

14. Mass-Storage Systems

Sungyoung Lee

*College of Engineering
KyungHee University*

Contents

- n *Disk Structure*
- n *Disk Scheduling*
- n *Disk Management*
- n *Swap-Space Management*
- n *RAID Structure*
- n *Disk Attachment*
- n *Stable-Storage Implementation*
- n *Tertiary Storage Devices*
- n *Operating System Issues*
- n *Performance Issues*

Secondary Storage

n Secondary storage usually

- ü is anything that is outside of “primary memory”
- ü does not permit direct execution of instructions or data retrieval via machine load/store instructions

n Characteristics

- ü It's large: 100GB and more
- ü It's cheap: 120GB IDE disk costs \ 250,000
- ü It's persistent: data survives power loss
- ü It's slow: milliseconds to access

Disk Structure

- n Disk drives are addressed as large 1-dimensional arrays of *logical blocks*, where the logical block is the smallest unit of transfer
- n The 1-dimensional array of logical blocks is mapped into the sectors of the disk sequentially
 - ü Sector 0 is the first sector of the first track on the outermost cylinder
 - ü 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

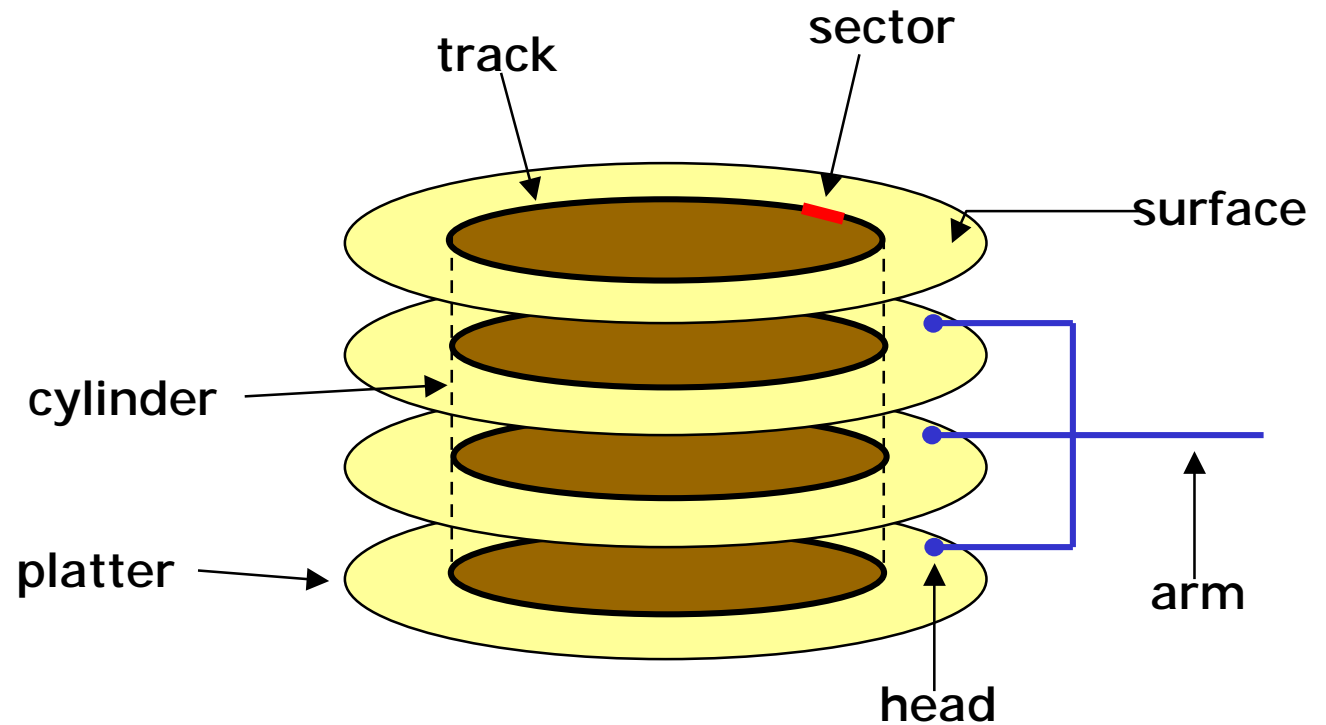
n Disks and the OS

- ü Disks are messy physical devices:
 - § Errors, bad blocks, missed seeks, etc.
- ü The job of the OS is to hide this mess from higher-level software
 - § Low-level device drivers (initiate a disk read, etc)
 - § Higher-level abstractions (files, databases, etc.)
- ü The OS may provide different levels of disk access to different clients
 - § Physical disk block (surface, cylinder, sector)
 - § Disk logical block (disk block #)
 - § Logical file (filename, block or record or byte #)

Disks (Cont'd)

n Physical disk structure

- ü platters
- ü surfaces
- ü tracks
- ü sectors
- ü cylinders
- ü arm
- ü heads



n Interacting with disks

- ü Specifying disk requests requires a lot of info:
 - § Cylinder #, surface #, track #, sector #, transfer size, etc.
- ü Older disks required the OS to specify all of this
 - § The OS needs to know all disk parameters
- ü Modern disks are more complicated
 - § Not all sectors are the same size, sectors are remapped, etc.
- ü Current disks provide a higher-level interface (e.g. SCSI)
 - § The disks exports its data as a logical array of blocks [0..N]
 - § Disk maps logical blocks to cylinder/surface/track/sector
 - § Only need to specify the logical block # to read/write
 - § As a result, physical parameters are hidden from OS

n Disk performance

ü Performance depends on a number of steps

§ **Seek**: moving the disk arm to the correct cylinder

→ depends on how fast disk arm can move (increasing very slowly)

§ **Rotation**: waiting for the sector to rotate under head

→ depends on rotation rate of disk (increasing, but slowly)

§ **Transfer**: transferring data from surface into disk controller, sending it back to the host

→ depends on density of bytes on disk (increasing, and very quickly)

ü Disk scheduling:

§ Because seeks are so expensive, the OS tries to schedule disk requests that are queued waiting for the disk

Disk Scheduling

- n The operating system is responsible for using hardware efficiently
 - ü for the disk drives, this means having a fast access time and disk bandwidth
- n Access time has two major components
 - ü *Seek time* is the time for the disk are to move the heads to the cylinder containing the desired sector
 - ü *Rotational latency* is the additional time waiting for the disk to rotate the desired sector to the disk head
- n Minimize seek time
 - ü Seek time \approx seek distance
- n 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'd)

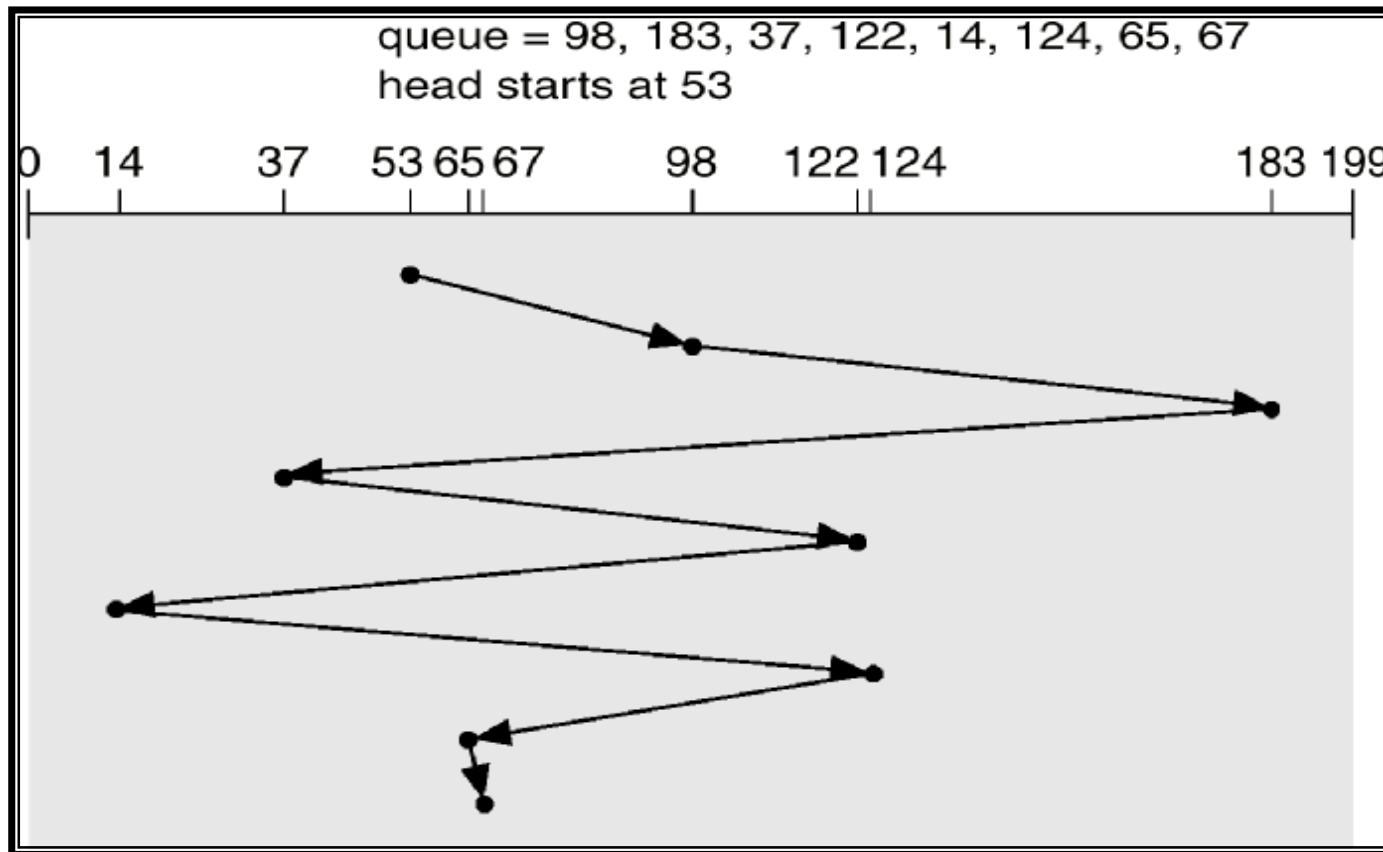
- n Several algorithms exist to schedule the servicing of disk I/O requests
- n We illustrate them with a request queue (0-199)

