# Chapter 12: Mass-Storage Systems
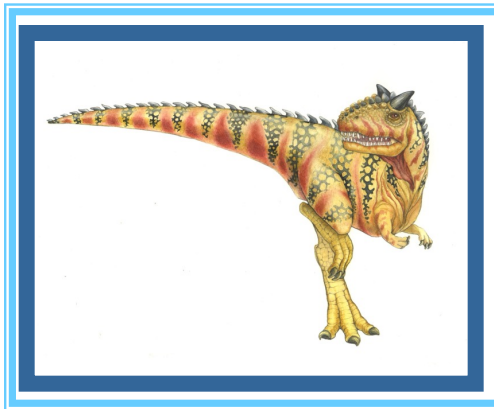
# Chapter 12:  Mass-Storage Systems

- Overview of Mass Storage Structure

- Disk Structure

- Disk Attachment

- Disk Scheduling

- Disk Management

- Swap-Space Management

- RAID Structure

- Stable-Storage Implementation

- Tertiary Storage Devices
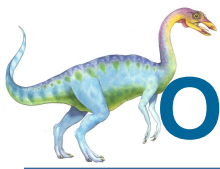
- Operating System Support

- Performance Issues

# 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;

◆ Discuss operating-system services provided for mass storage, including RAID and HSM;

# Overview of Mass Storage Structure

◆ Magnetic disks provide bulk of secondary storage of modern computers

  ➢ Drives rotate at 5400rpm、7200rpm、10000rpm、15000rpm

  ➢ 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)（寻道与旋转）

  ➢ Head crash results from disk head making contact with the disk surface

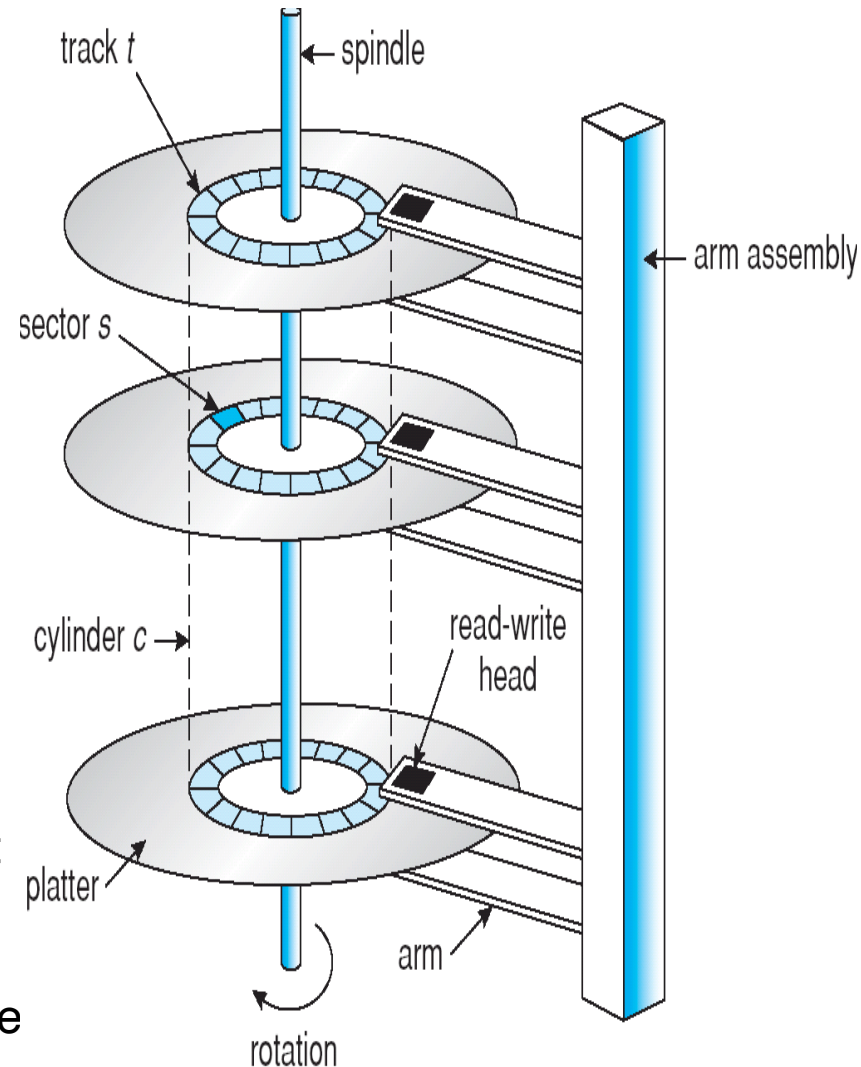    ◆ That's bad

# Overview of Mass Storage Structure

◆ Disks can be removable

◆ Drive attached to computer via I/O bus

  ➢ Busses vary, including EIDE, ATA, SATA, USB, Fibre Channel, SCSI

  ➢ Host controller in computer uses bus to talk to disk controller built into drive or storage array
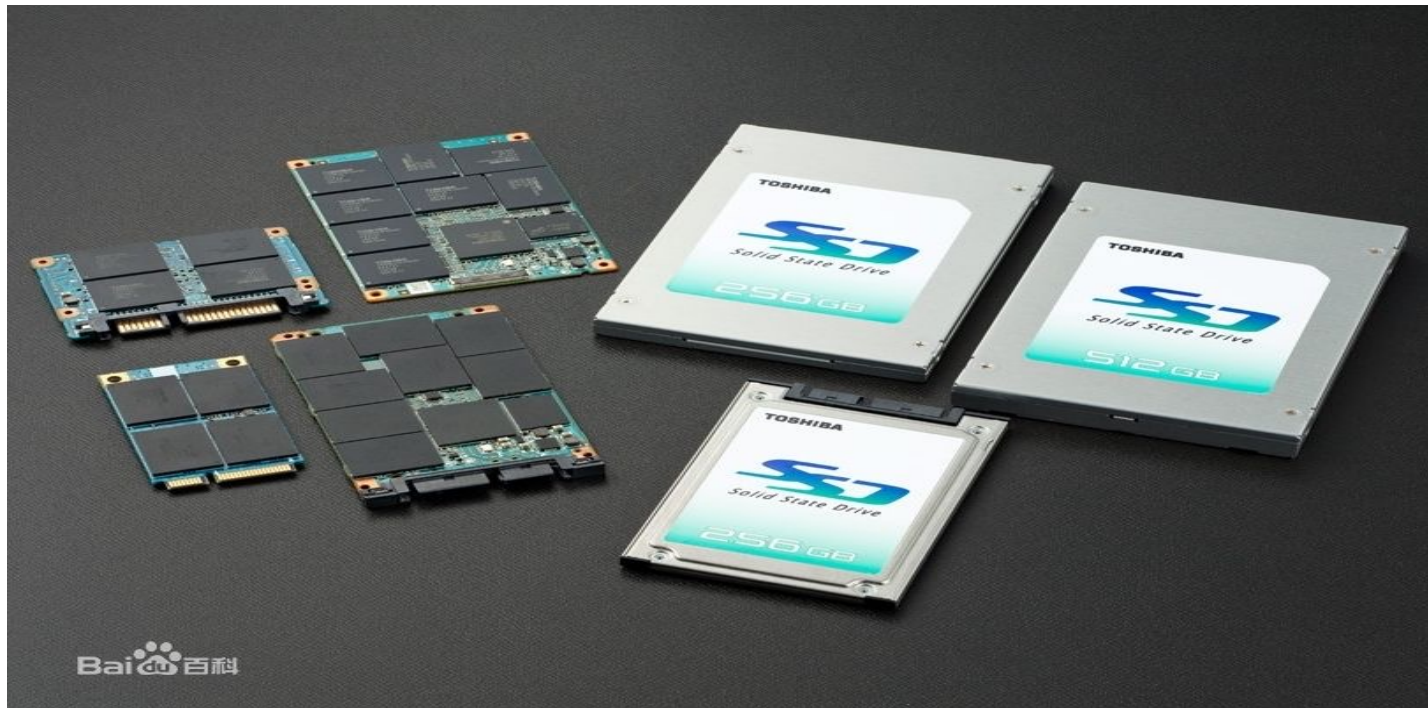
# Disk Structure

◆ 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

# SSD（Solid State Drives）

◆ 固态硬盘（Solid State Drives），简称固盘，用固态电子存储芯片阵列而制成的硬盘，包括控制单元和存储单元（FLASH芯片、DRAM芯片）。

◆ 新一代的固态硬盘普遍采用SATA-2接口、SATA-3接口、SAS接口、MSATA接口、PCI-E接口、NGFF接口、CFast接口和SFF-8639接口。**读写速度快、防震抗摔性、低功耗、无噪音、温度范围大，轻便等**

