

# Lecture 11

# Storage

Prof. Yinqian Zhang

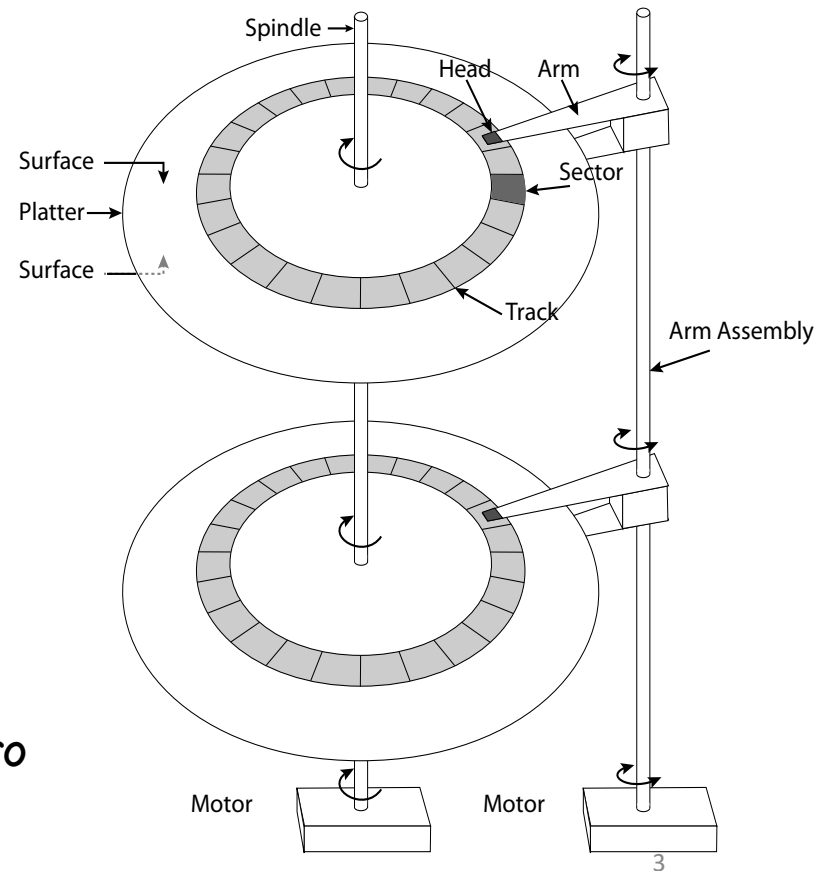
Spring 2023

# Storage Devices

- Magnetic disks
  - Storage that rarely becomes corrupted
  - Large capacity at low cost
  - Block level random access
  - Poor performance for random access
  - Better performance for sequential access
- Flash memory
  - Storage that rarely becomes corrupted
  - Capacity at intermediate cost (5–20x disk)
  - Block level random access
  - Good performance for reads; worse for random writes

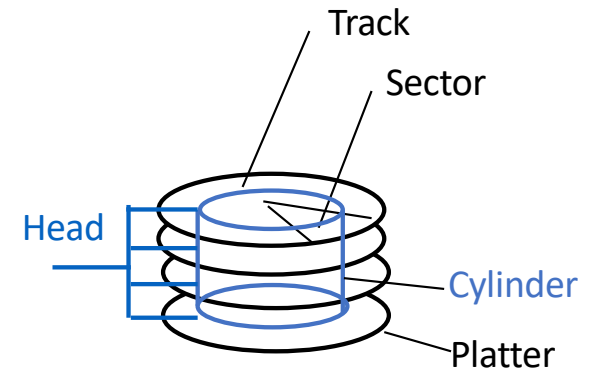
# The Amazing Magnetic Disk

- Unit of Transfer: Sector
  - Ring of sectors form a track
  - Stack of tracks form a cylinder
  - Heads position on cylinders
- Disk Tracks  $\sim 1\mu\text{m}$  (micron) wide
  - Wavelength of light is  $\sim 0.5\mu\text{m}$
  - 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)