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

Head pointer 53

FCFS

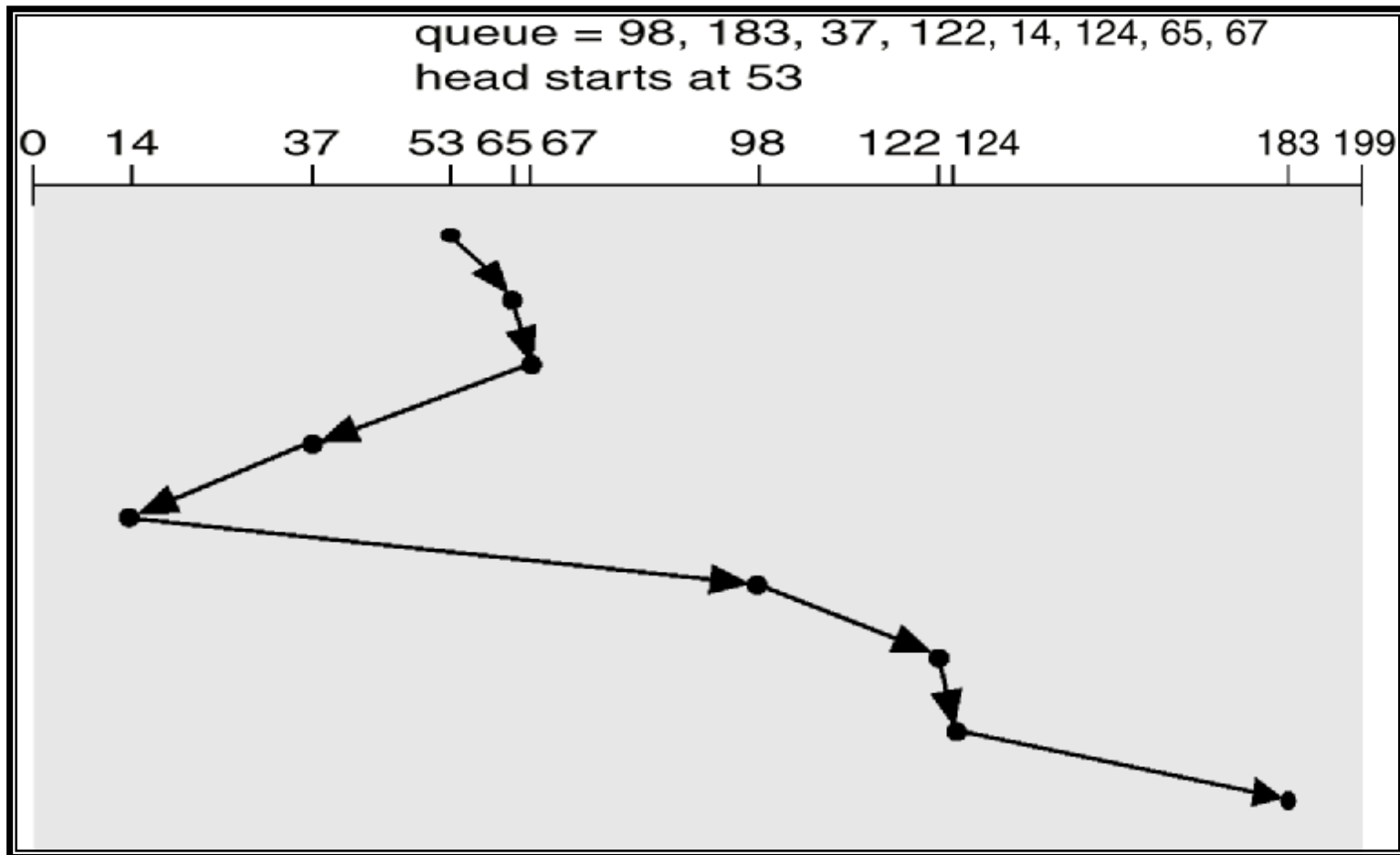
n Illustration shows total head movement of 640 cylinders



SSTF

- n Selects the request with the minimum seek time from the current head position
- n SSTF scheduling is a form of SJF scheduling
 - ü may cause starvation of some requests
- n Illustration shows total head movement of 236 cylinders

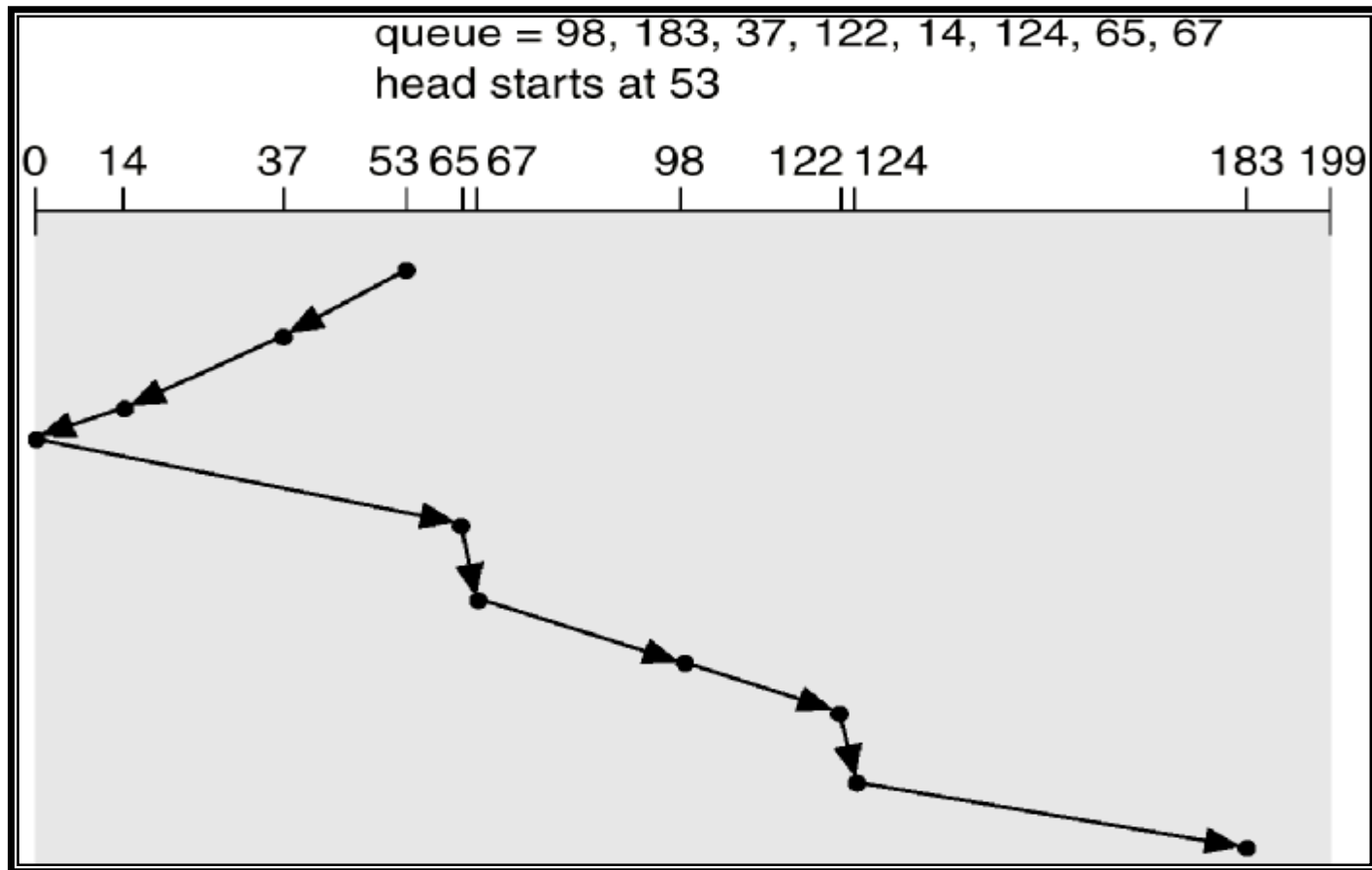
SSTF (Cont'd)