# Magnetic tape

◆ Was early secondary-storage medium

◆ 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

◆ 20-200GB typical storage

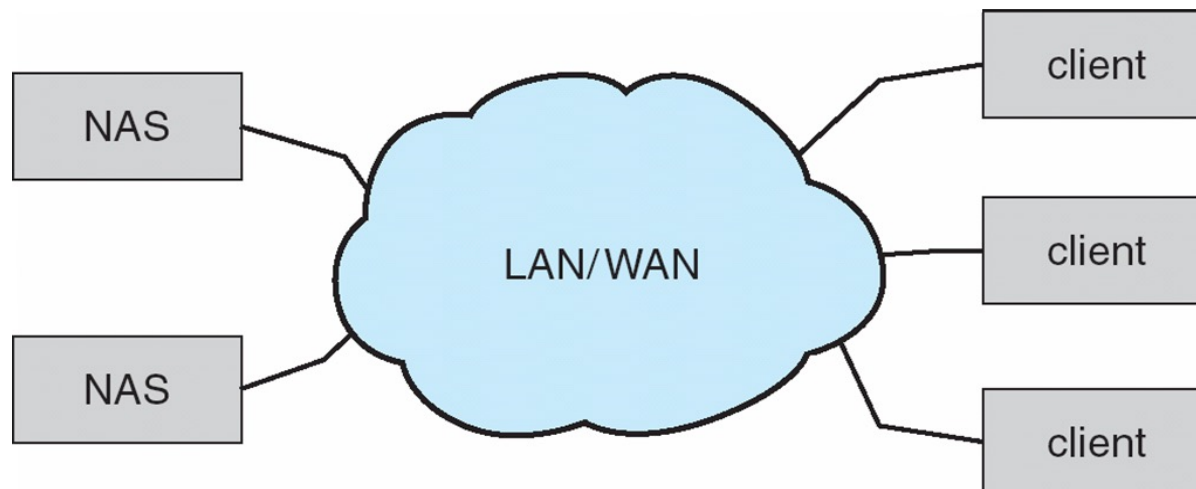◆ Common technologies are 4mm, 8mm, 19mm, LTO-2 and SDLT

# Disk Attachment

◆ Host-attached storage accessed through I/O ports talking to I/O busses

◆ SCSI itself is a bus, up to 16 devices on one cable, SCSI initiator requests operation and SCSI targets perform tasks

  ➢ Each target can have up to 8 logical units (disks attached to device controller

◆ FC is high-speed serial architecture

  ➢ Can be switched fabric with 24-bit address space – the basis of storage area networks (SANs) in which many hosts attach to many storage units

  ➢ Can be arbitrated loop (FC-AL) of 126 devices

# **Network-Attached Storage**

◆ Network-attached storage (NAS) is storage made available over a network rather than over a local connection (such as a bus)

◆ NFS and CIFS (Common Internet File System) are common protocols

◆ Implemented via remote procedure calls (RPCs) between host and storage

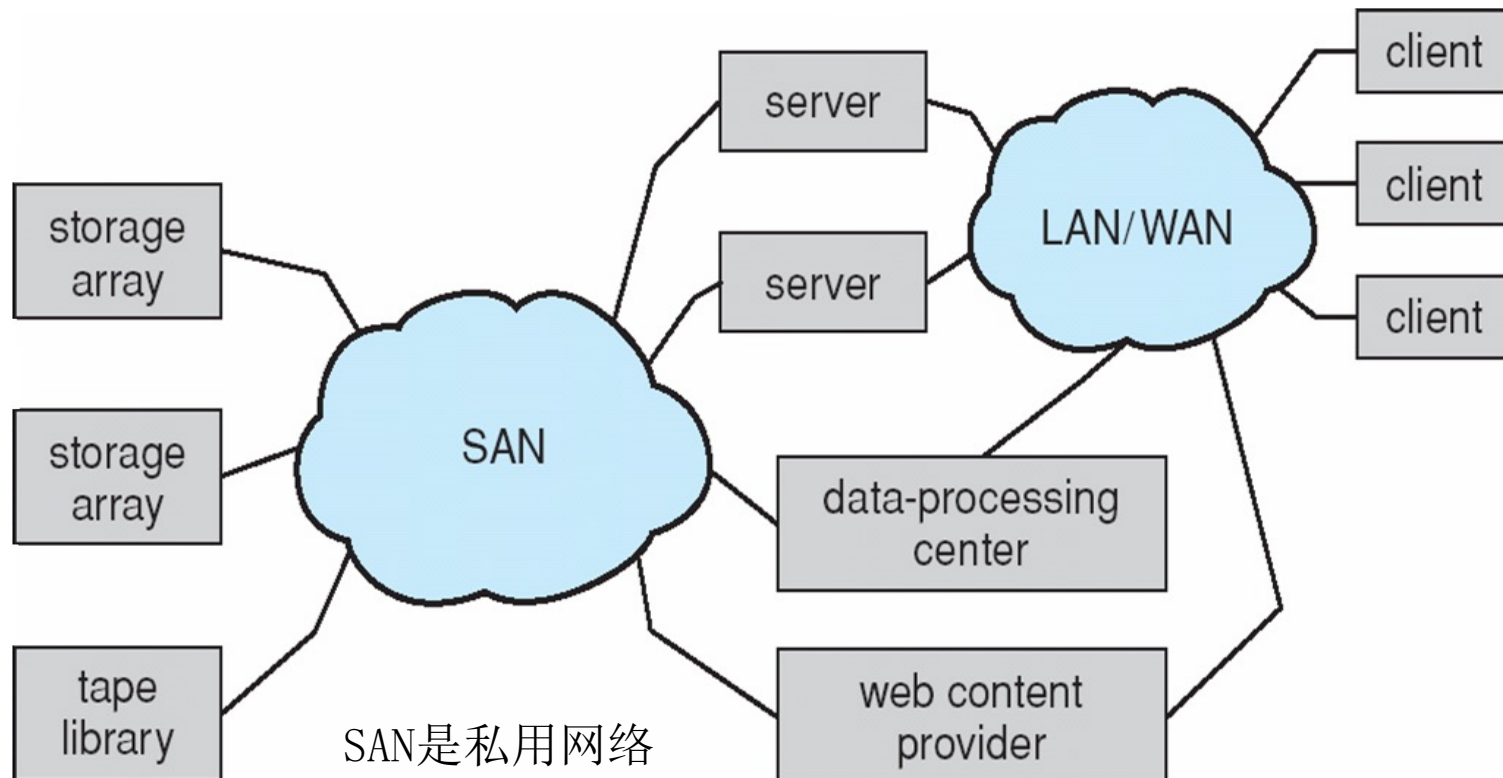◆ New iSCSI protocol uses IP network to carry the SCSI protocol

占用公网资源，增加通信延迟

# Storage Area Network

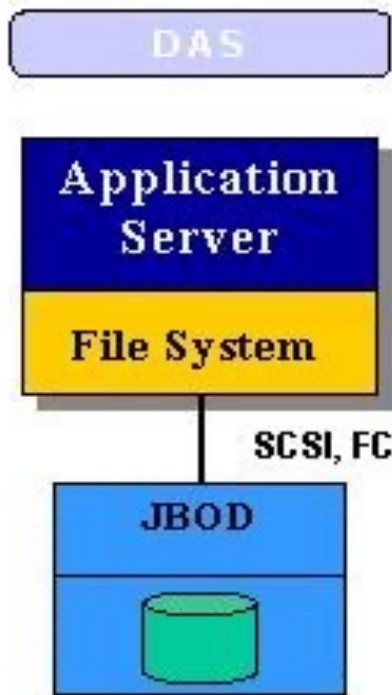◆ Common in large storage environments (and becoming more common)