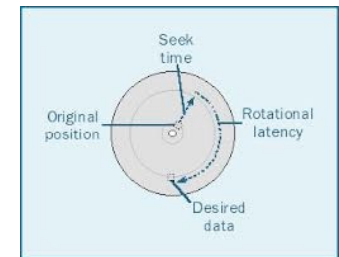


# Magnetic Disks

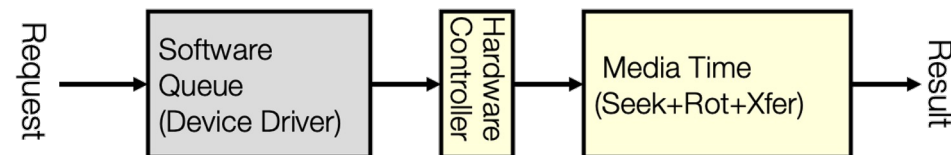
- Cylinders: all the tracks under the head at a given point on all surface
- Read/write data is a three-stage process:
  - Seek time: position the head/arm over the proper track
  - Rotational latency: wait for desired sector to rotate under r/w head
  - Transfer time: transfer a block of bits (sector)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



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



Disk Latency = Queuing Time + Controller time +  
Seek Time + Rotation Time + Transfer Time



# Disk Performance Example

- Assumptions:
  - Ignoring queuing and controller times for now
  - Avg seek time of 5ms,
  - 7200RPM  $\Rightarrow$  Time for rotation:  $60000 \text{ (ms/minute)} / 7200 \text{ (rev/min)} \cong 8\text{ms}$
  - Transfer rate 4MByte/s, sector size 1 Kbyte  $\Rightarrow 1024 \text{ bytes} / 4 \times 10^6 \text{ (bytes/s)} = 256 \times 10^{-6} \text{ sec} \cong .26 \text{ ms}$
- Read sector from random place on disk:
  - Seek (5ms) + Rot. Delay (4ms) + Transfer (0.26ms)
  - *Approx* 10ms to fetch/put data: 100 KByte/sec
- Read sector from random place in the same cylinder:
  - Rot. Delay (4ms) + Transfer (0.26ms)
  - *Approx* 5ms to fetch/put data: 200 KByte/sec
- Read next sector on same track:
  - Transfer (0.26ms): 4 MByte/sec

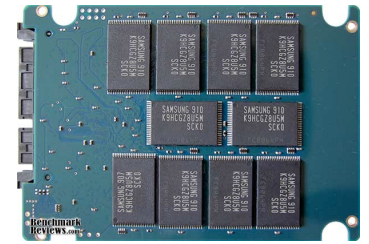
**Key to using disk effectively (especially for file systems) is to *minimize seek and rotational delays***

# Typical Numbers for Magnetic Disk

| Parameter                  | Info / Range   |
|----------------------------|--|
| Space/Density              | Space: 8TB (Seagate), 10TB (Hitachi) in 3½ inch form factor!<br>Areal Density: $\geq 1$ Terabit/square inch! (SMR, Helium, ...)                                    |
| Average seek time          | Typically 5-10 milliseconds.<br>Depending on reference locality, actual cost may be 25-33% of this number.   |
| Average rotational latency | Most laptop/desktop disks rotate at 3600-7200 RPM (16-8 ms/rotation). Server disks up to 15,000 RPM.<br>Average latency is halfway around disk so 8-4 milliseconds |
| Controller time            | Depends on controller hardware   |
| Transfer rate              | Typically 50 to 100 MB/s.  |
| Cost                       | Used to drop by a factor of two every 1.5 years (or even faster); now slowing down   |

# Solid State Disks (SSDs)

- 1995 – Replace rotating magnetic media with non-volatile memory
- 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
  - Data erased at the block level
- 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!