SCAN

- n 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
- n Sometimes called the *elevator algorithm*
- n Illustration shows total head movement of 208 cylinders

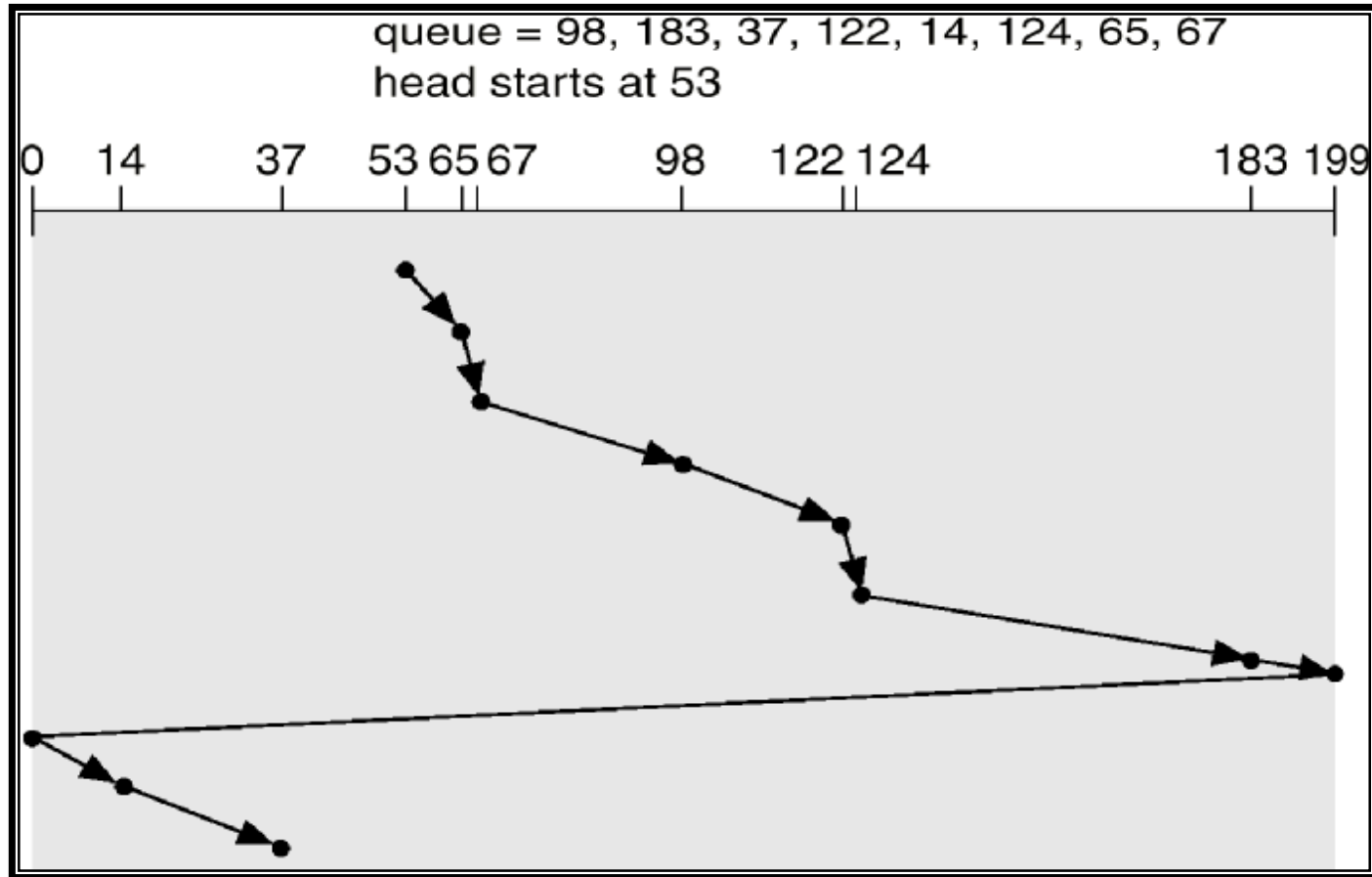
SCAN (Cont'd)



C-SCAN

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

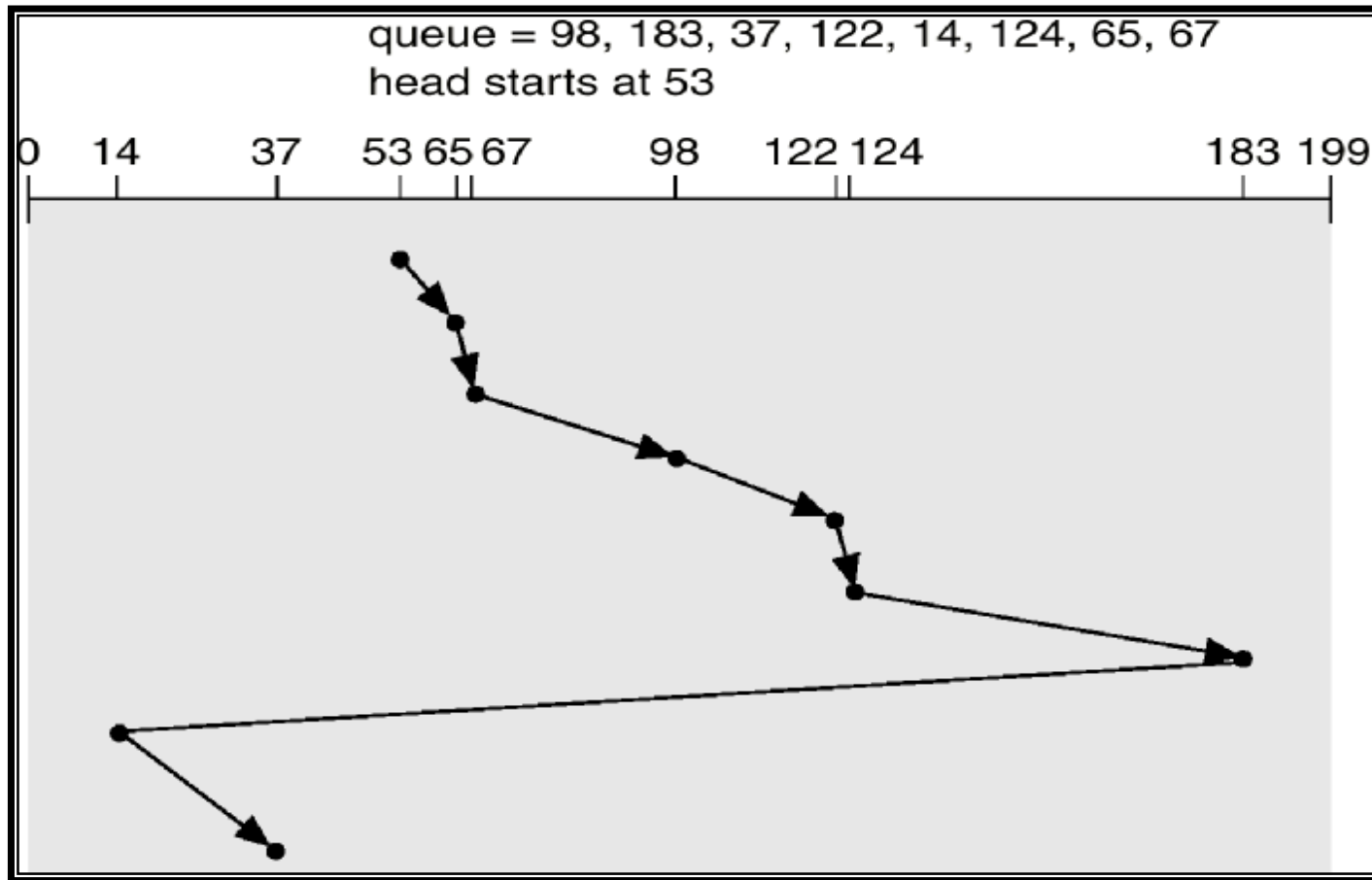
C-SCAN (Cont'd)



