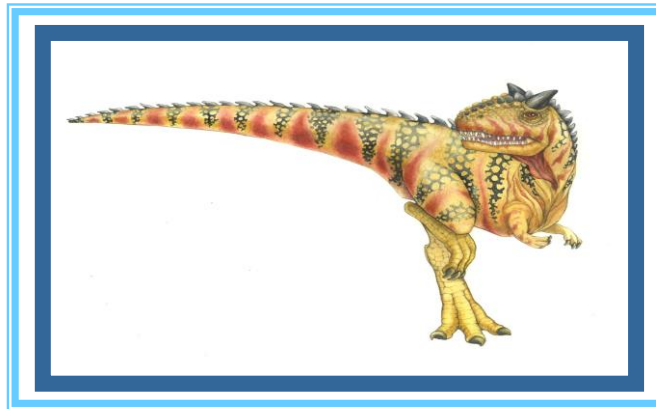# Chapter 10: Mass-Storage Systems

# Chapter 10:  Mass-Storage Systems

- Overview of Mass Storage Structure

- Disk Structure

- Disk Scheduling

- Disk Management

- Swap-Space Management

- RAID Structure

- Stable-Storage Implementation

# Objectives

- To describe the physical structure of secondary storage devices and its effects on the uses of the devices

- To explain the performance characteristics of mass-storage devices

- To evaluate disk scheduling algorithms

- To discuss operating-system services provided for mass storage, including RAID

# Storage Structure Hierarchy

Register

↓

Cache

↓

Main Memory ← *volatile*

↓

Electronic Disk

↓

Magnetic Disk

↓

Optical Disk

↓

Magnetic Tape

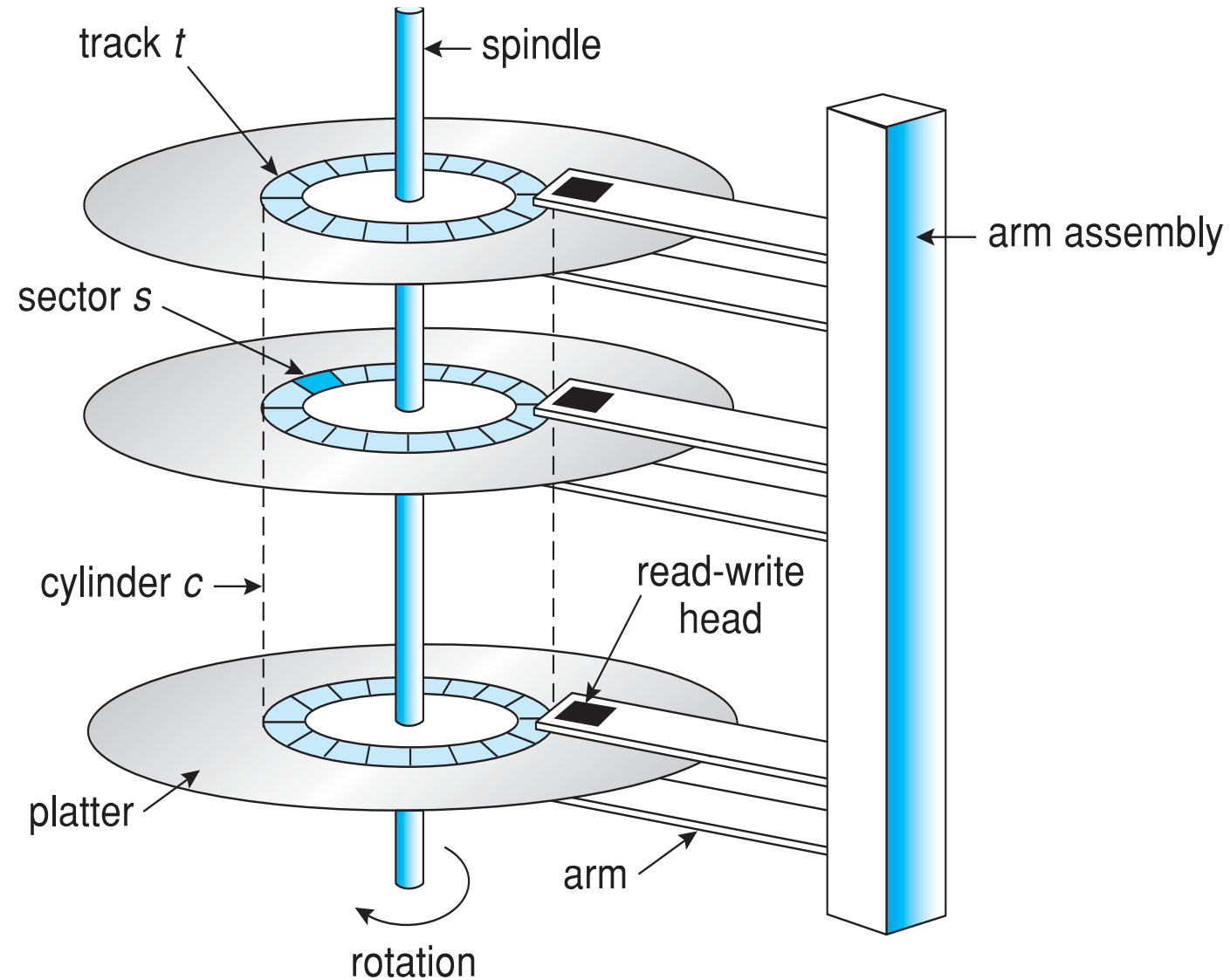} *Mass storage / Permanent storage*

Speed

Size
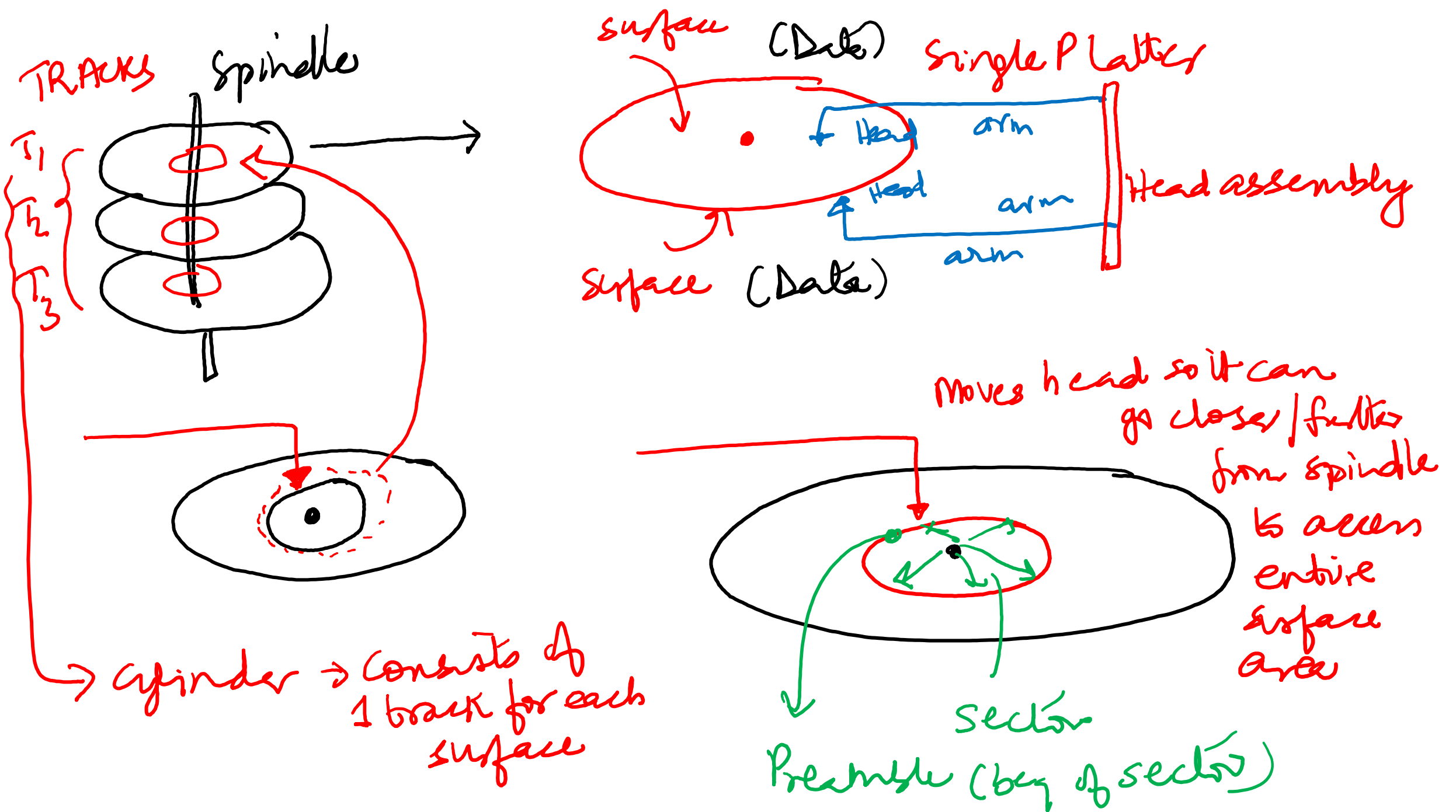
# Overview of Mass Storage Structure

- **Magnetic disks** provide bulk of secondary storage of modern computers
  - Drives rotate at 60 to 250 times per second
  - **Transfer rate** is rate at which data flow between drive and computer