# HDD vs SSD Comparison

SSD prices drop much faster than HDD

Price Crossover Point for HDD and SSD

|          | 2012 | 2013 | 2014 | 2015E | 2016F | 2017F |
|----------|------|------|------|-------|-------|-------|
| HDD      | 0.09 | 0.08 | 0.07 | 0.06  | 0.06  | 0.06  |
| 2.5" SSD | 0.99 | 0.68 | 0.55 | 0.39  | 0.24  | 0.17  |



## 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**

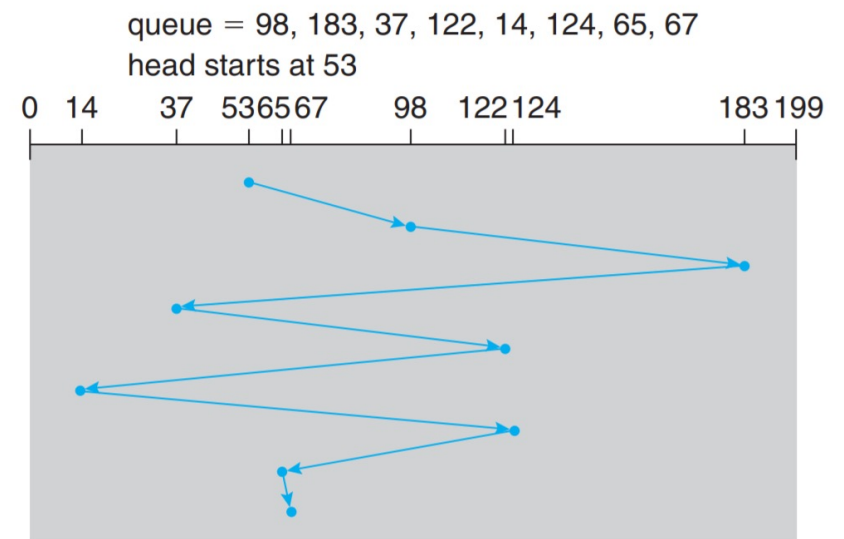


# Disk Scheduling

- There are many sources of disk I/O request
  - OS
  - System processes
  - User processes
- OS should think how to use hardware efficiently?
  - Disk bandwidth
  - Access time
- Given a sequence of access cylinders in the HDD
  - 98, 183, 37, 122, 14, 124, 65, 67
  - Head point: 53
  - Pages: 0 ~ 199
- Minimize seek time
  - Seek time  $\approx$  seek distance
- How to minimize the total head movement distance?
  - Minimize the total number of cylinders.

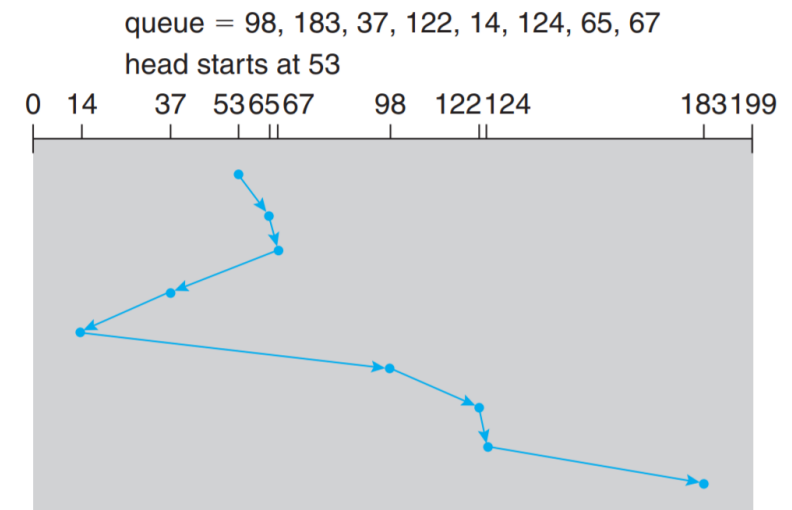
# Disk Scheduling: FIFO

- FIFO Order
  - Fair among requesters, but order of arrival may be to random spots on the disk  $\Rightarrow$  Very long seeks
- The head movement distance = ?