◆ Multiple hosts attached to multiple storage arrays - flexible
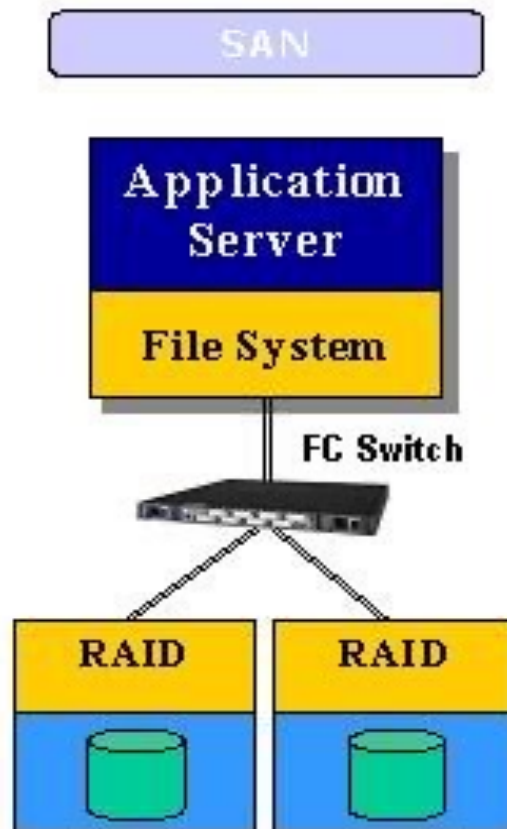
SAN是私用网络
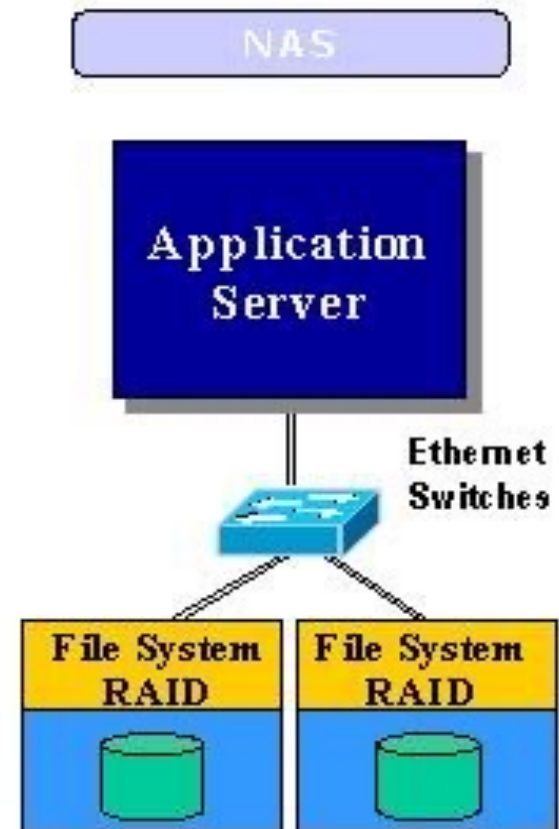
# DAS,NAS,SAN

Direct-Attached Storage      Storage Area Network      Network-attached storage

# 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

  ➢ 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

◆ Minimize seek time: Seek time ≈ seek distance

◆ 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(单位时间内传输的量)

# Disk Scheduling (Cont)

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

◆ We illustrate them with a request queue (0-199)

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

Head pointer 53
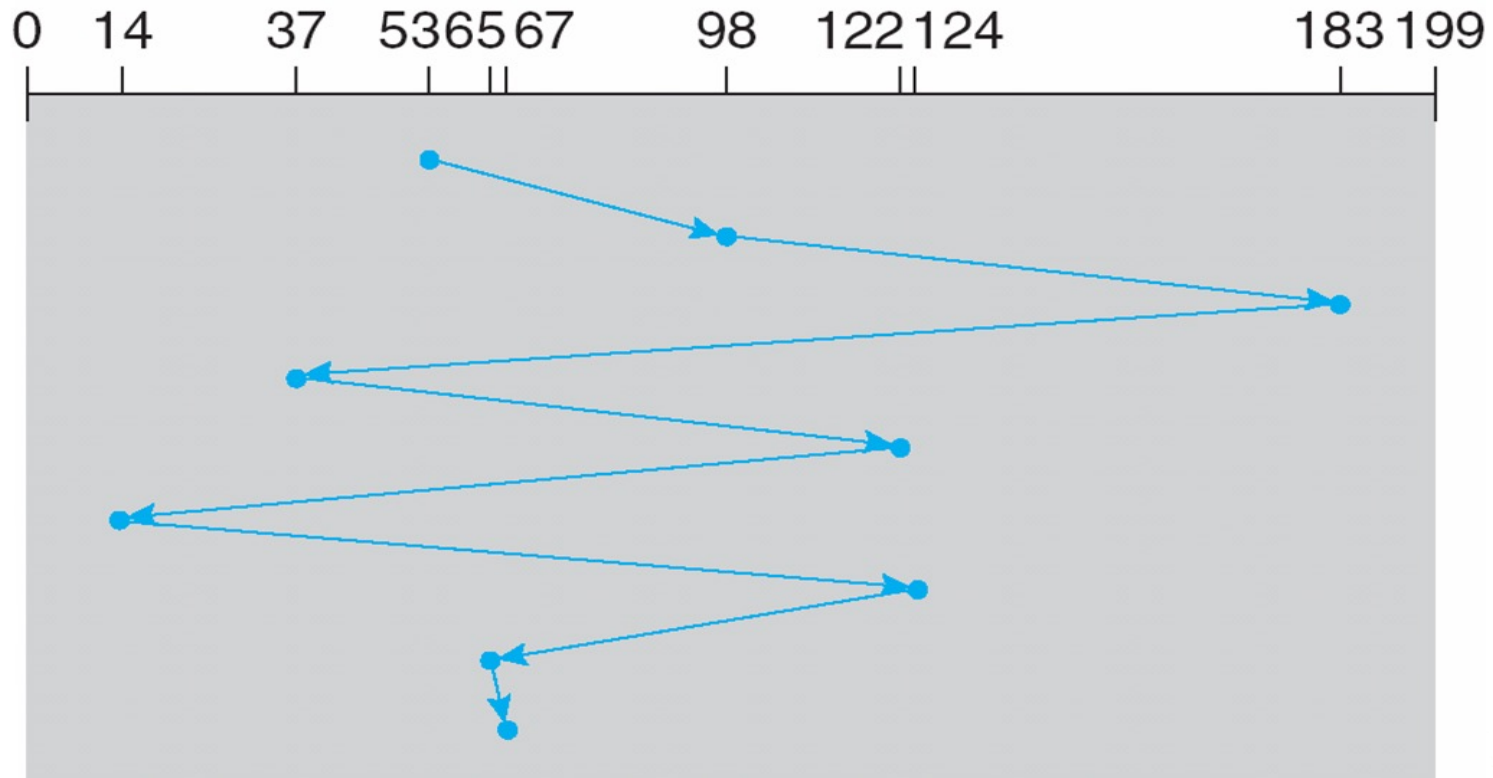
# FCFS

Illustration shows total head movement of 640 cylinders



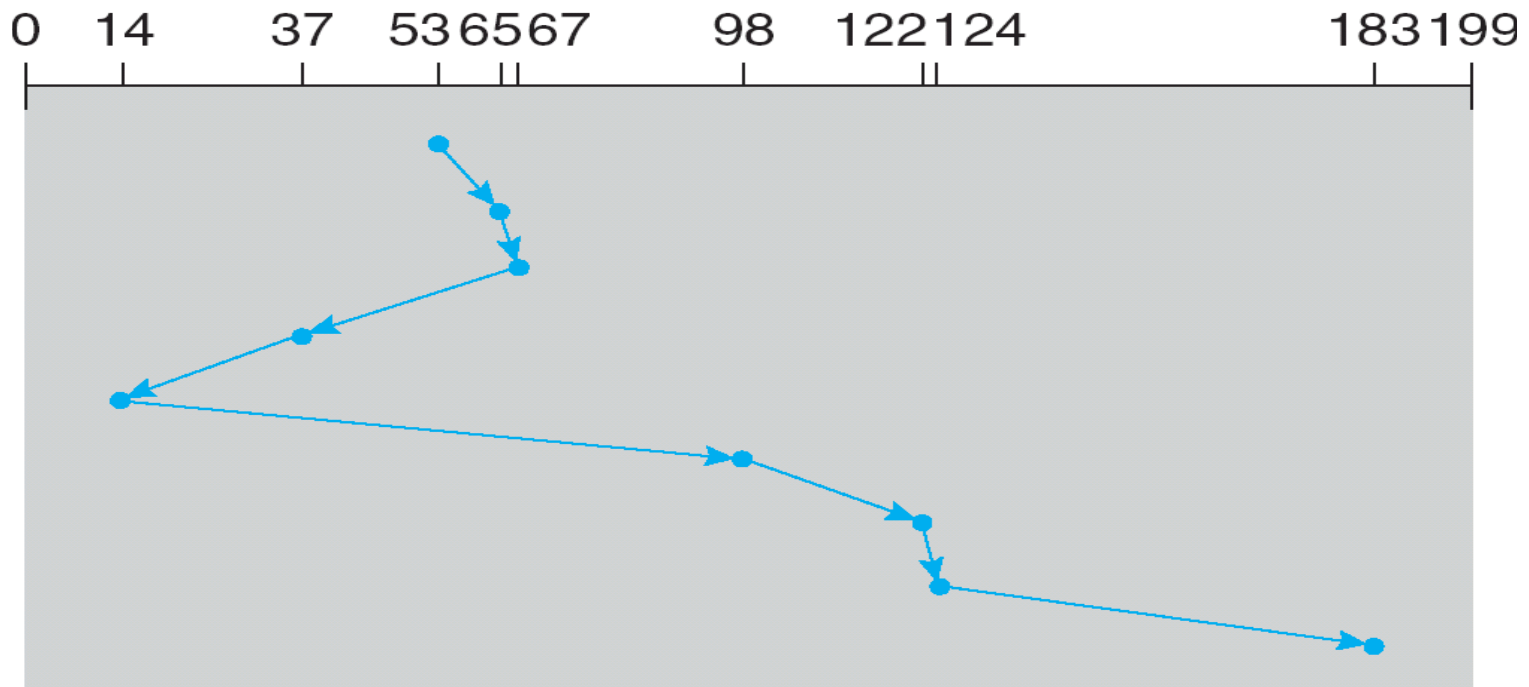queue = 98, 183, 37, 122, 14, 124, 65, 67
head starts at 53

