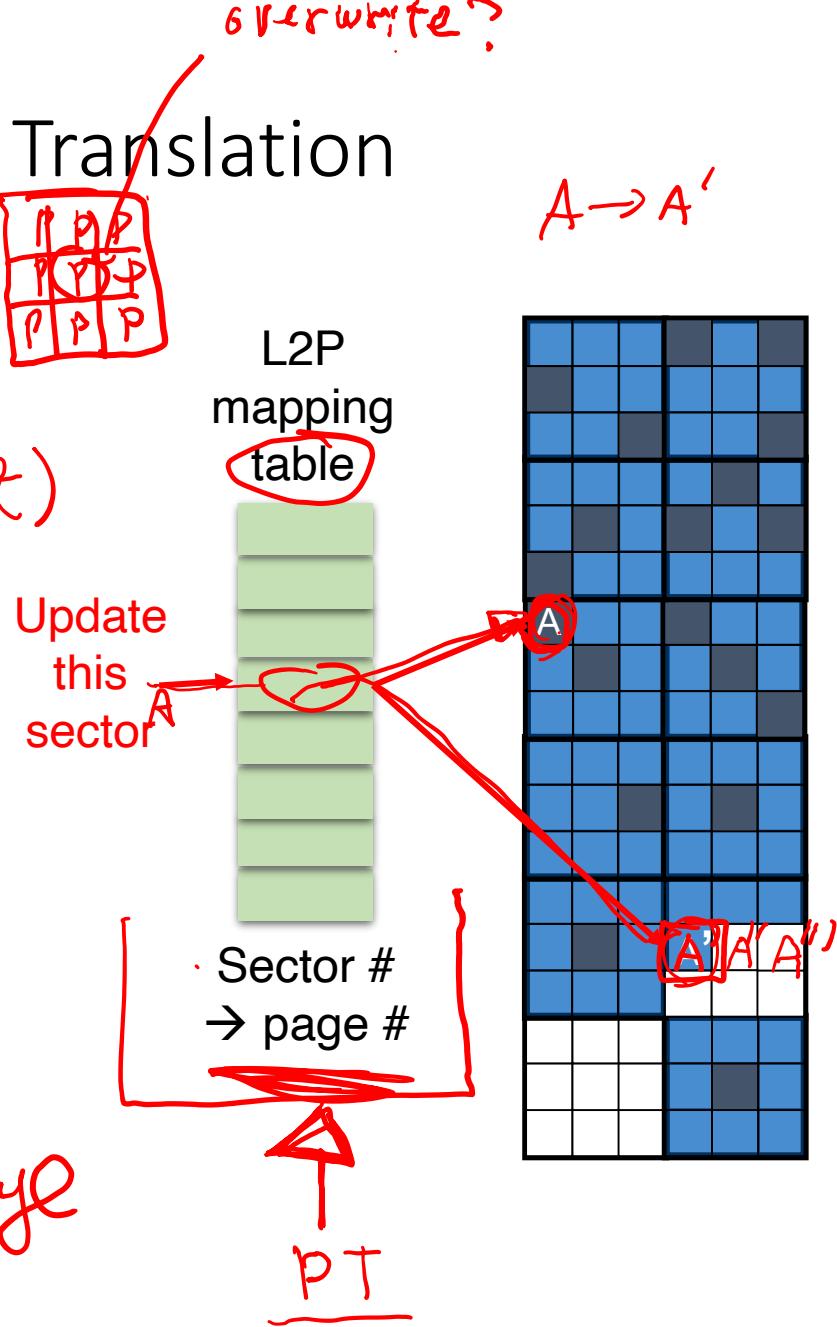
- Erase a block every time a page is overwritten
 - Too inefficient

- Out-of-place update; mark old data invalid

LFS

- Need logical-to-physical address translation
 - From logical sector # to physical page #

garbage



$A \rightarrow A'$

overwrite?