  - **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**)
  - **Head crash** results from disk head making contact with the disk surface -- That's bad
- Disks can be removable
- Drive attached to computer via **I/O bus**
  - Busses vary, including **EIDE**, **ATA**, **SATA**, **USB**, **Fibre Channel**, **SCSI, SAS, Firewire**
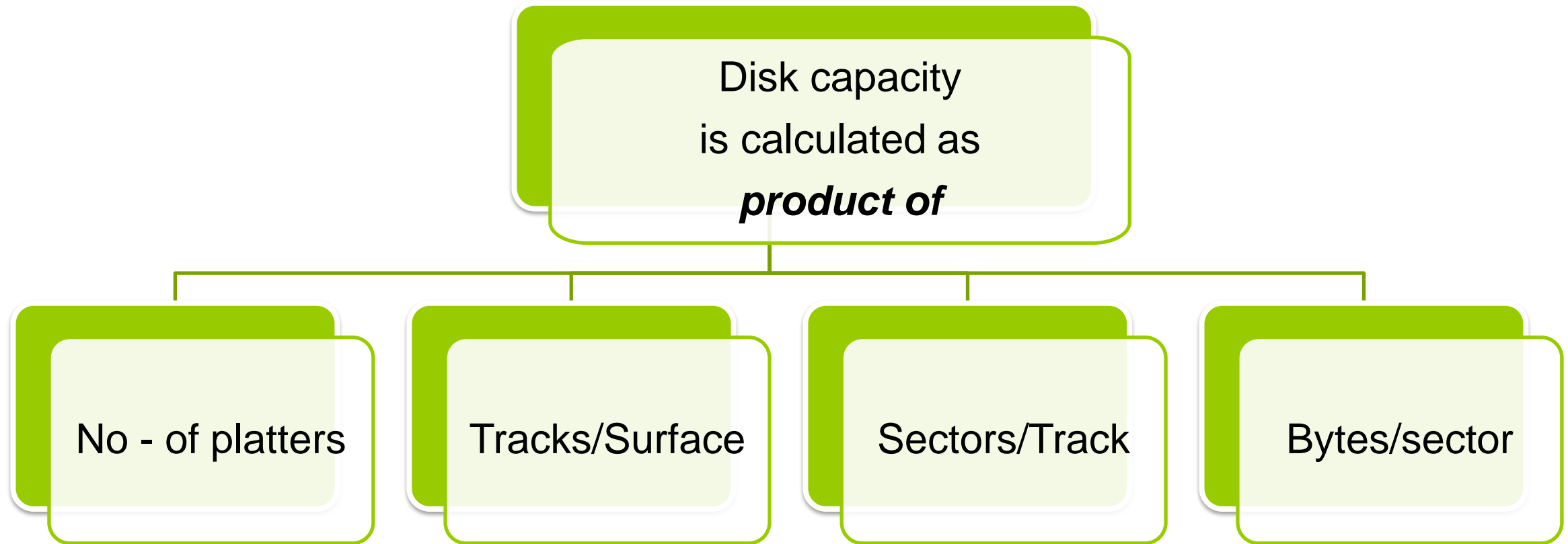  - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array

# Moving-head Disk Mechanism

TRACKS  Spindle

$T_1$
$T_2$
$T_3$

surface (Data)    Single Platter

Head   arm

Head   arm          Head assembly

arm

Surface (Data)

→ Cylinder → Consists of 1 track for each surface

Moves head so it can go closer/further from spindle to access entire surface area

sector

Preamble (beg of sector)

Disk capacity
is calculated as
*product of*

| No - of platters | Tracks/Surface | Sectors/Track | Bytes/sector |

# Hard Disks

| Spindle [rpm] | Average latency [ms] |
|---|---|
| 4200 | 7.14 |
| 5400 | 5.56 |
| 7200 | 4.17 |
| 10000 | 3 |
| 15000 | 2 |

- Platters range from .85" to 14" (historically)
  - Commonly 3.5", 2.5", and 1.8"
- Range from 30GB to 3TB per drive
- Performance

(From Wikipedia)