# SSTF

◆ Selects the request with the minimum seek time from the current head position. total head movement of 236 cylinders

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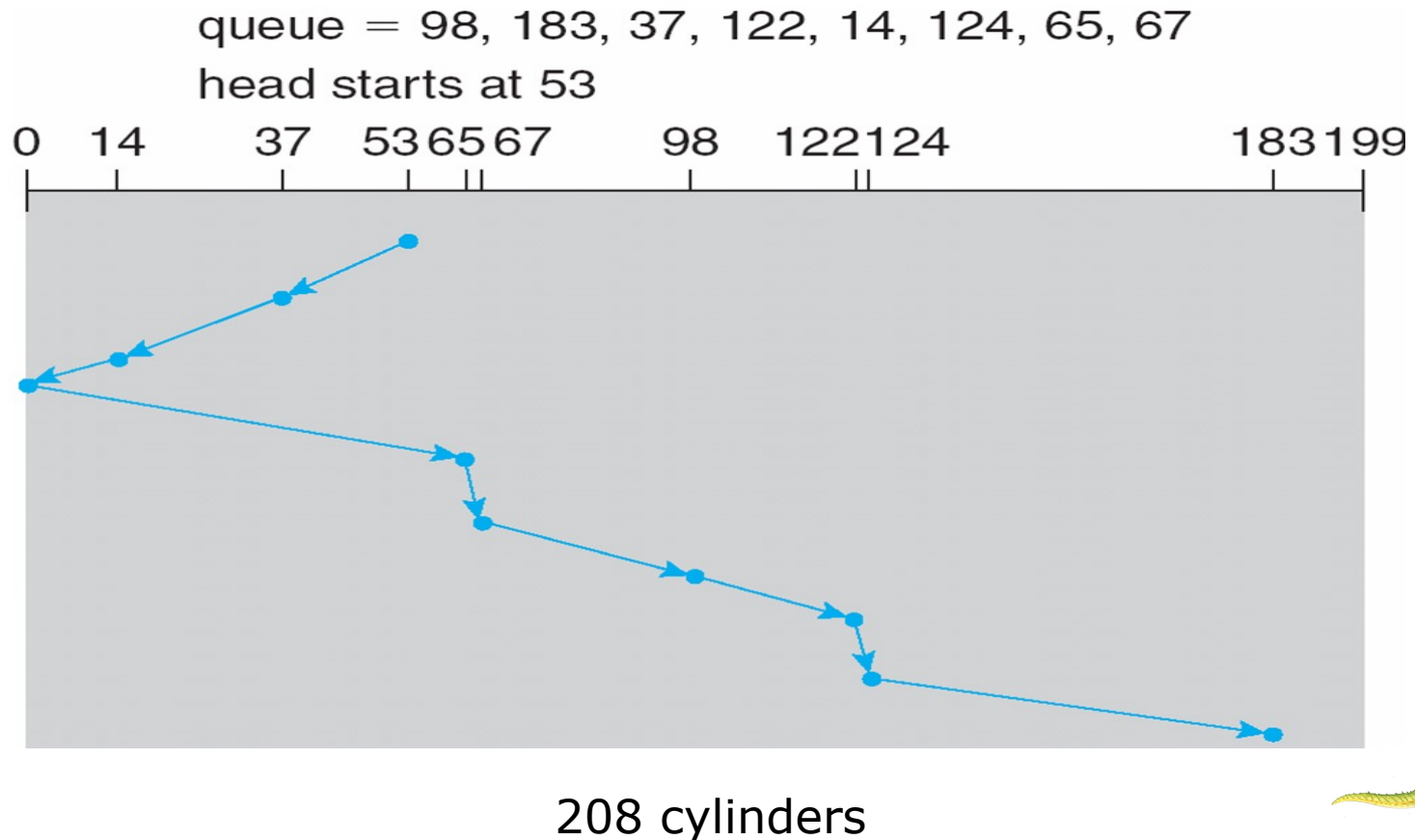
head starts at 53



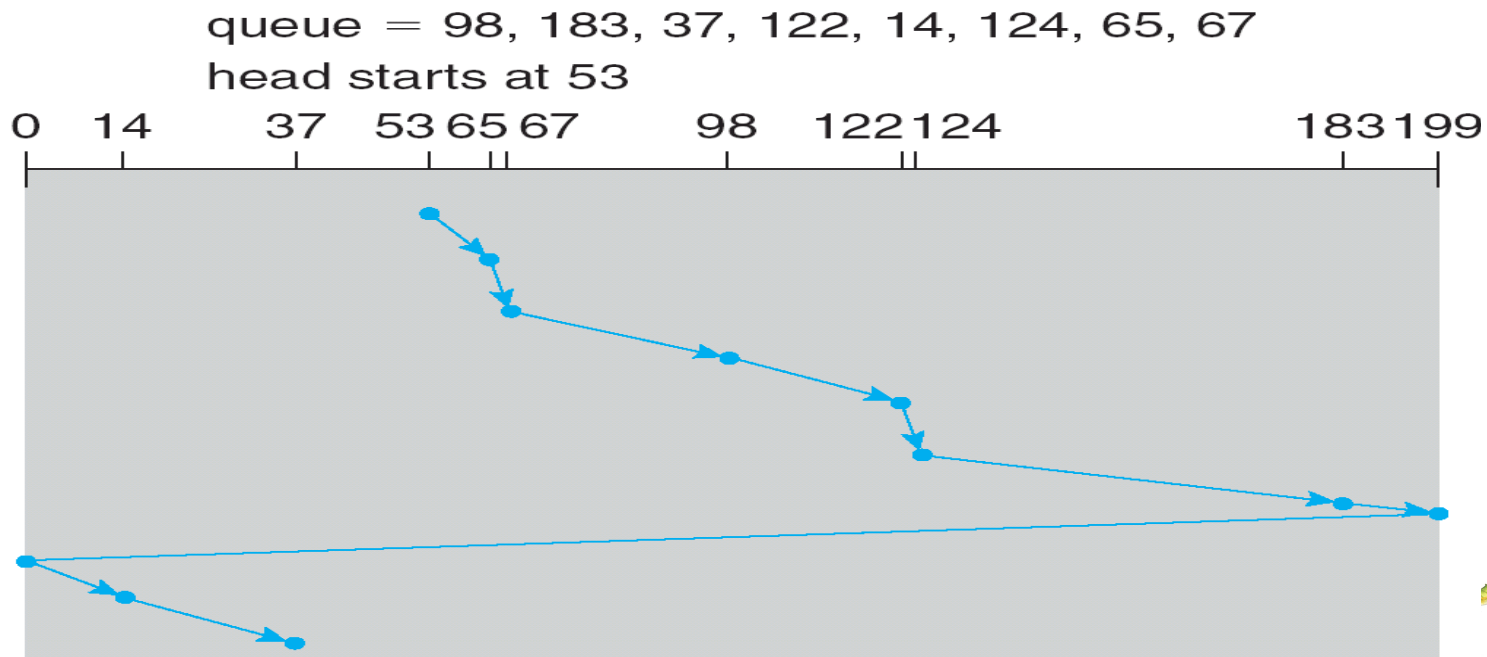◆ SSTF scheduling is a form of SJF scheduling; may cause starvation of some requests