C-LOOK

- n Version of C-SCAN
- n 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
- n Cf) Look

C-LOOK (Cont'd)



Selecting a Disk-Scheduling Algorithm

- n SSTF is common and has a natural appeal
- n SCAN and C-SCAN perform better for systems that place a heavy load on the disk
- n Performance depends on the number and types of requests
- n Requests for disk service can be influenced by the file-allocation method
- n 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
- n Either SSTF or LOOK is a reasonable choice for the default algorithm
- n In general, unless there are request queues, disk scheduling does not have much impact
 - ü Important for servers, less so for PCs
- n Modern disks often do the disk scheduling themselves
 - ü Disks know their layout better than OS, can optimize better
 - ü Ignores, undoes any scheduling done by OS

Disk Management

- n *Low-level formatting, or physical formatting*
 - ü Dividing a disk into sectors that the disk controller can read and write

- n To use a disk to hold files, the operating system still needs to record its own data structures on the disk
 - ü *Partition* the disk into one or more groups of cylinders
 - ü *Logical formatting* or “making a file system”

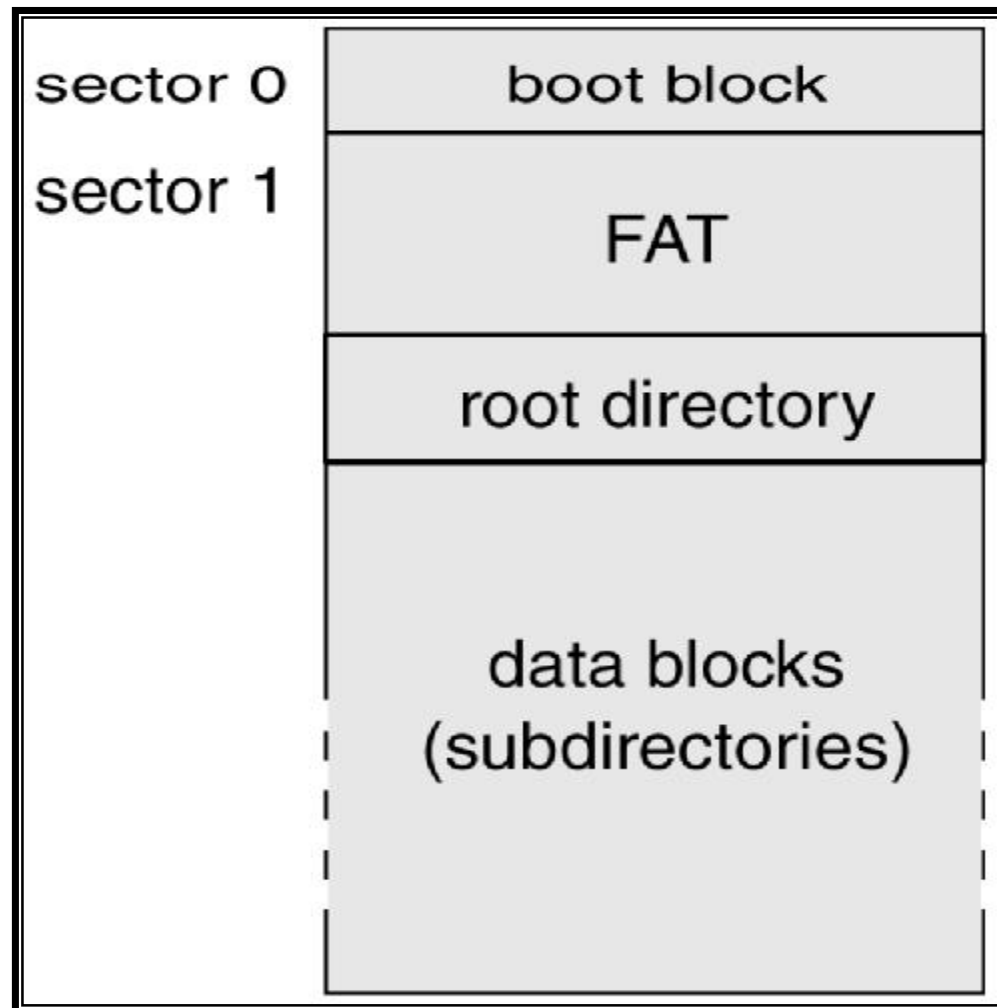
- n **Boot block initializes system**
 - ü The bootstrap is stored in ROM
 - ü *Bootstrap loader* program

- n Methods such as *sector sparing* used to handle bad blocks

n Intelligent controllers

- ü Nowadays, most disk controllers are built around a small CPU and have many kilobytes of memory
- ü They run a program written by the controller manufacturer to process I/O requests from the CPU and satisfy them
- ü Intelligent features:
 - § Read-ahead: the current track
 - § Caching: frequently-used blocks
 - § Request reordering: for seek and/or rotational optimality
 - § Request retry on hardware failure
 - § Bad block identification
 - § Bad block remapping: onto spare blocks and/or tracks

MS-DOS Disk Layout



Swap-Space Management

n Swap-space

- ü Virtual memory uses disk space as an extension of main memory

n Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition

n Swap-space management

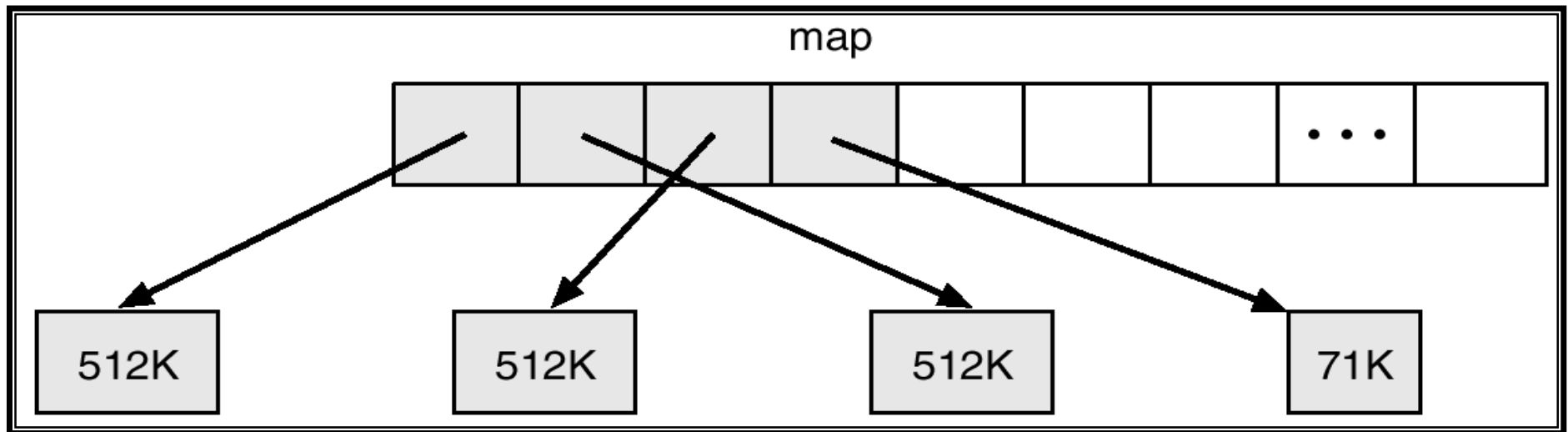
- ü 4.3BSD allocates swap space when process starts