  - Transfer Rate – theoretical – 6 Gb/sec
  - Effective Transfer Rate – real – 1Gb/sec
  - Seek time from 3ms to 12ms – 9ms common for desktop drives
  - Average seek time measured or calculated based on 1/3 of tracks
  - Latency based on spindle speed
    - $1 / (RPM / 60) = 60 / RPM$
  - Average latency = ½ latency

# Hard Disk Performance

☐ **Access Latency** = **Average access time** = average seek time + average latency

  ☐ For fastest disk 3ms + 2ms = 5ms

  ☐ For slow disk 9ms + 5.56ms = 14.56ms

☐ Average I/O time = average access time + (amount to transfer / transfer rate) + controller overhead

☐ For example to transfer a 4KB block on a 7200 RPM disk with a 5ms average seek time, 1Gb/sec transfer rate with a .1ms controller overhead =

  ☐ 5ms + 4.17ms + 0.1ms + transfer time =

  ☐ Transfer time = 4KB / 1Gb/s * 8Gb / GB * 1GB / $1024^2$KB = 32 / $(1024^2)$ = 0.031 ms

  ☐ Average I/O time for 4KB block = 9.27ms + .031ms = 9.301ms

# The First Commercial Disk Drive



1956
IBM RAMDAC computer
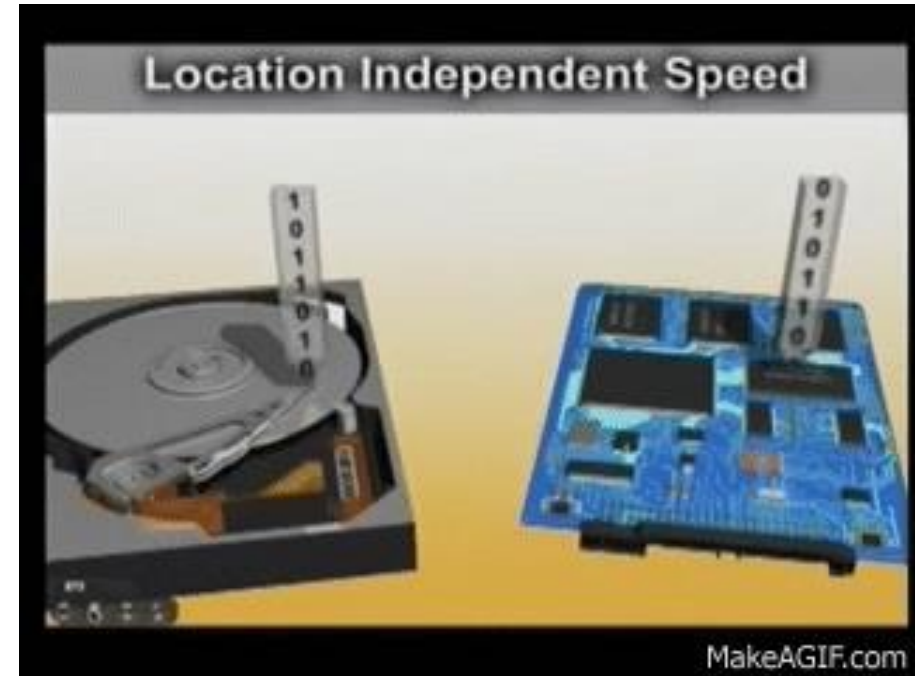included the IBM Model 350
disk storage system

5M (7 bit) characters
50 x 24" platters
Access time = < 1 second

# Solid-State Disks

- Nonvolatile memory used like a hard drive
  - Many technology variations
- Can be more reliable than HDDs
- More expensive per MB
- Maybe have shorter life span
- Less capacity
- But much faster
- Busses can be too slow -> connect directly to PCI for example
- No moving parts, so no seek time or rotational latency

# Magnetic Tape

- Was early secondary-storage medium

  - Evolved from open spools to cartridges

- Relatively permanent and holds large quantities of data

- Access time slow

- Random access ~1000 times slower than disk

- Mainly used for backup, storage of infrequently-used data, transfer medium between systems

- Kept in spool and wound or rewound past read-write head

- Once data under head, transfer rates comparable to disk

  - 140MB/sec and greater

- 200GB to 1.5TB typical storage

- Common technologies are LTO-{3,4,5} and T10000

# After the Break!

# 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

- Minimize seek time

- Seek time ≈ seek distance: Defined as **the time required by the read/write head to move from one track to another**.

- Disk **bandwidth** is the total number of bytes transferred, divided by the total time between the first request for service and the completion of the last transfer

# Disk Scheduling (Cont.)

- There are many sources of disk I/O request

  - OS

  - System processes

  - Users processes

- I/O request includes input or output mode, disk address, memory address, number of sectors to transfer

- OS maintains queue of requests, per disk or device

- Idle disk can immediately work on I/O request, busy disk means work must queue

  - Optimization algorithms only make sense when a queue exists

# Disk Scheduling (Cont.)

- Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying "depth")

- Several algorithms exist to schedule the servicing of disk I/O requests

- The analysis is true for one or many platters

- Algorithms:
  - FCFS
  - SSTS
  - SCAN
  - C-SCAN
  - LOOK
  - C-LOOK

# Disk Scheduling (Cont.)

☐ Note that drive controllers have small buffers and can manage a queue of I/O requests (of varying "depth")

☐ Several algorithms exist to schedule the servicing of disk I/O requests

☐ The analysis is true for one or many platters

☐ We illustrate scheduling algorithms with a request queue (0-199)