# 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

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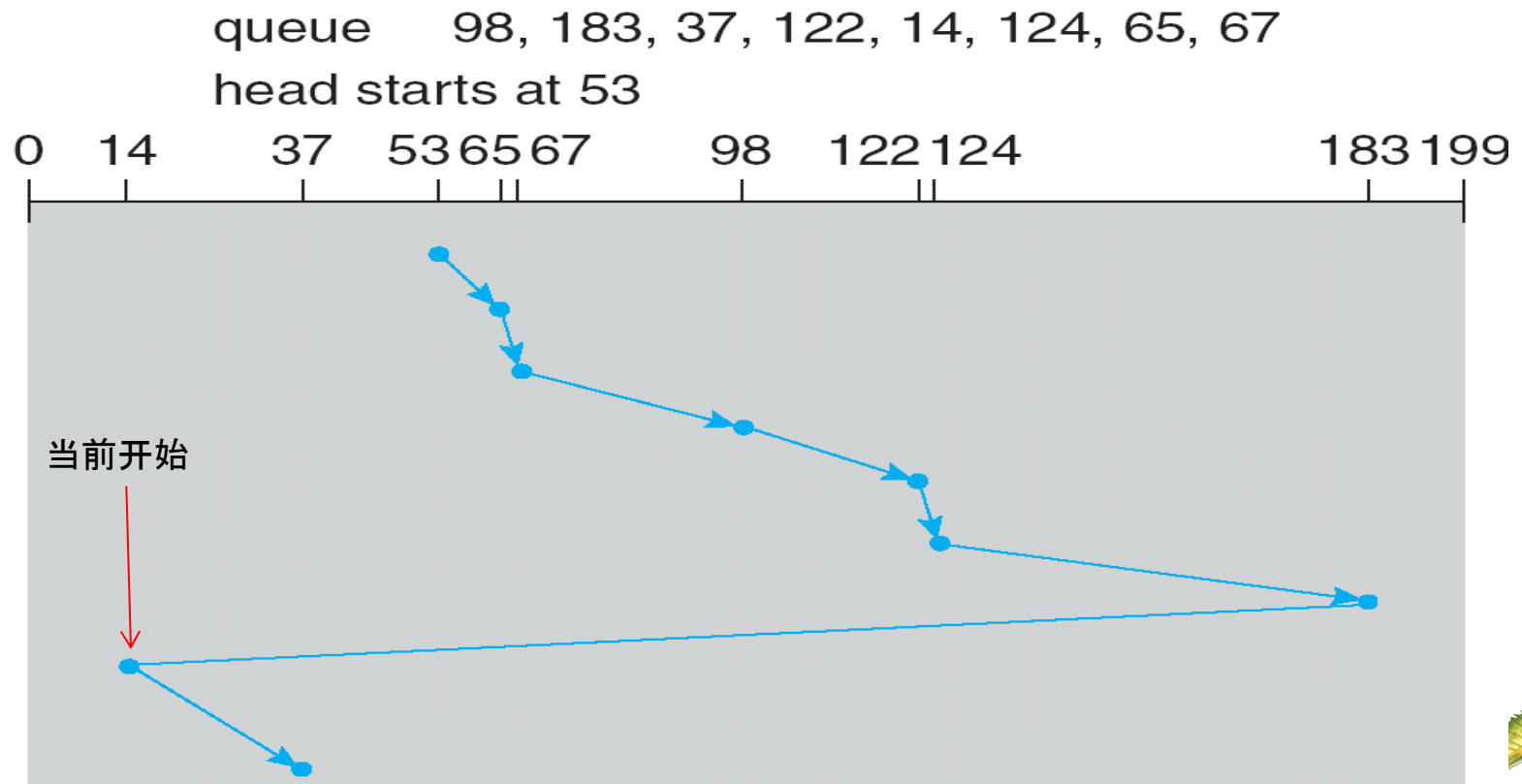head starts at 53



208 cylinders

# 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

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

# 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

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

当前开始

# Selecting a Disk-Scheduling Algorithm

◆ SSTF is common and has a natural appeal(比FCFS好)

◆ SCAN and C-SCAN perform better for systems that place a heavy load on the disk(重载时不会出现"饥饿"状态)

◆ Performance depends on the number and types of requests

◆ Requests for disk service can be influenced by the file-allocation method

◆ 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(算法可置换)

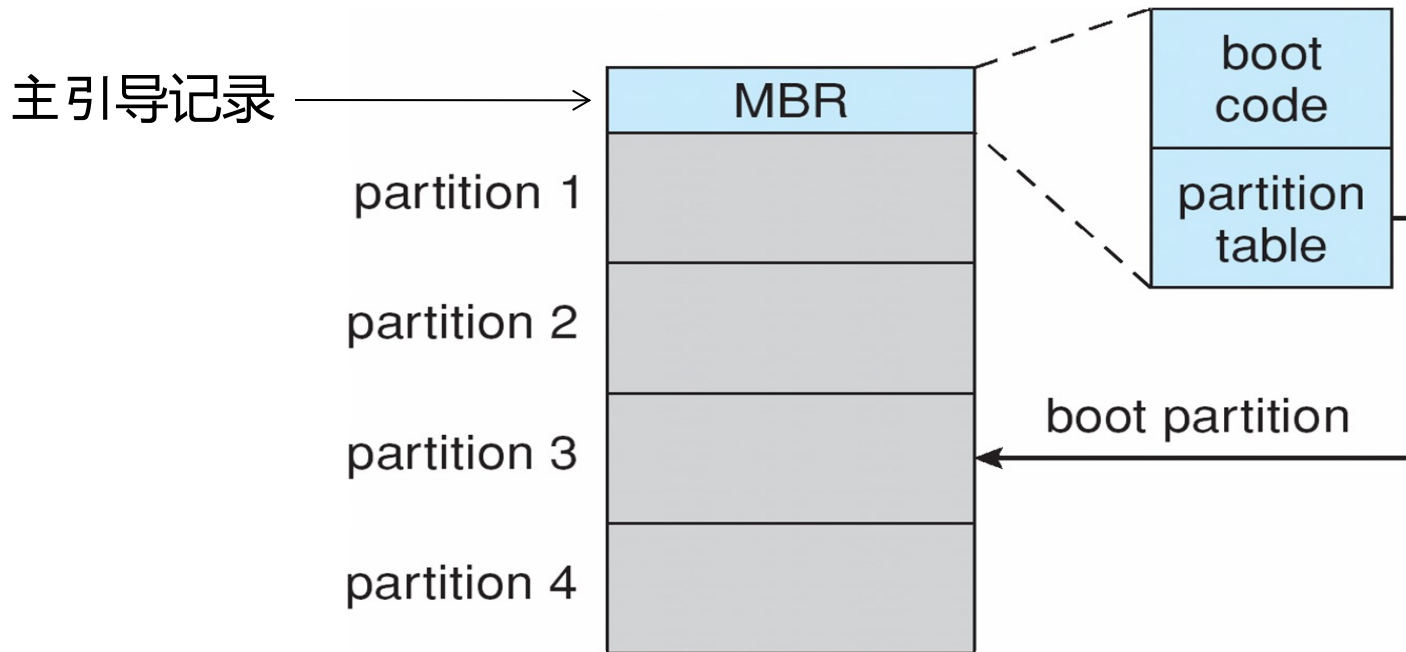◆ Either SSTF or LOOK is a reasonable choice for the default algorithm

# Disk Management

◆ Disk Formatting:

   ❑ Low-level formatting, or physical formatting

   ❑ Logical formatting

◆ Low-level formatting, or physical formatting — Dividing a disk into sectors that the disk controller can read and write

◆ Logical formatting-- To use a disk to hold files, the operating system still needs to record its own data structures on the disk

   ➢ Partition the disk into one or more groups of cylinders

   ➢ Logical formatting or "making a file system"

   ➢ To increase efficiency most file systems group blocks into clusters

      ◆ Disk I/O done in blocks

      ◆ File I/O done in clusters

# Disk Management

◆ Boot block initializes system

  ➢ The bootstrap is stored in ROM

  ➢ Bootstrap loader program

主引导记录 ⟶



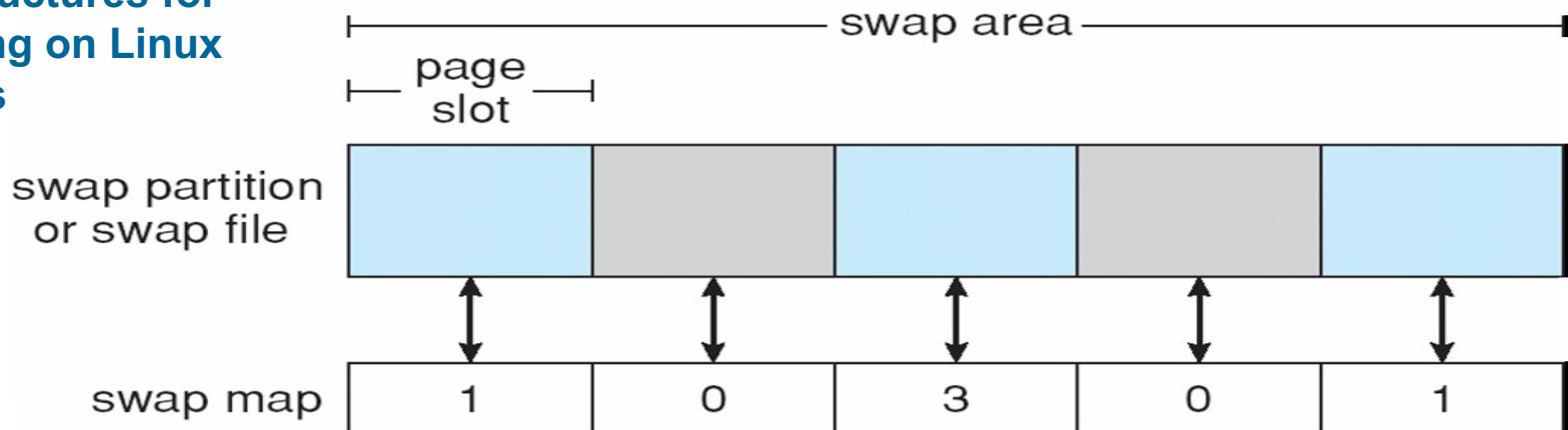◆ Methods such as sector sparing（备用）used to handle bad blocks