# Disk Scheduling: SSTF

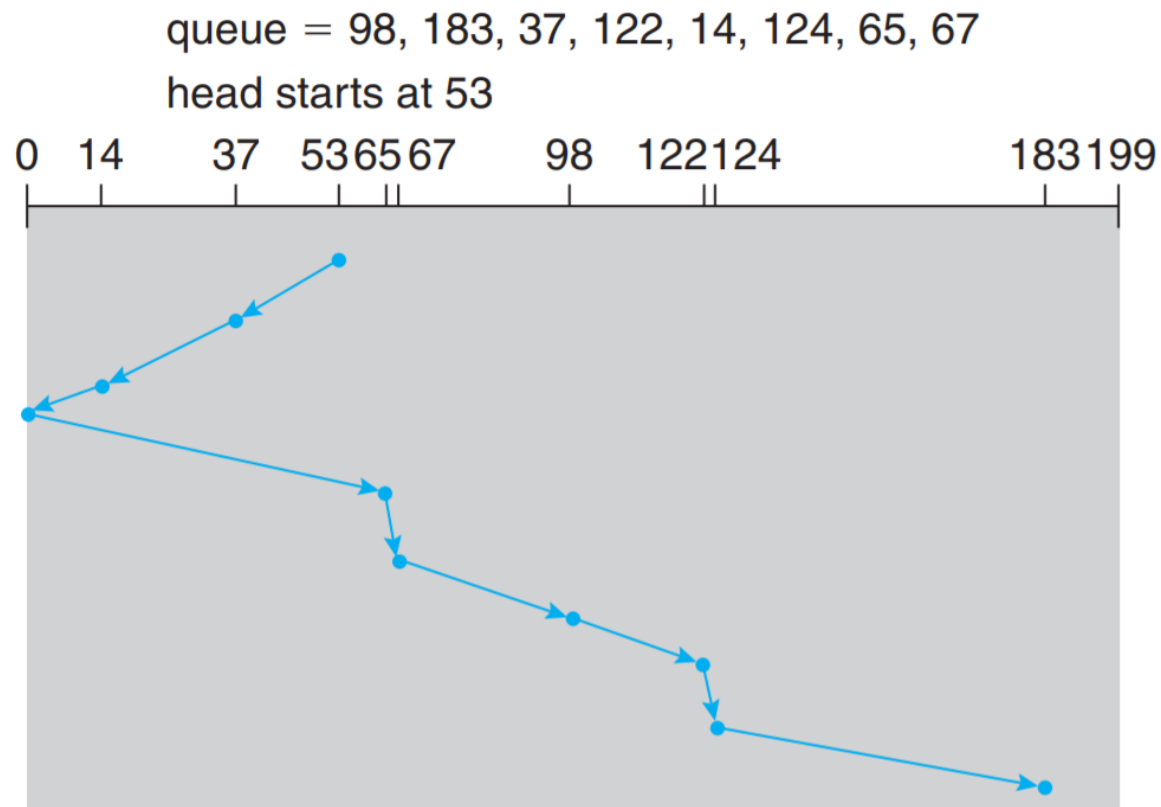
- Shortest Seek Time First order
  - Shortest Seek Time First selects the request with the minimum seek time from the current head position
  - SSTF scheduling is a form of **SJF** scheduling; may cause **starvation** of some requests
- The head movement distance = ?



# Disk Scheduling: SCAN

- SCAN order
  - SCAN algorithm a.k.a., elevator algorithm
  - 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.
  - But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest
- The head movement distance = ?