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

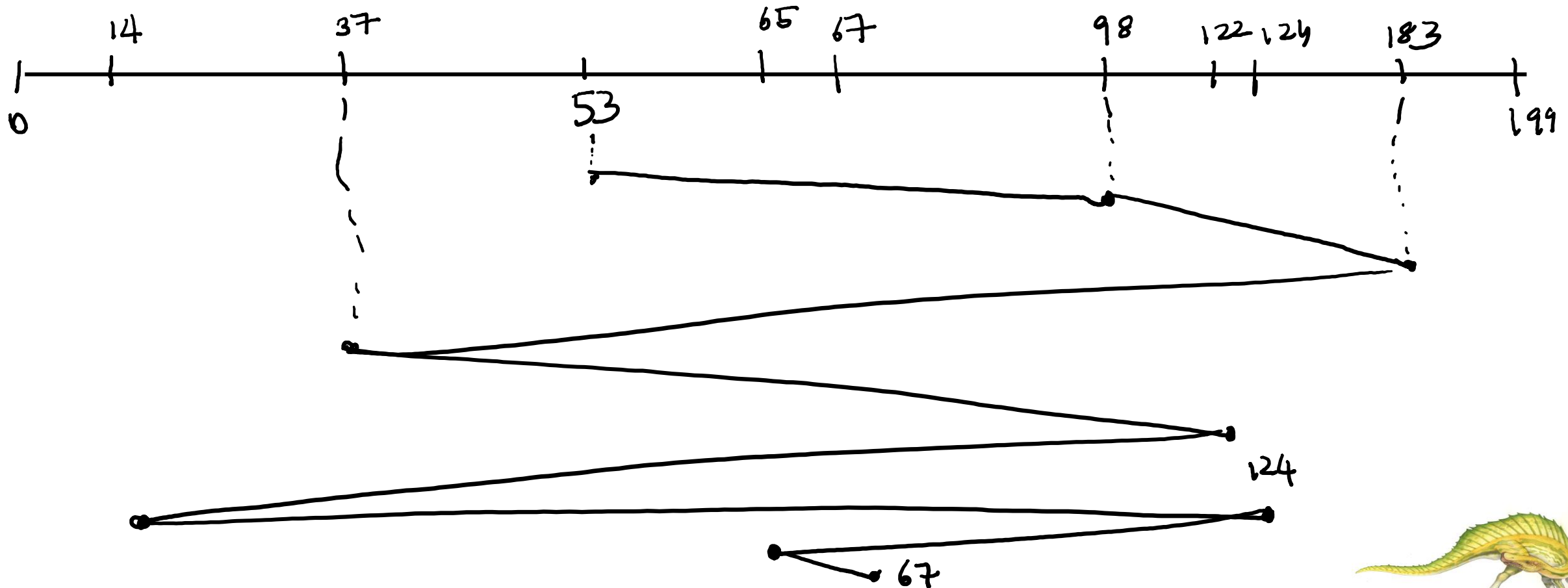Head pointer 53

Total head movement of 640 cylinders

queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

# FCFS

Request queue (0-199)
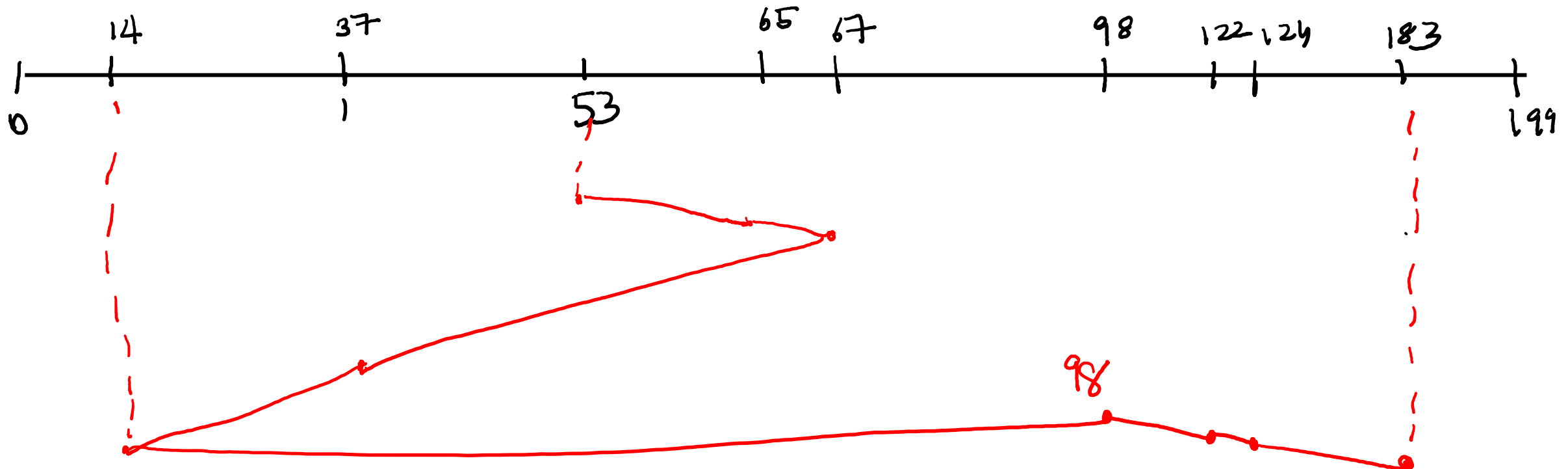98, 183, 37, 122, 14, 124, 65, 67
Head pointer 53

# SSTF Scheduling

☐ Shortest Seek Time First selects the request with the minimum seek time from the current head position

☐ Form of SJF scheduling; may cause starvation of some requests

☐ Illustration shows total head movement of 236 cylinders



queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

# SSTF (shortest-seek-time-first)

Request queue (0-199)
98, 183, 37, 122, 14, 124, 65, 67
Head pointer 53

# SCAN

☐ 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.

☐ **SCAN algorithm** Sometimes called the **elevator algorithm**

☐ Illustration shows total head movement of 208 cylinders

☐ But note that if requests are uniformly dense, largest density at other end of disk and those wait the longest
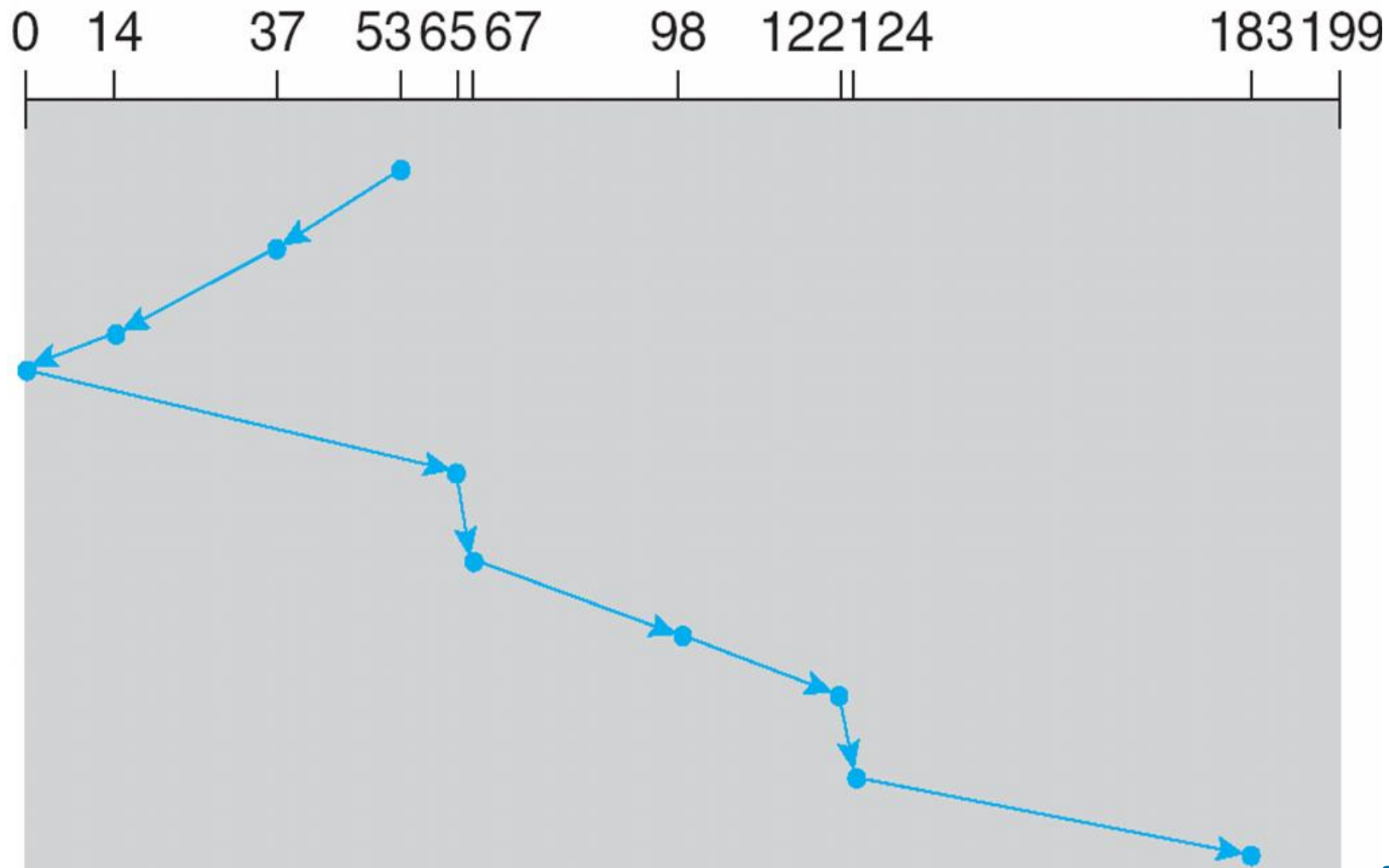
queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