# Swap-Space Management

◆ Swap-space — Virtual memory uses disk space as an extension of main memory

◆ 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 page is forced out of physical memory, not when the virtual memory page is first created

**Data Structures for Swapping on Linux Systems**

swap area

page slot

swap partition or swap file

swap map

| 1 | 0 | 3 | 0 | 1 |
|---|---|---|---|---|

0表示对应的页槽可用

3表示被换出的页映射到
3个进程（内存共享）

# RAID Structure

◆ RAID(**R**edundant **A**rray of **I**ndependent **D**isks) – multiple disk drives provides reliability via redundancy;

◆ Increases the mean time to failure(平均故障间隔期);

◆ Frequently combined with NVRAM (Non-Volatile RAM) to improve write performance;

◆ RAID is arranged into six different levels;

# RAID (Cont)

◆ Disk striping uses a group of disks as one storage unit

◆ RAID schemes improve performance and improve the reliability of the storage system by storing redundant data

- ➢ Mirroring or shadowing (RAID 1) keeps duplicate of each disk

- ➢ Striped mirrors (RAID 1+0) or mirrored stripes (RAID 0+1) provides high performance and high reliability

- ➢ Block interleaved parity (RAID 4, 5, 6) uses much less redundancy

◆ RAID within a storage array can still fail if the array fails, so automatic replication of the data between arrays is common

◆ Frequently, a small number of hot-spare disks are left unallocated, automatically replacing a failed disk and having data rebuilt onto them.

# RAID Levels



(a) RAID 0: non-redundant striping.

(b) RAID 1: mirrored disks.

(c) RAID 2: memory-style error-correcting codes.

(d) RAID 3: bit-interleaved parity. 位交错奇偶校验

(e) RAID 4: block-interleaved parity. 块交错奇偶校验
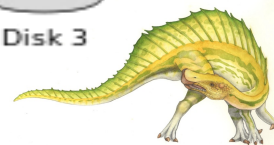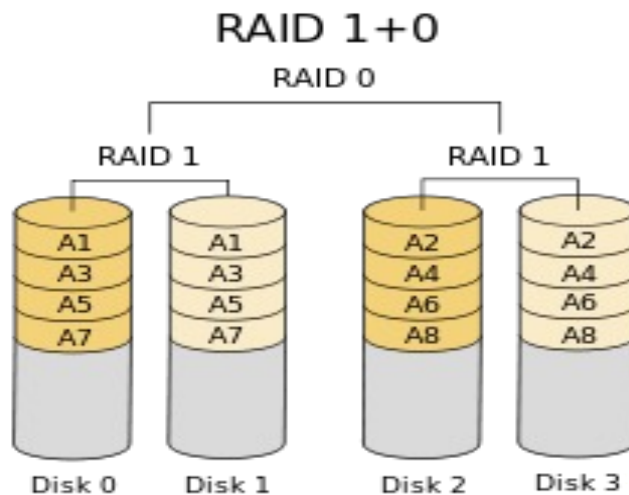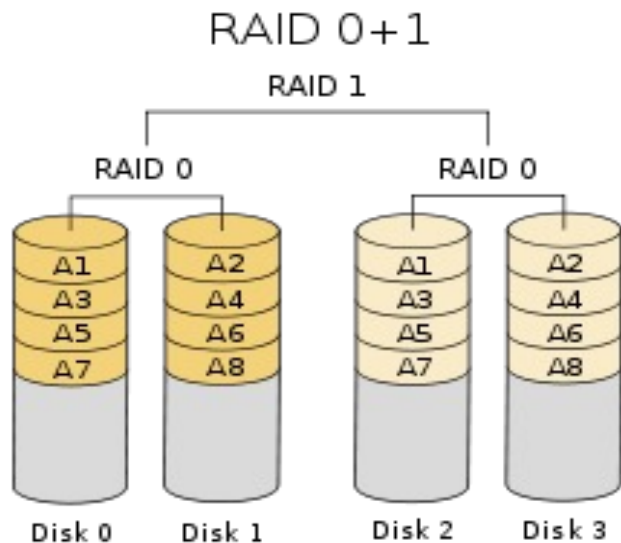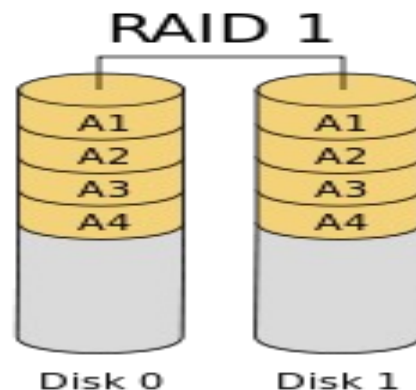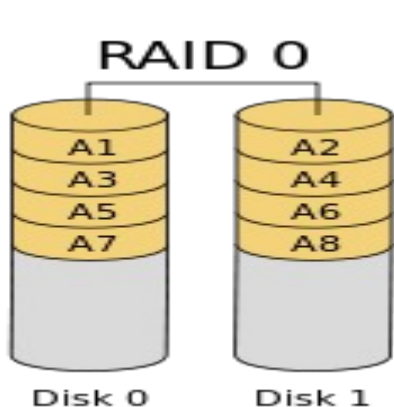
(f) RAID 5: block-interleaved distributed parity.

(g) RAID 6: P + Q redundancy.
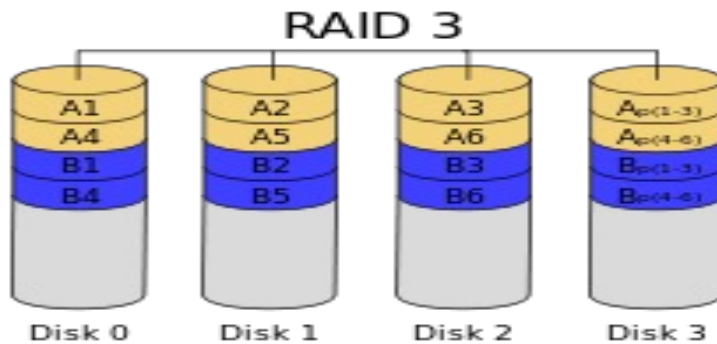
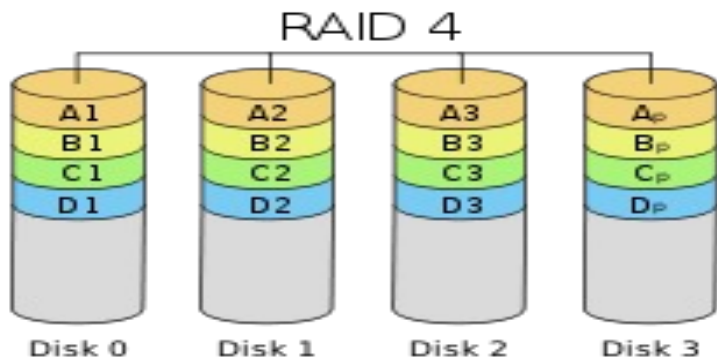# RAID0 and RAID1

# RAID (2--4)

RAID 2



RAID2：海明码纠错、数据分块、并行访问、适合大批量数据、使用少；

RAID 3



RAID3：采用Bit－nterleaving（数据交错存储）技术，通过编码再将数据比特分区后分别存在硬盘中；

RAID 4


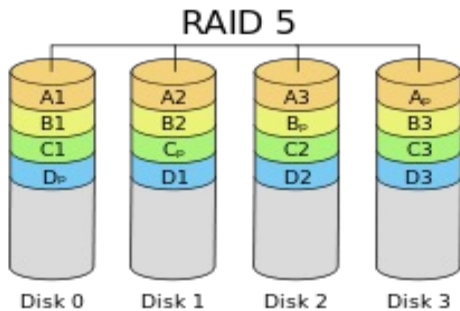
RAID4：与RAID 3不同它以块为单位分区；

# RAID (5,50,6,60)
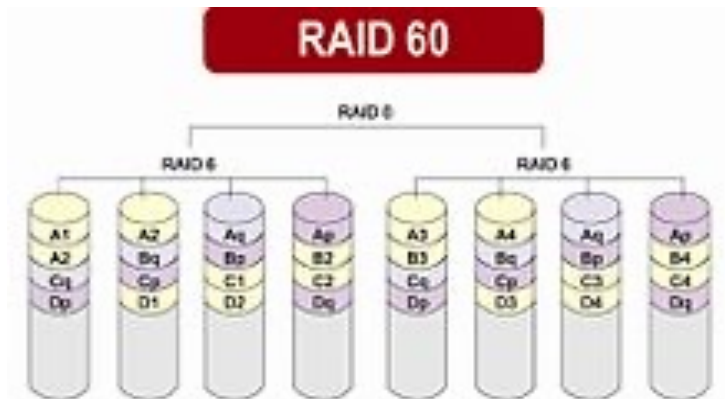
RAID5：独立存取、无单独校验盘
、适合访问频繁、传输率低；

RAID6:与RAID 5相比，RAID 6增加
了第二个独立的奇偶校验信息块



RAID50

| metadata block 1 | |
|---|---|
| address 1 | address 2 |
| checksum MB2 | checksum |

| metadata block 2 | |
|---|---|
| address | address |
| checksum D1 | checksum D2 |

| data 1 |
|---|

| data 2 |
|---|

它是一个标准的POSIX文件系统，一种汇集存储模式，

(a) Traditional volumes and file systems.