Garbage Collection

~~Copy → Erase~~



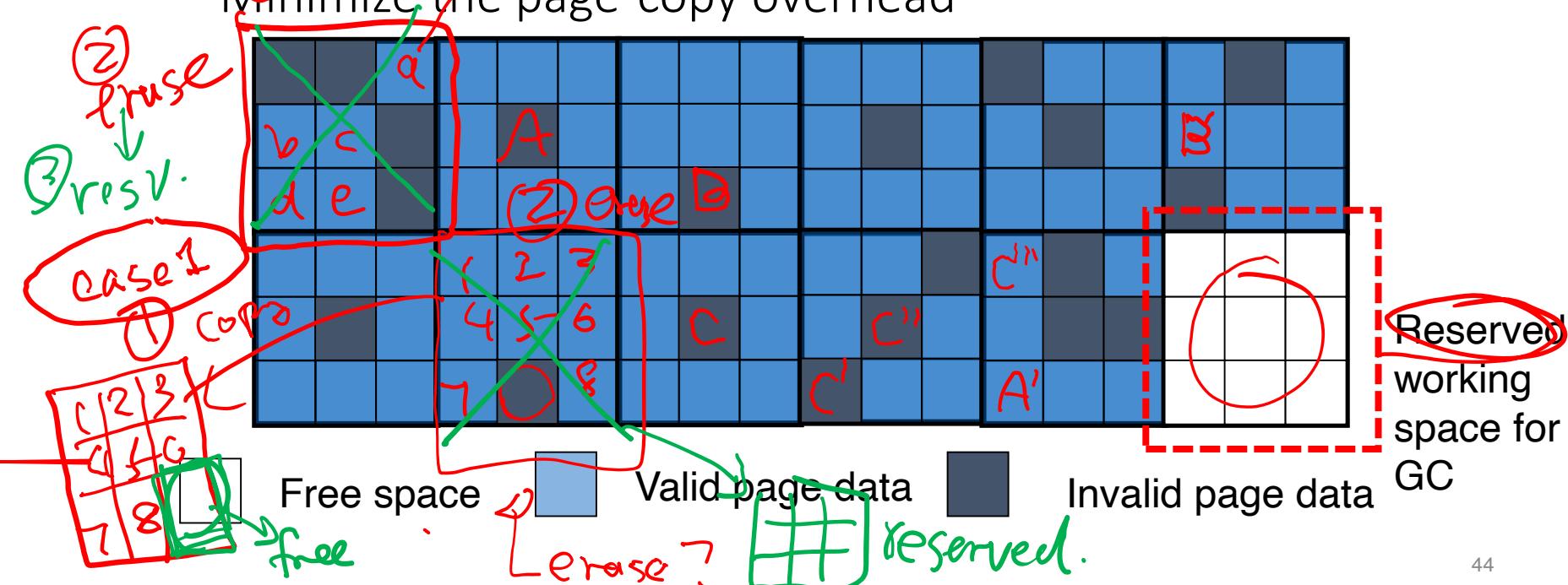
LFS

Compaction

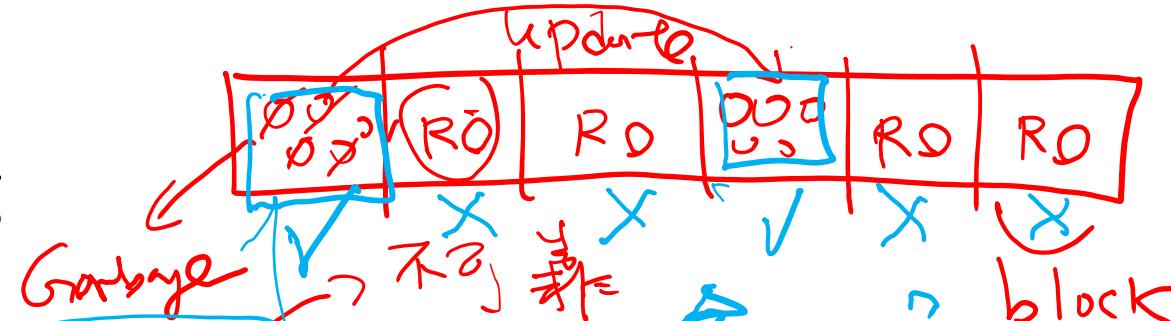
- Recycle memory space occupied by invalid data through block erasure
 - Victim selection
 - Minimize the page-copy overhead