- § holds *text segment* (the program) and *data segment*

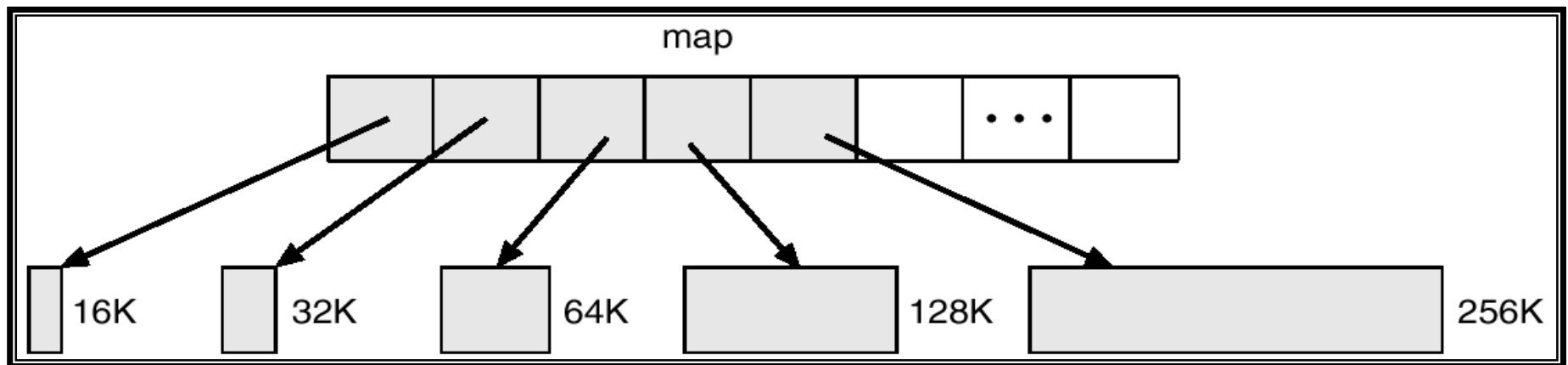
- ü Kernel uses *swap maps* to track swap-space use

- ü Solaris 2 allocates swap space only when a page is forced out of physical memory, not when the virtual memory page is first created

4.3 BSD Text-Segment Swap Map



4.3 BSD Data-Segment Swap Map



RAID Structure

- n **RAID** – multiple disk drives provides **reliability** via **redundancy**
- n RAID is arranged into six different levels
- n Several improvements in disk-use techniques involve the use of multiple disks working cooperatively
- n Disk striping uses a group of disks as one storage unit
- n RAID schemes improve performance and improve the reliability of the storage system by storing redundant data
 - ü *Mirroring* or *shadowing* keeps duplicate of each disk
 - ü *Block interleaved parity* uses much less redundancy

n Redundant Array of Inexpensive Disks

- ü A storage system, not a file system

n Motivations

- ü Use a small cheap disks as a cost-effective alternative to large, expensive disks (I = Inexpensive)
- ü Provide higher reliability and higher data-transfer (I = Independent)

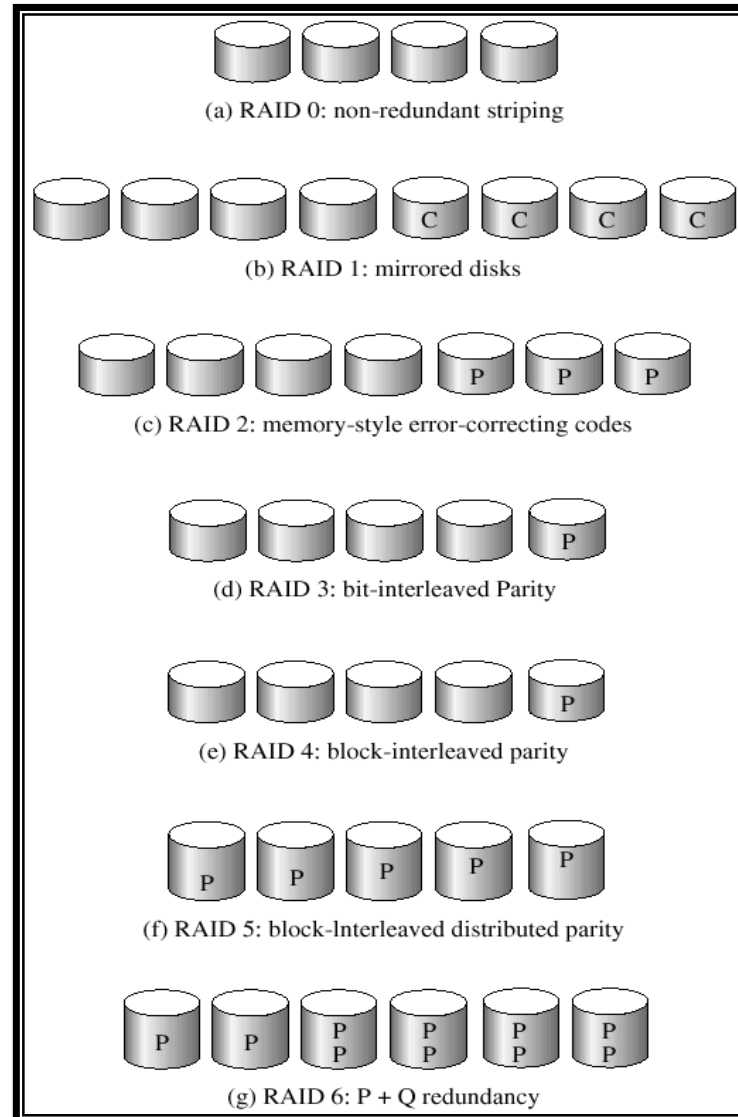
n Improving reliability via redundancy

- ü Mirroring (shadowing)
- ü Parity or error-correcting codes

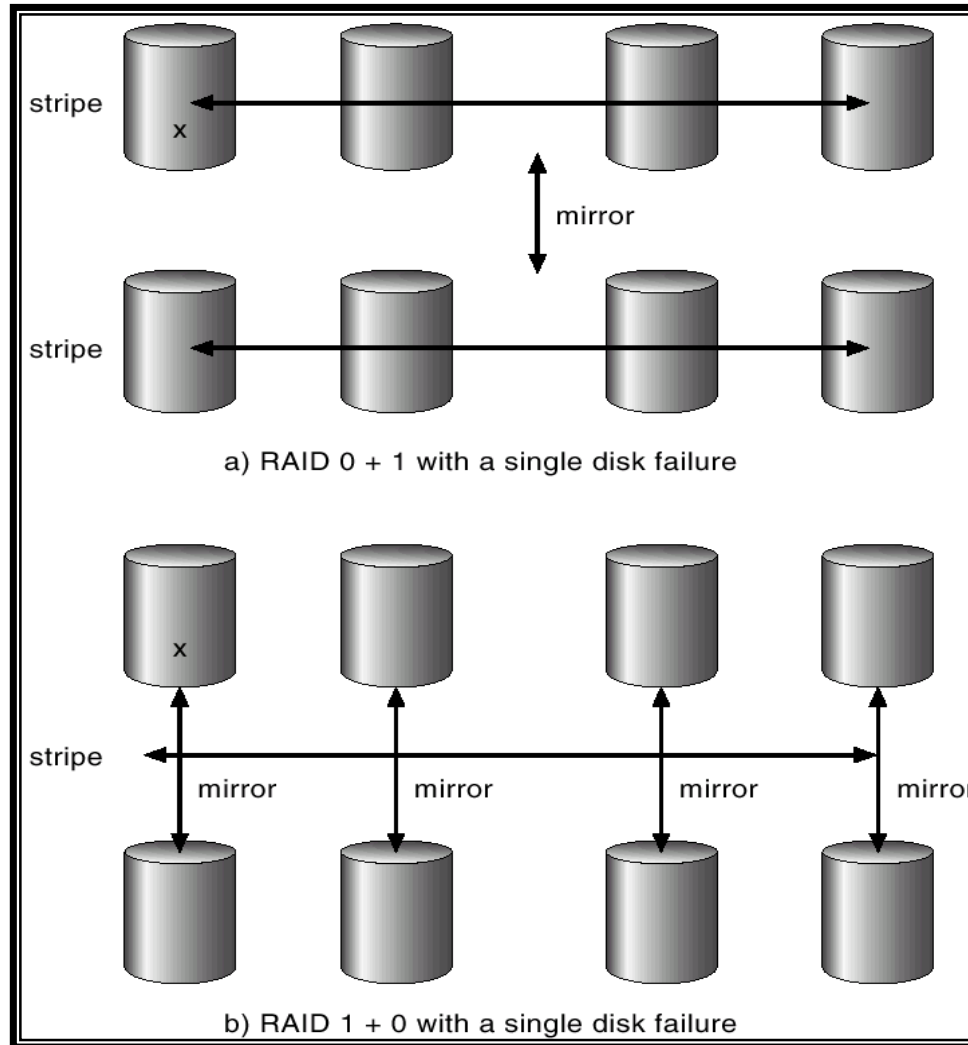
n Improving performance via parallelism

- ü Data striping: bit-level vs. block-level

RAID Levels



RAID (0 + 1) and (1 + 0)



n RAID 0

- ü Non-redundant striping (no reliability)
- ü Data is broken down into blocks and each block is striped across multiple disks
- ü I/O performance is greatly improved by spreading the I/O load across many channels and drives
- ü Typically used in data rate intensive applications (e.g. video editing)

0	1	2	3
A0	A1	A2	A3
A4	A5	A6	A7
A8	A9	A10	A11
...