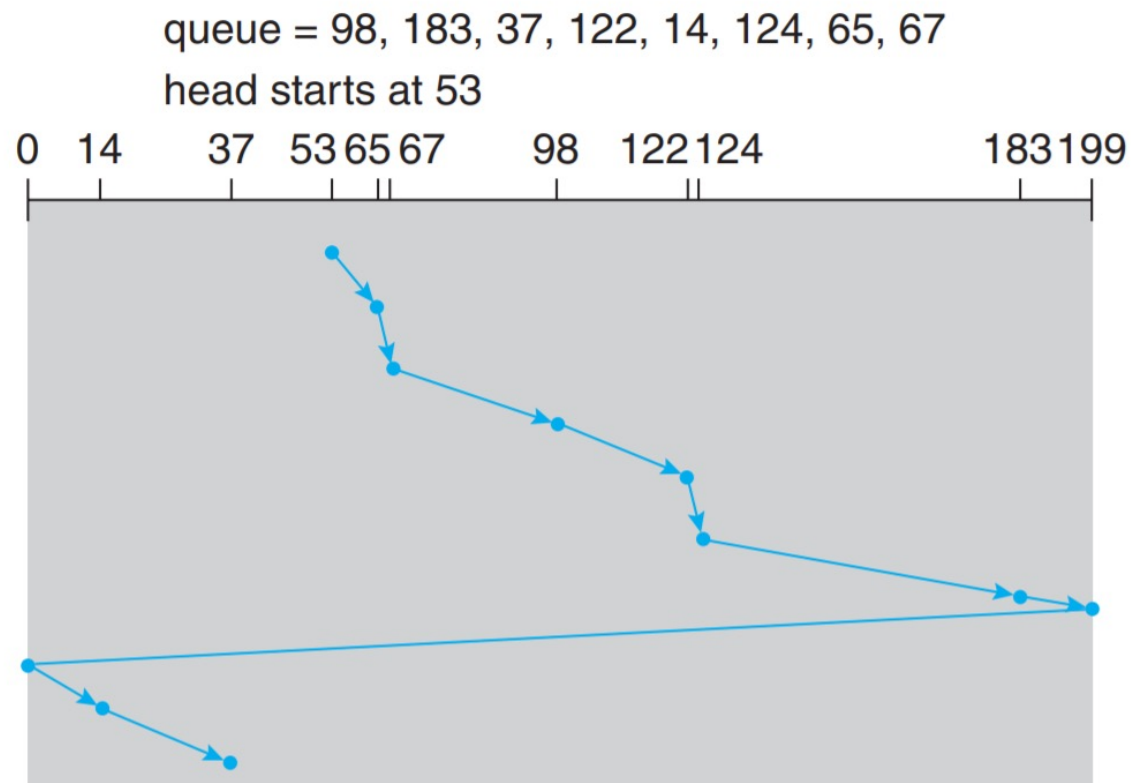
# Disk Scheduling: SCAN



# Disk Scheduling: C-SCAN

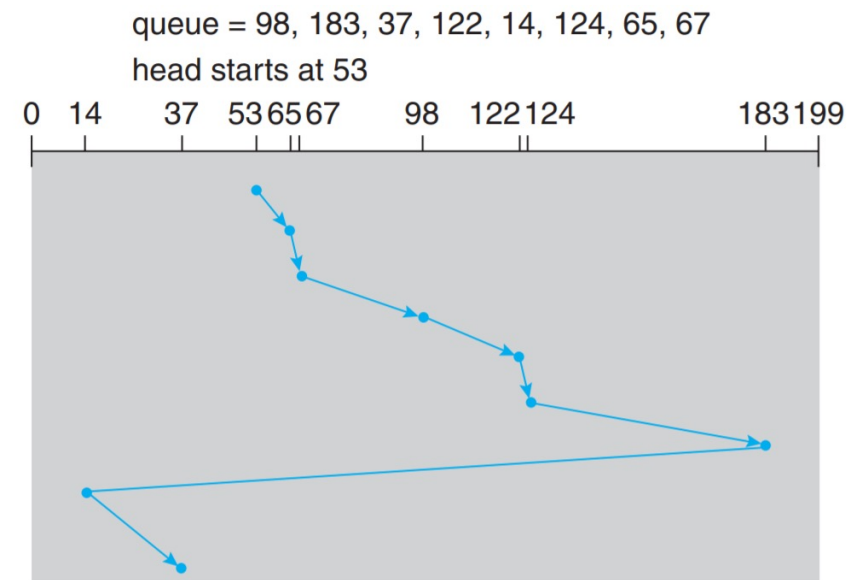
- C-SCAN order
  - Provides a more uniform wait time than SCAN
  - The head moves from one end of the disk to the other, servicing requests as it goes
    - When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip
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

# Disk Scheduling: C-SCAN



# Disk Scheduling: LOOK, C-LOOK

- Look and C-LOOK order
  - LOOK a variant of SCAN, C-LOOK a variant 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: the head movement distance = ?





# Selecting a Disk-Scheduling Algorithm

- SSTF is common and has a natural appeal
- SCAN and C-SCAN perform better for systems that place a heavy load on the disk
  - Less starvation
- Either SSTF or LOOK is a reasonable choice for the default algorithm
- Performance depends on the number and types of requests
- The disk-scheduling algorithm should be written as a separate module of the operating system, allowing it to be replaced with a different algorithm if necessary

**Thank you!**