(b) ZFS and pooled storage.

# Stable-Storage Implementation

◆ Write-ahead log scheme requires stable storage

◆ To implement stable storage:

  ➢ Replicate information on more than one nonvolatile storage media with independent failure modes

  ➢ Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery



Stable Storage, Crash after drive 1 is updated Bad spot

# Tertiary Storage Devices

◆ Low cost is the defining characteristic of tertiary storage

◆ Generally, tertiary storage is built using removable media

◆ Common examples of removable media are floppy disks and CD-ROMs; other types are available

# Removable Disks

◆ Floppy disk — thin flexible disk coated with magnetic material, enclosed in a protective plastic case

➢ Most floppies hold about 1 MB; similar technology is used for removable disks that hold more than 1 GB

➢ Removable magnetic disks can be nearly as fast as hard disks, but they are at a greater risk of damage from exposure

# Removable Disks (Cont.)

◆ A magneto-optic disk records data on a rigid platter coated with magnetic material

  ➢ Laser heat is used to amplify a large, weak magnetic field to record a bit

  ➢ Laser light is also used to read data (Kerr effect)

  ➢ The magneto-optic head flies much farther from the disk surface than a magnetic disk head, and the magnetic material is covered with a protective layer of plastic or glass; resistant to head crashes

◆ Optical disks do not use magnetism; they employ special materials that are altered by laser light

# WORM Disks

◆ The data on read-write disks can be modified over and over

◆ WORM ("Write Once, Read Many Times") disks can be written only once

◆ Thin aluminum film sandwiched between two glass or plastic platters

◆ To write a bit, the drive uses a laser light to burn a small hole through the aluminum; information can be destroyed by not altered

◆ Very durable and reliable

◆ Read-only disks, such ad CD-ROM and DVD, com from the factory with the data pre-recorded

# Operating System Support

◆ Major OS jobs are to manage physical devices and to present a virtual machine abstraction to applications

◆ For hard disks, the OS provides two abstraction:

➢ Raw device – an array of data blocks

➢ File system – the OS queues and schedules the interleaved requests from several applications

# Application Interface

◆ Most OSs handle removable disks almost exactly like fixed disks — a new cartridge is formatted and an empty file system is generated on the disk

◆ Tapes are presented as a raw storage medium, i.e., and application does not not open a file on the tape, it opens the whole tape drive as a raw device

◆ Usually the tape drive is reserved for the exclusive use of that application

◆ Since the OS does not provide file system services, the application must decide how to use the array of blocks

◆ Since every application makes up its own rules for how to organize a tape, a tape full of data can generally only be used by the program that created it

# Tape Drives

◆ The basic operations for a tape drive differ from those of a disk drive

◆ `locate()` positions the tape to a specific logical block, not an entire track (corresponds to `seek()`)

◆ The `read position()` operation returns the logical block number where the tape head is

◆ The `space()` operation enables relative motion

◆ Tape drives are "append-only" devices; updating a block in the middle of the tape also effectively erases everything beyond that block

◆ An EOT mark is placed after a block that is written

# File Naming

◆ The issue of naming files on removable media is especially difficult when we want to write data on a removable cartridge on one computer, and then use the cartridge in another computer

◆ Contemporary OSs generally leave the name space problem unsolved for removable media, and depend on applications and users to figure out how to access and interpret the data

◆ Some kinds of removable media (e.g., CDs) are so well standardized that all computers use them the same way

# Hierarchical Storage Management (HSM)

◆ A hierarchical storage system extends the storage hierarchy beyond primary memory and secondary storage to incorporate tertiary storage — usually implemented as a jukebox of tapes or removable disks

◆ Usually incorporate tertiary storage by extending the file system

  ➢ Small and frequently used files remain on disk

  ➢ Large, old, inactive files are archived to the jukebox

◆ HSM is usually found in supercomputing centers and other large installations that have enormous volumes of data
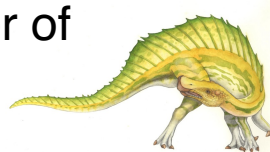
# Speed

◆ Two aspects of speed in tertiary storage are bandwidth and latency

◆ Bandwidth is measured in bytes per second

  ◆ Sustained bandwidth – average data rate during a large transfer; # of bytes/transfer time
    Data rate when the data stream is actually flowing

  ◆ Effective bandwidth – average over the entire I/O time, including `seek()` or `locate()`, and cartridge switching
    Drive's overall data rate

# Speed (Cont)

◆ Access latency – amount of time needed to locate data

  ➢ Access time for a disk – move the arm to the selected cylinder and wait for the rotational latency; < 35 milliseconds

  ➢ Access on tape requires winding the tape reels until the selected block reaches the tape head; tens or hundreds of seconds

  ➢ Generally say that random access within a tape cartridge is about a thousand times slower than random access on disk

◆ The low cost of tertiary storage is a result of having many cheap cartridges share a few expensive drives

◆ A removable library is best devoted to the storage of infrequently used data, because the library can only satisfy a relatively small number of I/O requests per hour

# Reliability

◆ A fixed disk drive is likely to be more reliable than a removable disk or tape drive

◆ An optical cartridge is likely to be more reliable than a magnetic disk or tape

◆ A head crash in a fixed hard disk generally destroys the data, whereas the failure of a tape drive or optical disk drive often leaves the data cartridge unharmed
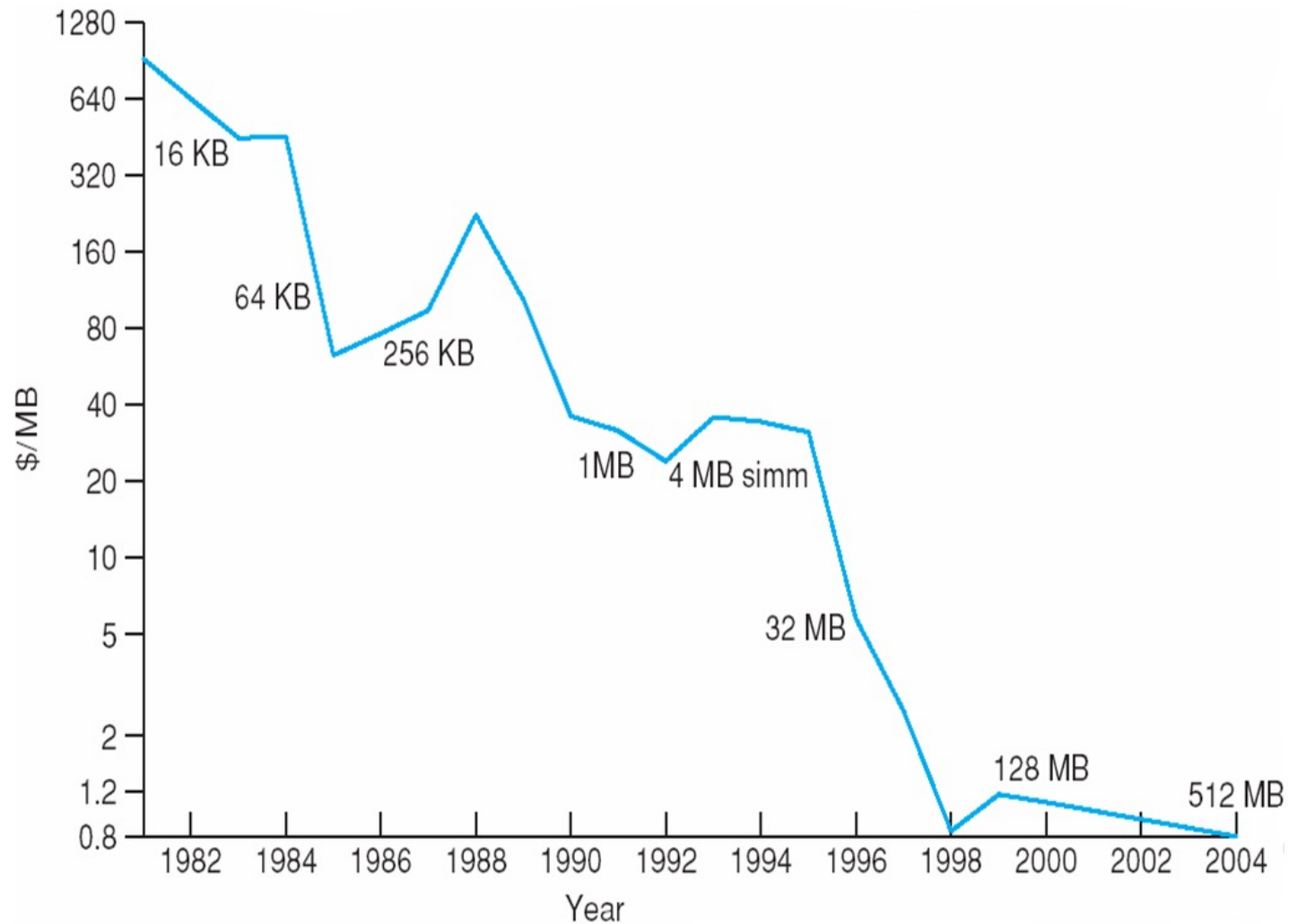
# Cost