~~case 2~~ → Simple: Greedy

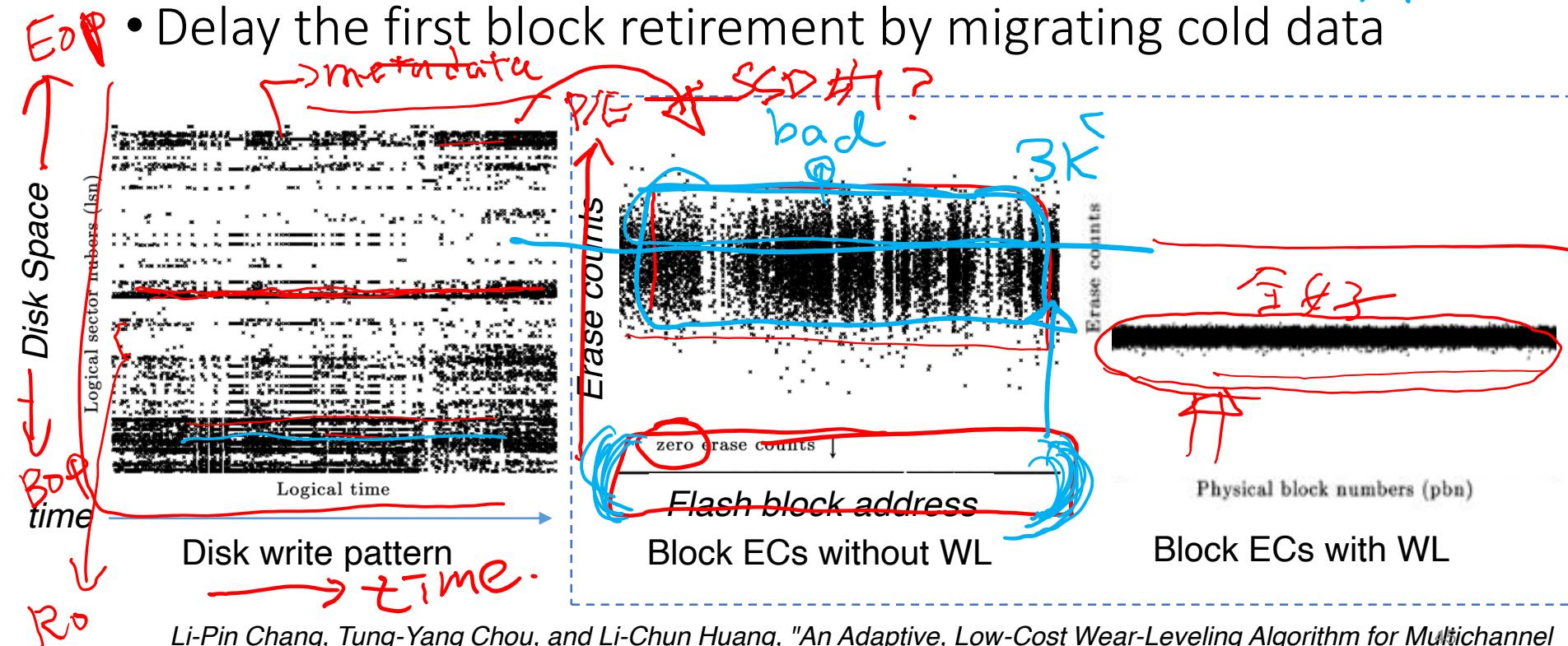
Simple: Greedy



Wear Leveling



- Typically a (MLC) block endures 3000 cycles of program-erase operations (P/E cycles)
 - Locality of write creates frequently written blocks
 - Delay the first block retirement by migrating cold data



End of Chapter 12