RAID Levels (Cont'd)

n RAID 1

- ü Mirrored disks
- ü Twice the read transaction rate of single disks, same write transaction rate as single disks
- ü Expensive, highest disk overhead
- ü No rebuild is necessary in case of a disk failure

0	1	2	3	4	5	6	7
A0	A10	A20	A30	A0	A10	A20	A30
A1	A11	A21	A31	A1	A11	A21	A31
A2	A12	A22	A32	A2	A12	A22	A32
...

RAID Levels (Cont'd)

n RAID 2

- ü Memory-style error-correcting codes (ECC)
- ü Each data word has its Hamming Code ECC word recorded on the ECC disks
- ü On read, the ECC code verifies correct data or corrects single bit errors
- ü ECC is embedded in almost all modern disk drives (e.g. SCSI)

0	1	2	3	4	5	6
A0	A0	A0	A0	Px	Py	Pz
A1	A1	A1	A1	Px	Py	Pz
A2	A2	A2	A2	Px	Py	Pz
A3	A3	A3	A3	Px	Py	Pz
...

RAID Levels (Cont'd)

n RAID 3

- ü Bit-interleaved parity
- ü Stripe parity is generated on writes, recorded on the parity disk and checked on reads
- ü Less storage overhead than RAID 2
- ü Requires hardware support for efficient use

0	1	2	3	4
A0	A0	A0	A0	P
A1	A1	A1	A1	P
A2	A2	A2	A2	P
A3	A3	A3	A3	P
...

RAID Levels (Cont'd)

n RAID 4

- ü Block-interleaved parity
- ü Parity for same rank blocks is generated on writes, recorded on the parity disk and checked on reads
- ü Very good read performance (same as RAID 0)
- ü Writes, however, require parity data be updated each time
- ü No advantages over RAID 5 and does not support multiple simultaneous write operations

0	1	2	3	4
A0	A1	A2	A3	P
A4	A5	A6	A7	P
A8	A9	A10	A11	P
...

RAID Levels (Cont'd)

n RAID 5

- ü Block-interleaved distributed parity
- ü Parity for blocks in the same rank is generated on writes, recorded in a distributed location
- ü Can speed up small writes in multiprocessing systems, since the parity disk does not become a bottleneck

0	1	2	3	4
A0	A1	A2	A3	P
A4	A5	A6	P	A7
A8	A9	P	A10	A11
A12	P	A13	A14	A15
...

RAID Levels (Cont'd)

n RAID 6

- ü P+Q redundancy scheme
- ü RAID 5 + extra redundancy information to guard against multiple disk failure
- ü Use error-correcting codes instead of parity

0	1	2	3	4	5
A0	A1	A2	A3	Px	Py
A4	A5	A6	Px	Py	A7
A8	A9	Px	Py	A10	A11
A12	Px	Py	A13	A14	A15
...

RAID Levels (Cont'd)

n RAID 0+1

- ü A set of disks are striped, and then the stripe is mirrored to another, equivalent stripe
- ü RAID 0 provides performance, while RAID 1 provides the reliability
- ü A single drive failure will cause the whole array to become, in essence, a RAID 0 array

0	1	2	3	4	5	6	7
A0	A1	A2	A3	A0	A1	A2	A3
A4	A5	A6	A7	A4	A5	A6	A7
A8	A9	A10	A11	A8	A9	A10	A11
...

RAID Levels (Cont'd)

n RAID 10 (or RAID 1+0)

ü Disks are mirrored in pairs, and then the resulting mirror pairs are striped

ü Reliable better than RAID 0+1

§ RAID 0+1 is fault tolerant as long as the second through n-th disk is on the same stripe

§ RAID 1+0 is fault tolerant as long as no two disks are part of the same mirror

0	1	2	3	4	5	6	7
A0	A1	A2	A3	A4	A5	A6	A7
A8	A9	A10	A11	A12	A13	A14	A15
...
A4	A5	A6	A7	A0	A1	A2	A3
A12	A13	A14	A15	A8	A9	A10	A11
...

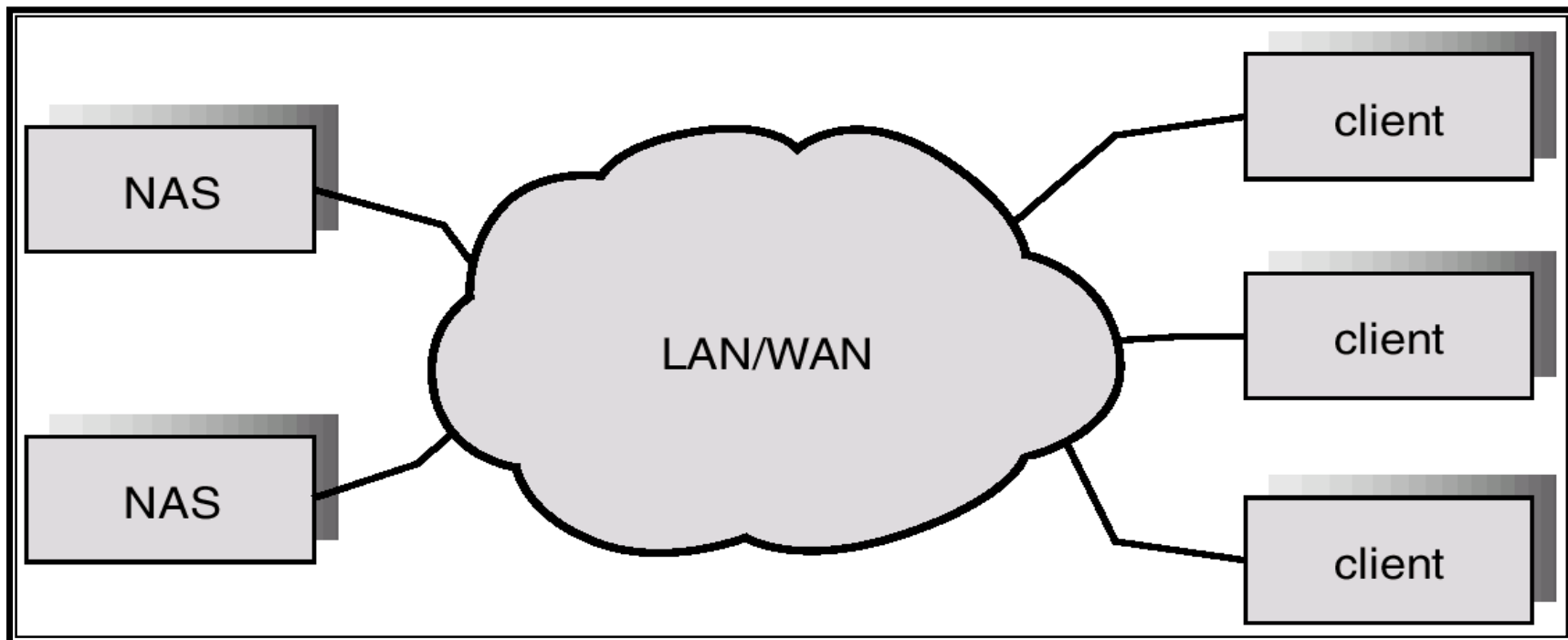
Disk Attachment

n Disks may be attached one of two ways:

1. **Host attached** via an I/O port
2. **Network attached** via a network connection