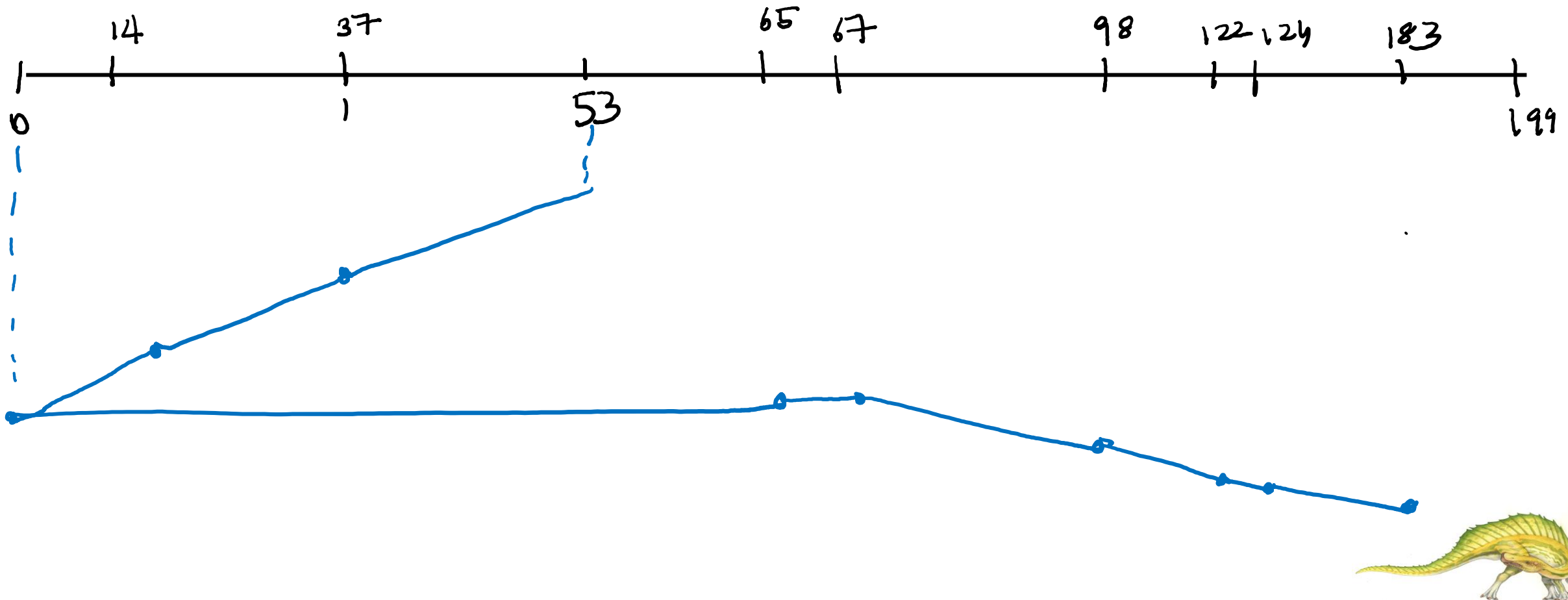0   14      37   53 65 67      98   122 124           183 199

(Left)

Direction (provided)

Request queue (0-199)
98, 183, 37, 122, 14, 124, 65, 67
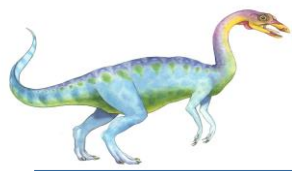Head pointer 53