◆ Main memory is much more expensive than disk storage

◆ The cost per megabyte of hard disk storage is competitive with magnetic tape if only one tape is used per drive

◆ The cheapest tape drives and the cheapest disk drives have had about the same storage capacity over the years

◆ Tertiary storage gives a cost savings only when the number of cartridges is considerably larger than the number of drives
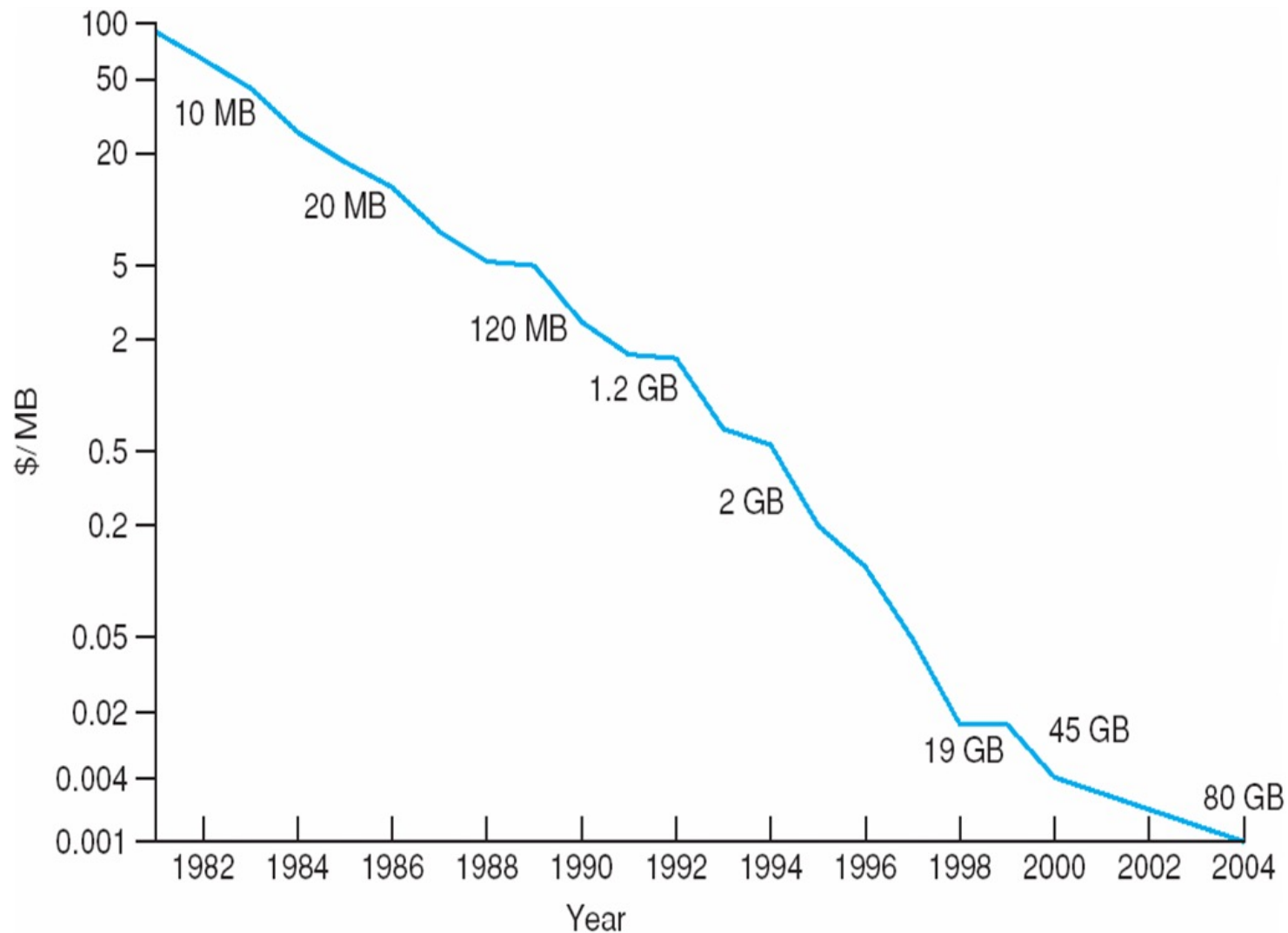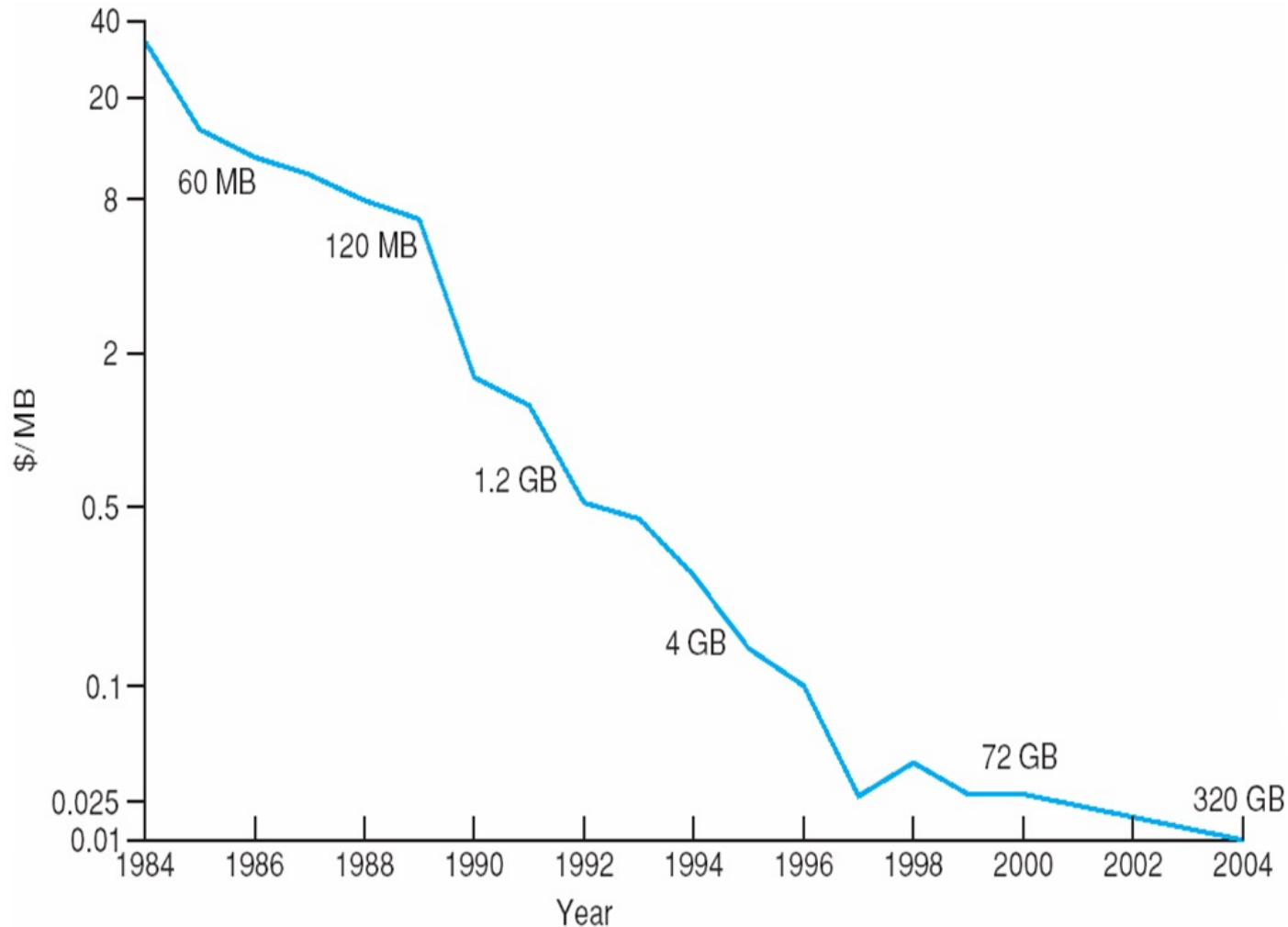
# Price per Megabyte of DRAM, From 1981 to 2004

# Price per Megabyte of a Tape Drive, From 1984-2000

# assignment

- ➢ 12.2
- ➢ 12.3

# End of Chapter 12