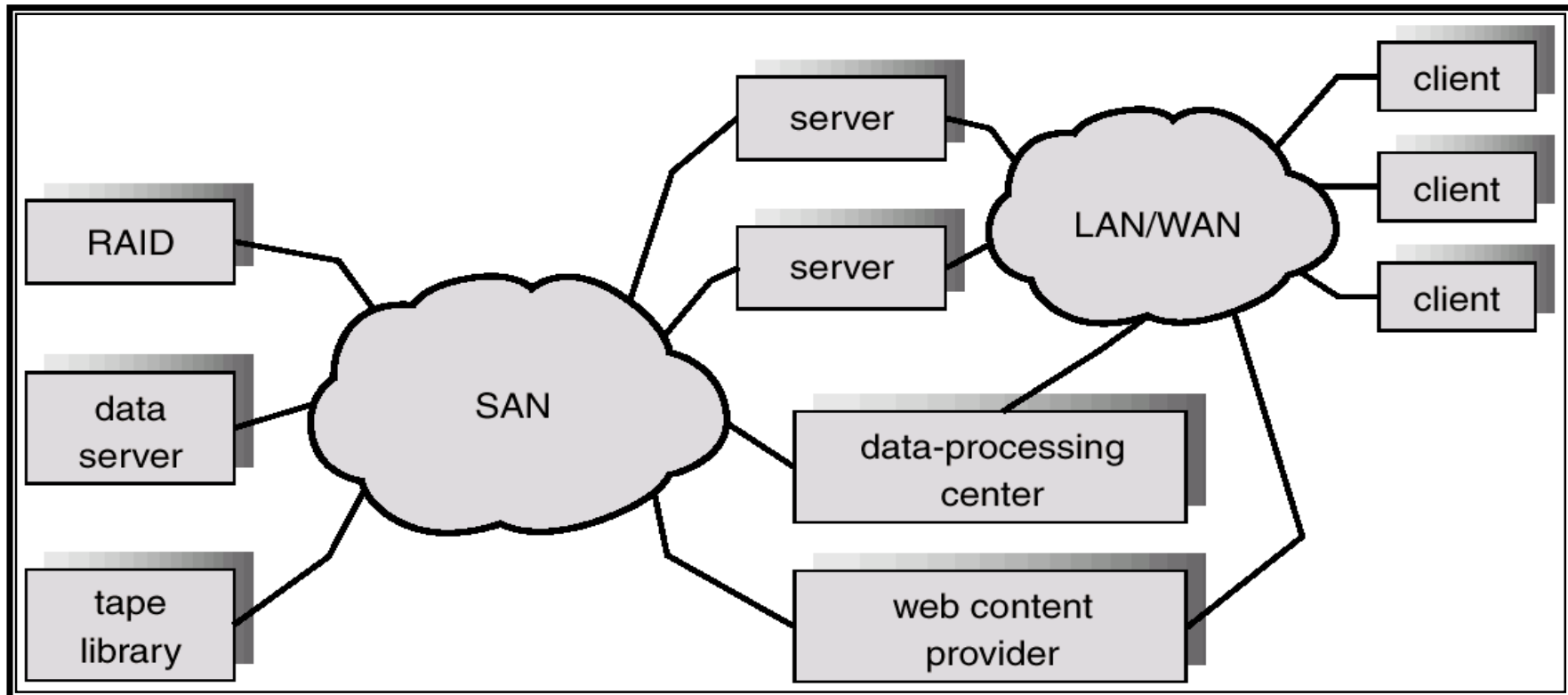
Network-Attached Storage

- n Accessed via TCP/IP- or UDP/IP-based networks
- n NFS, CIFS, etc.
- n File-level access



Storage-Area Network

- n Accessed via private network dedicated for storages
- n Use storage protocols such as SCSI or Fibre channel
- n Block-level access
- n File systems for SAN is another story (e.g. GFS)



Storage Architecture

Protocols

		Block	File
Interconnection	Non-IP Networks	IDE / SCSI (Direct) SAN (Storage Area Network: Fibre Channel)	DAFS (Direct Access File System: NFS over VIA)
	IP-based Networks	NBD (Network Block Device) iSCSI (SCSI over TCP/IP)	NAS (Network Attached Storage: NFS, CIFS)

Stable-Storage Implementation

- n Write-ahead log scheme requires stable storage
- n To implement stable storage:
 - ü Replicate information on more than one nonvolatile storage media with independent failure modes
 - ü Update information in a controlled manner to ensure that we can recover the stable data after any failure during data transfer or recovery

Tertiary Storage Devices

- n Low cost is the defining characteristic of tertiary storage
- n Generally, tertiary storage is built using *removable media*
- n Common examples of removable media are floppy disks and CD-ROMs
 - ü other types are available

Removable Disks

n Floppy disk

- ü Thin flexible disk coated with magnetic material, enclosed in a protective plastic case
- ü Most floppies hold about 1 MB
- ü Similar technology is used for removable disks that hold more than 1 GB
- ü 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'd)

- n A magneto-optic disk records data on a rigid platter coated with magnetic material
 - ü Laser heat is used to amplify a large, weak magnetic field to record a bit
 - ü Laser light is also used to read data (Kerr effect)
 - ü 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
 - ü Resistant to head crashes

- n Optical disks do not use magnetism
 - ü They employ special materials that are altered by laser light

WORM Disks

- n The data on read-write disks can be modified over and over
- n WORM (“Write Once, Read Many Times”) disks can be written only once
- n Thin aluminum film sandwiched between two glass or plastic platters
- n To write a bit, the drive uses a laser light to burn a small hole through the aluminum
 - ü Information can be destroyed by not altered
- n Very durable and reliable
- n *Read Only* disks, such as CD-ROM and DVD, come from the factory with the data pre-recorded

Tapes

- n Compared to a disk, a tape is less expensive and holds more data, but random access is much slower
- n Tape is an economical medium for purposes that do not require fast random access
 - ü e.g., backup copies of disk data, holding huge volumes of data
- n Large tape installations typically use robotic tape changers that move tapes between tape drives and storage slots in a tape library
 - ü stacker – library that holds a few tapes
 - ü silo – library that holds thousands of tapes
- n A disk-resident file can be *archived* to tape for low cost storage
 - ü The computer can *stage* it back into disk storage for active use

Operating System Issues

- n Major OS jobs are to manage physical devices and to present a virtual machine abstraction to applications
- n For hard disks, the OS provides two abstraction:
 - ü Raw device – an array of data blocks
 - ü File system – the OS queues and schedules the interleaved requests from several applications

Application Interface

- n Most OSs handle removable disks almost exactly like fixed disks
 - ü a new cartridge is formatted and an empty file system is generated on the disk
- n Tapes are presented as a raw storage medium
 - ü i.e., application does not open a file on the tape, it opens the whole tape drive as a raw device
- n Usually the tape drive is reserved for the exclusive use of that application
- n Since the OS does not provide file system services, the application must decide how to use the array of blocks
- n 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

- n The basic operations for a tape drive differ from those of a disk drive
- n **locate** positions the tape to a specific logical block, not an entire track (corresponds to **seek**)
- n The **read position** operation returns the logical block number where the tape head is
- n The **space** operation enables relative motion
- n Tape drives are “append-only” devices
 - ü Updating a block in the middle of the tape also effectively erases everything beyond that block
- n An EOT mark is placed after a block that is written

File Naming

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

Hierarchical Storage Management (HSM)

- n A hierarchical storage system extends the storage hierarchy beyond primary memory and secondary storage to incorporate tertiary storage
 - ü Usually implemented as a jukebox of tapes or removable disks
- n Usually incorporate tertiary storage by extending the file system
 - ü Small and frequently used files remain on disk
 - ü Large, old, inactive files are archived to the jukebox
- n HSM is usually found in supercomputing centers and other large installations that have enormous volumes of data

Speed

- n Two aspects of speed in tertiary storage are bandwidth and latency
- n Bandwidth is measured in bytes per second
 - ü Sustained bandwidth
 - § average data rate during a large transfer
 - § # of bytes/transfer time
 - § Data rate when the data stream is actually flowing
 - ü Effective bandwidth
 - § average over the entire I/O time, including **seek** or **locate**, and cartridge switching
 - § Drive's overall data rate

Speed (Cont'd)

n Access latency

- ü Amount of time needed to locate data

- ü Access time for a disk

 - § move the arm to the selected cylinder and wait for the rotational latency

 - § < 35 milliseconds

- ü Access on tape requires winding the tape reels until the selected block reaches the tape head