# C-SCAN

- 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
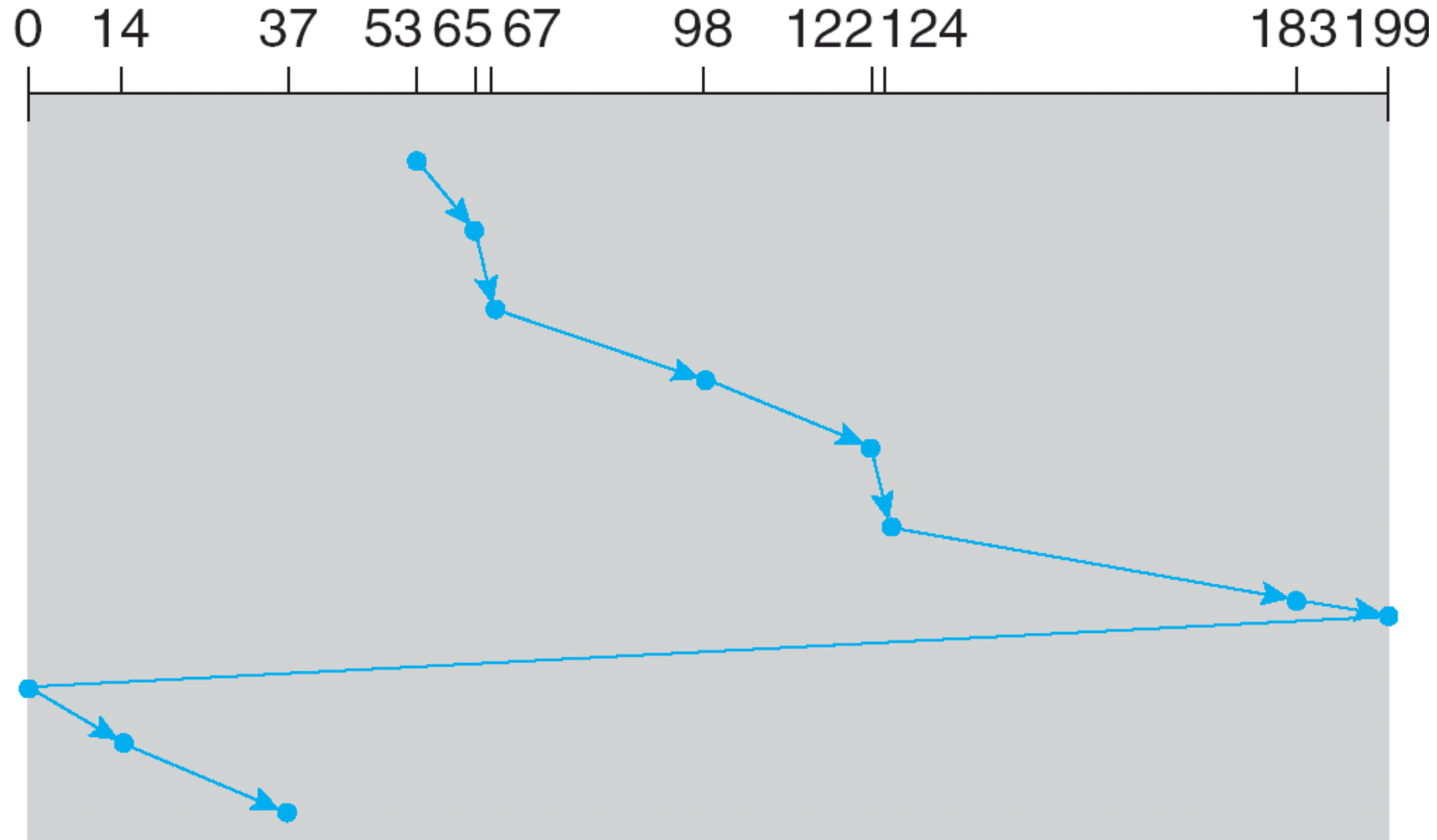
- Total number of cylinders?

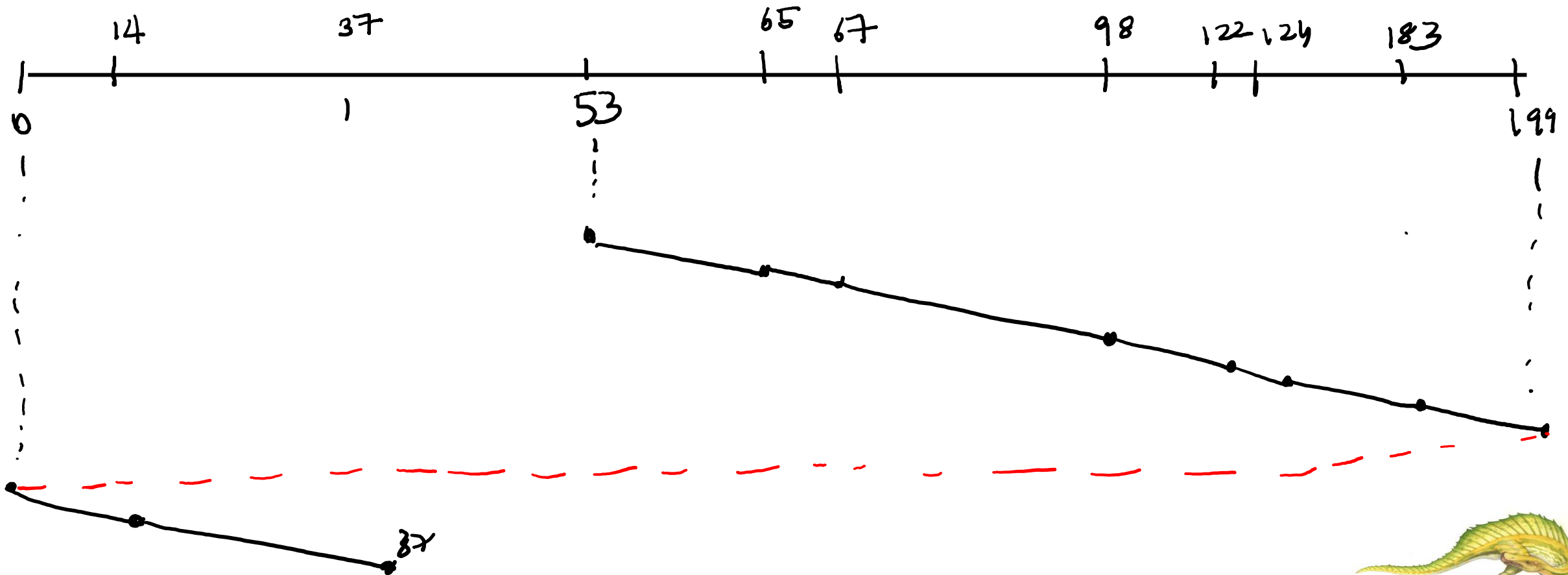queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

# C-SCAN (Circular)

Request queue (0-199)
98, 183, 37, 122, 14, 124, 65, 67
Head pointer 53

# C-LOOK

LOOK a version of SCAN, C-LOOK a 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
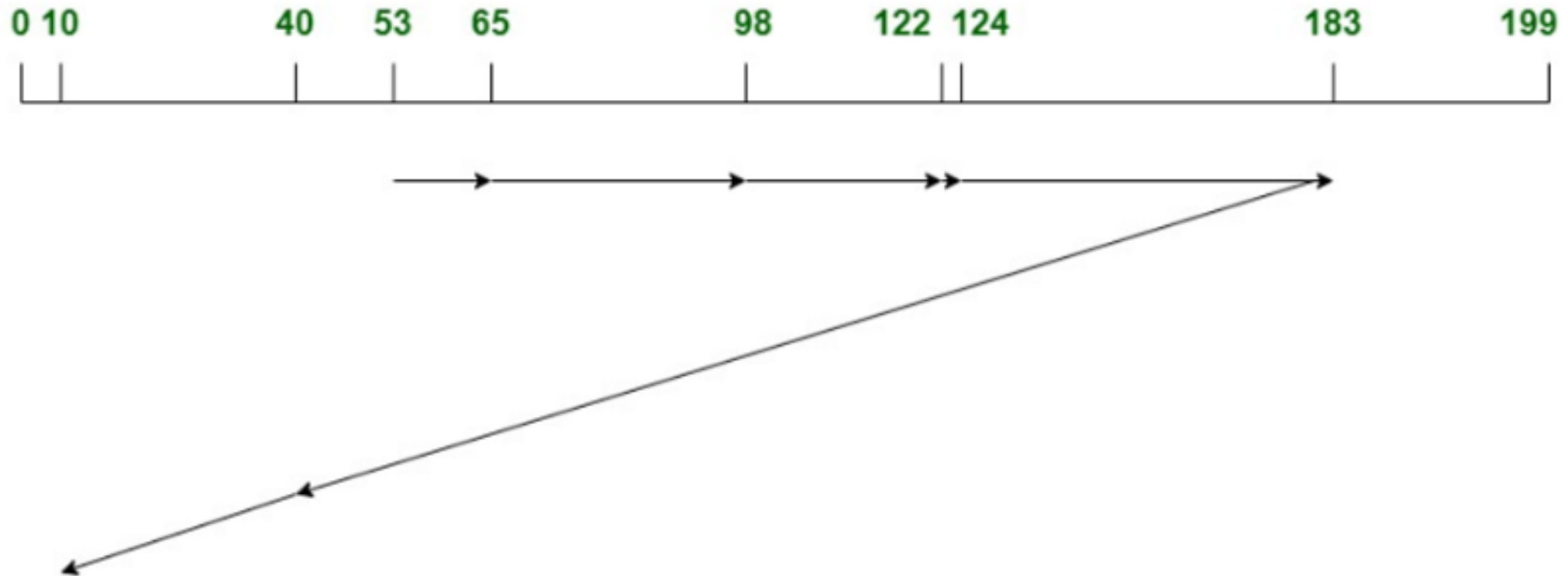
Total number of cylinders?

# LOOK

98, 183, 40, 122, 10, 124, 65. The current head position of the Read/Write head is 53 and will move in Right direction

QUEUE : 98,183,40,122,10,124,65      HEAD START AT 53 & moves in RIGHT Direction
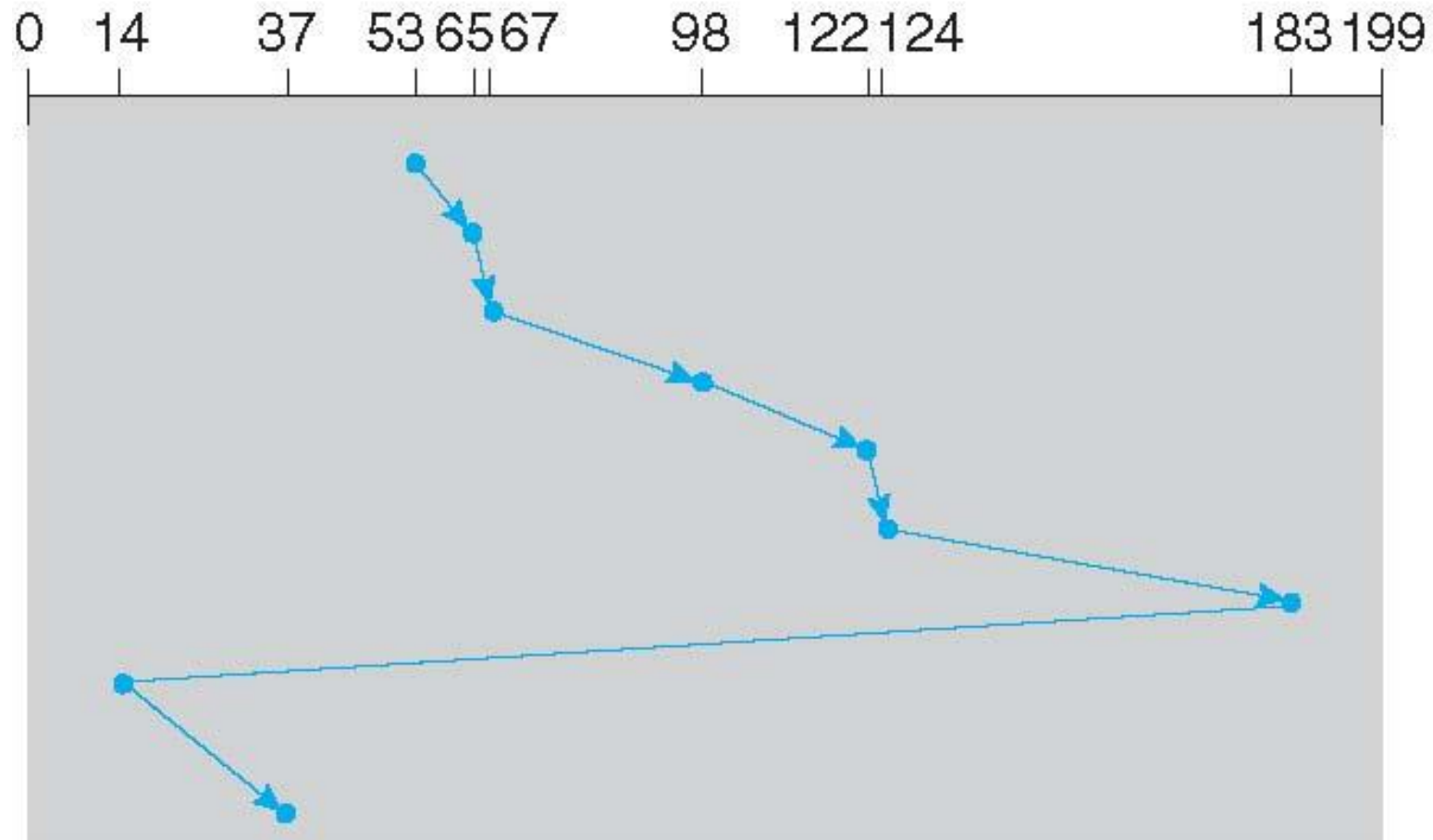
0 10          40    53    65          98      122 124          183        199
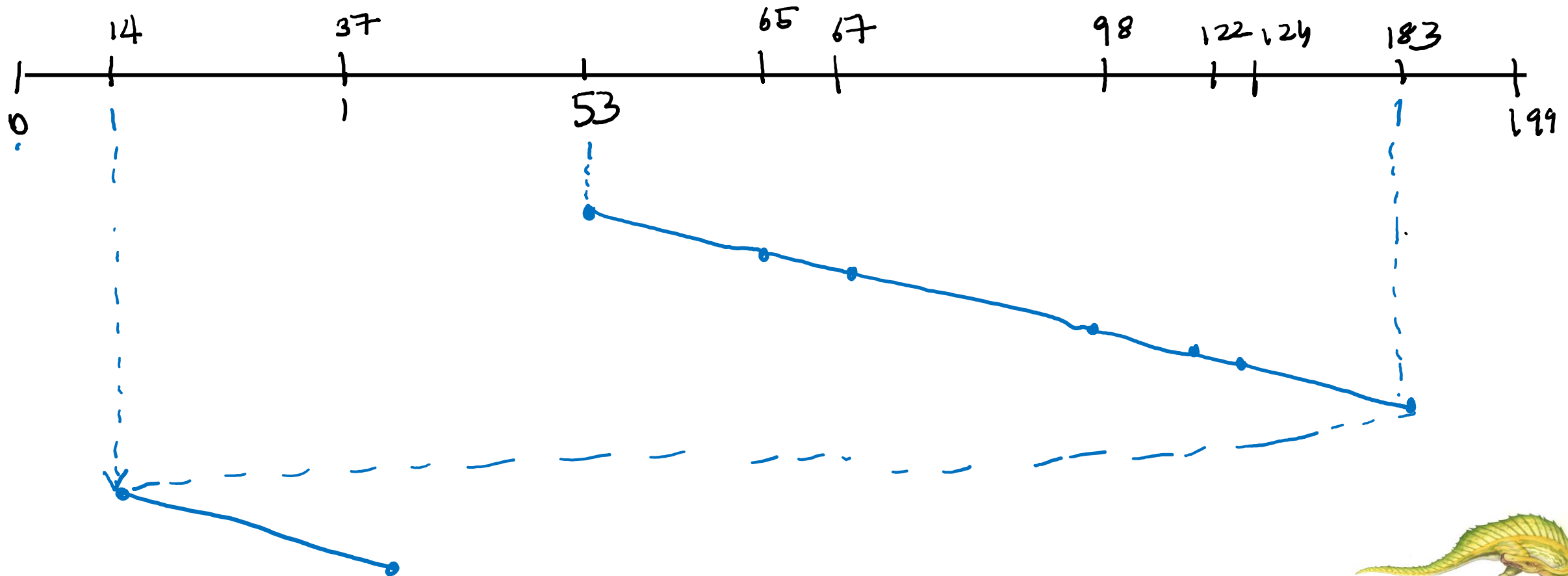
queue = 98, 183, 37, 122, 14, 124, 65, 67