 - § tens or hundreds of seconds

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

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

n 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

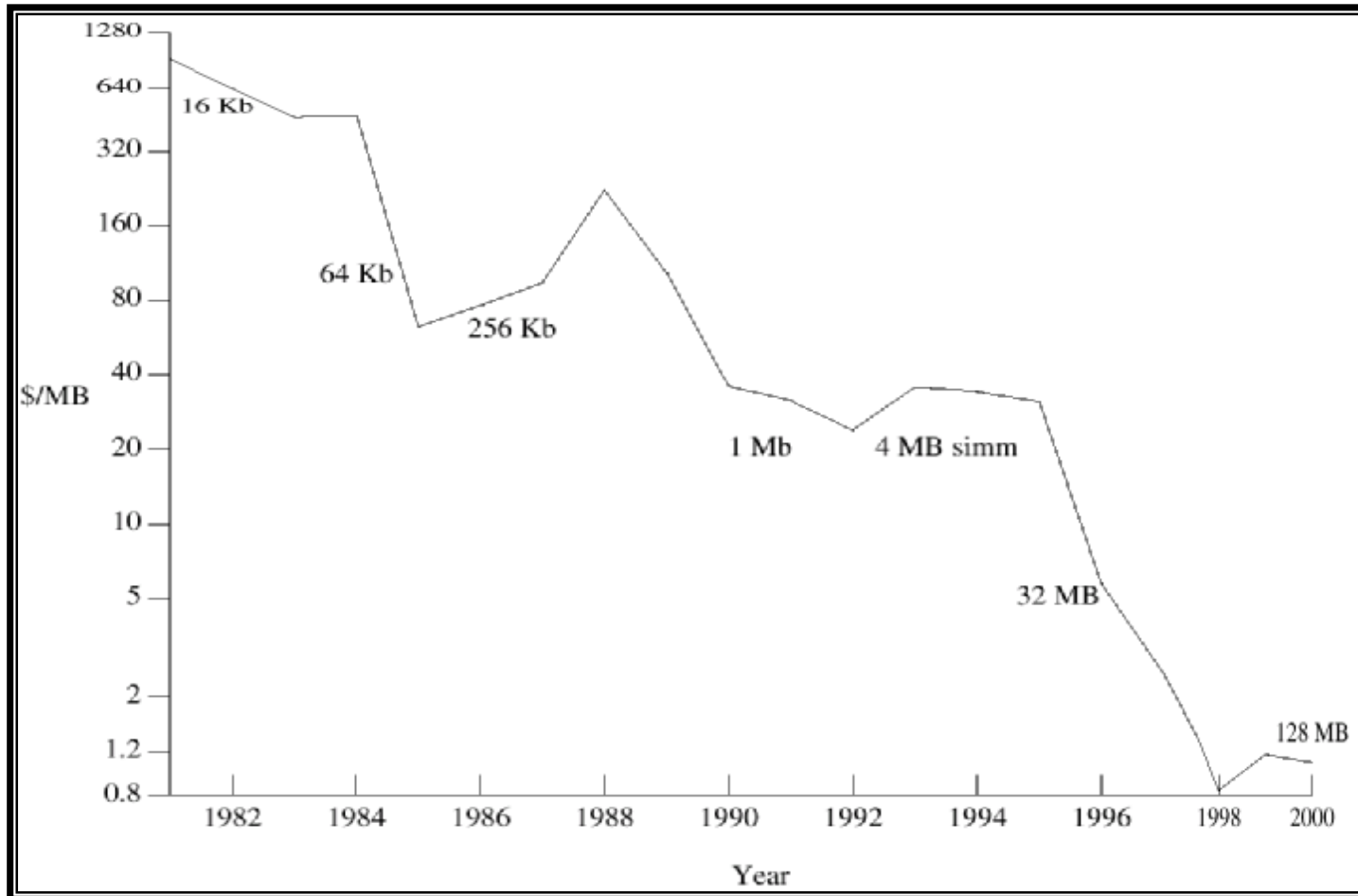
Reliability

- n A fixed disk drive is likely to be more reliable than a removable disk or tape drive
- n An optical cartridge is likely to be more reliable than a magnetic disk or tape
- n 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

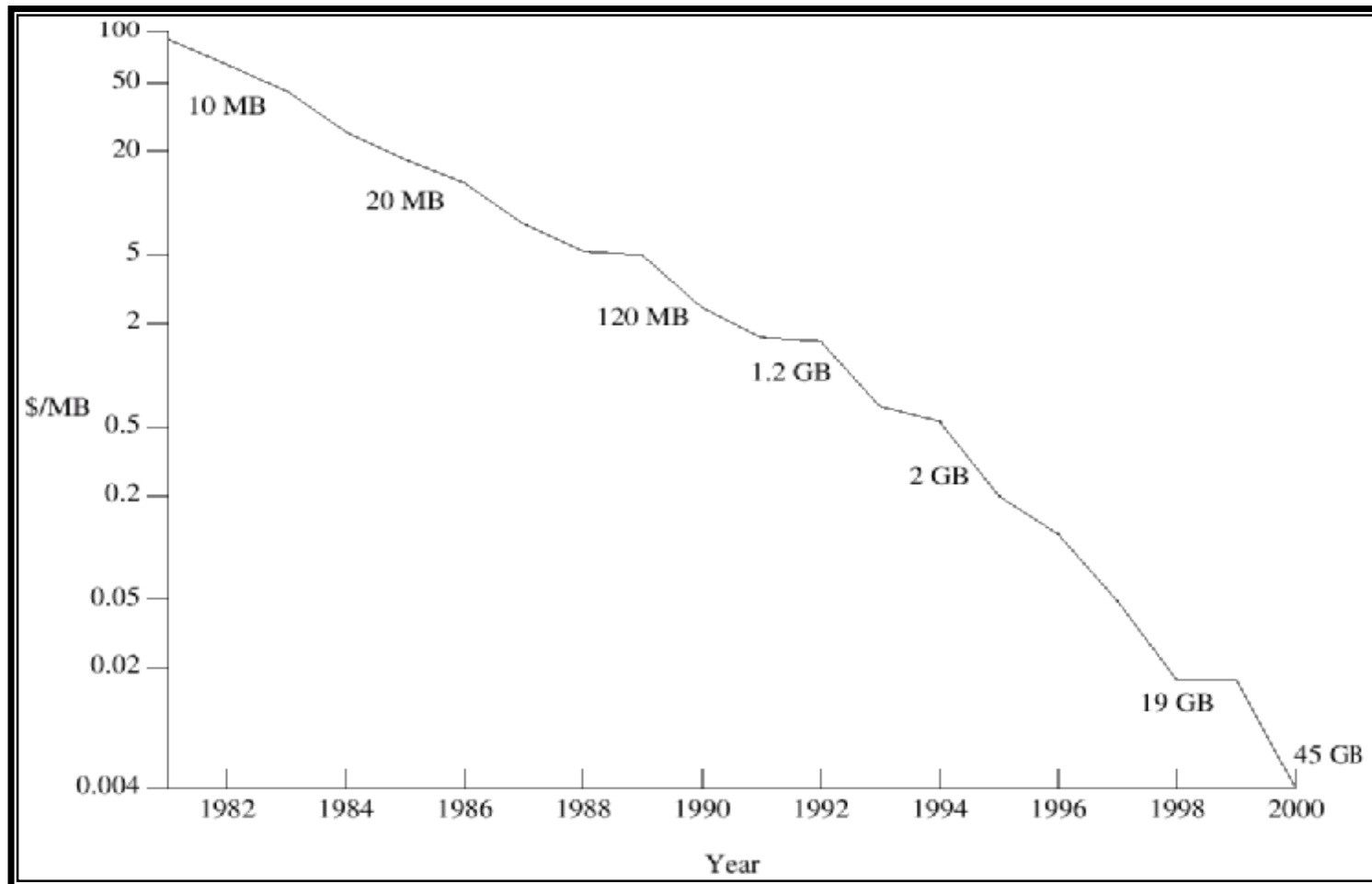
Cost

- n Main memory is much more expensive than disk storage
- n The cost per megabyte of hard disk storage is competitive with magnetic tape if only one tape is used per drive
- n The cheapest tape drives and the cheapest disk drives have had about the same storage capacity over the years
- n 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 2000



Price per Megabyte of Magnetic Hard Disk, From 1981 to 2000



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