head starts at 53

# C-LOOK

Request queue (0-199)
98, 183, 37, 122, 14, 124, 65, 67
Head pointer 53

# 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
- Performance depends on the number and types of requests
- Requests for disk service can be influenced by the file-allocation method
  - And metadata layout
- 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
- Either SSTF or LOOK is a reasonable choice for the default algorithm
- What about rotational latency?
  - Difficult for OS to calculate
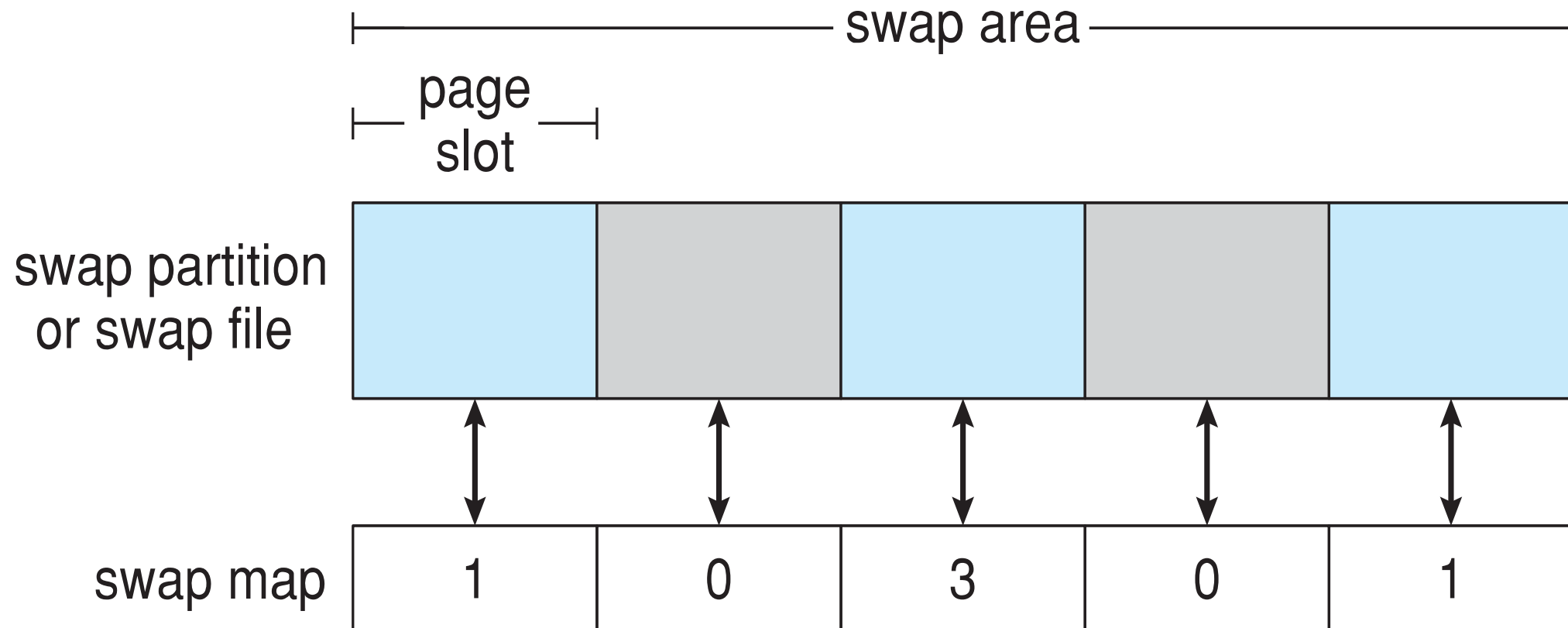- How does disk-based queueing effect OS queue ordering efforts?

# Swap-Space Management

- **Swap-space** — Virtual memory uses disk space as an extension of main memory
  - Less common now due to memory capacity increases
- Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition (raw)
- Swap-space management
  - 4.3BSD allocates swap space when process starts; holds text segment (the program) and data segment
  - Kernel uses **swap maps** to track swap-space use
  - Solaris 2 allocates swap space only when a dirty page is forced out of physical memory, not when the virtual memory page is first created
    - ▸ File data written to swap space until write to file system requested
    - ▸ Other dirty pages go to swap space due to no other home
    - ▸ Text segment pages thrown out and reread from the file system as needed
- What if a system runs out of swap space?
- Some systems allow multiple swap spaces

# End of Chapter 10