

## Chapter 12: Mass-Storage Systems

Operating System Concepts - 8<sup>th</sup> Edition

Silberschatz, Galvin and Gagne ©2009

12.1

Silberschatz, Galvin and Gagne ©2009

### 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

Operating System Concepts - 8<sup>th</sup> Edition

12.3

Silberschatz, Galvin and Gagne ©2009

12.3

Silberschatz, Galvin and Gagne ©2009

### Overview of Mass Storage Structure

- Magnetic disks provide bulk of secondary storage of modern computers
  - Drives rotate at 60 to 200 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**
  - **Host controller** in computer uses bus to talk to **disk controller** built into drive or storage array

Operating System Concepts - 8<sup>th</sup> Edition

12.4

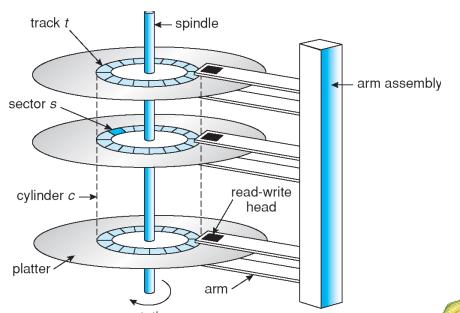
Silberschatz, Galvin and Gagne ©2009

Operating System Concepts - 8<sup>th</sup> Edition

12.5

Silberschatz, Galvin and Gagne ©2009

### Moving-head Disk Mechanism

Operating System Concepts - 8<sup>th</sup> Edition

12.5

Silberschatz, Galvin and Gagne ©2009

Operating System Concepts - 8<sup>th</sup> Edition

12.6

Silberschatz, Galvin and Gagne ©2009

### Overview of Mass Storage Structure (Cont.)

- 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 Structure



**Operating System Concepts - 8<sup>th</sup> Edition** 12.7 Silberschatz, Galvin and Gagne ©2009

- 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.



## Disks

### Disk Hardware (1)



**Operating System** 12.8 Silberschatz, Galvin and Gagne ©2009

Parameter	IBM 360-KB floppy disk	WD 18300 hard disk
Number of cylinders	40	10601
Tracks per cylinder	2	12
Sectors per track	9	281 (avg)
Sectors per disk	720	35742000
Bytes per sector	512	512
Disk capacity	360 KB	18.3 GB
Seek time (adjacent cylinders)	6 msec	0.8 msec
Seek time (average case)	77 msec	6.9 msec
Rotation time	200 msec	8.33 msec
Motor stop/start time	250 msec	20 sec
Time to transfer 1 sector	22 msec	17 $\mu$ sec

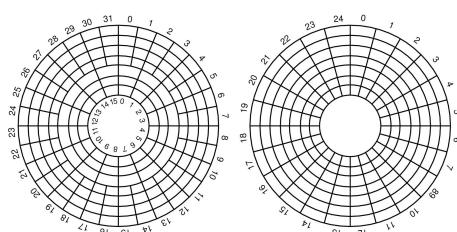
Disk parameters for the original IBM PC floppy disk and a Western Digital WD 18300 hard disk



## Disk Hardware (2)



**Operating System** 12.9 Silberschatz, Galvin and Gagne ©2009



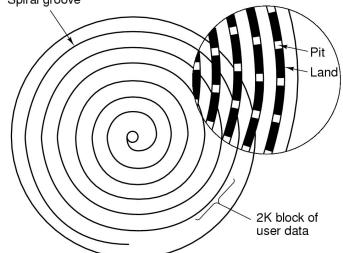
- Physical geometry of a disk with two zones
- A possible virtual geometry for this disk



## Disk Hardware (5)



**Operating System** 12.10 Silberschatz, Galvin and Gagne ©2009



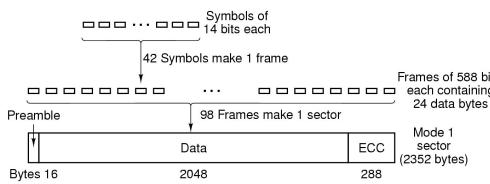
Recording structure of a CD or CD-ROM



## Disk Hardware (6)



**Operating System** 12.11 Silberschatz, Galvin and Gagne ©2009



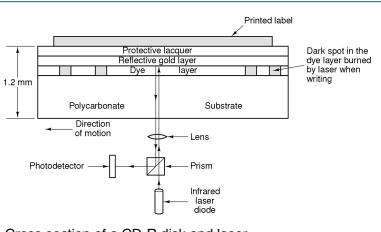
Logical data layout on a CD-ROM



## Disk Hardware (7)



**Operating System** 12.12 Silberschatz, Galvin and Gagne ©2009



- Cross section of a CD-R disk and laser
  - not to scale
- Silver CD-ROM has similar structure
  - without dye layer
  - with pitted aluminum layer instead of gold



**Disk Hardware (8)**

A double sided, dual layer DVD disk

Operating System 12.13 Silberschatz, Galvin and Gagne ©2009

**Question**

- A disk manufacturer has two 5.25 inch disks that each have 10,000 cylinders. The newer one has double the linear recording density of the older one.
- Which disks properties are better on the newer one and which are the same?

Operating System 12.14 Silberschatz, Galvin and Gagne ©2009

**Answer**

- Drive capacity and transfer rate are doubled.
- Seek time and average rotational delay are the same.

Operating System 12.15 Silberschatz, Galvin and Gagne ©2009

**Disk Management**

- *Low-level formatting*, or *physical formatting* — Dividing a disk into sectors that the disk controller can read and write.
- 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".
- Boot block initializes system.
  - The bootstrap is stored in ROM.
  - *Bootstrap loader* program.
- Methods such as *sector sparing* used to handle bad blocks.

Operating System 12.16 Silberschatz, Galvin and Gagne ©2009

**Error Handling**

- A disk track with a bad sector
- Substituting a spare for the bad sector
- Shifting all the sectors to bypass the bad one

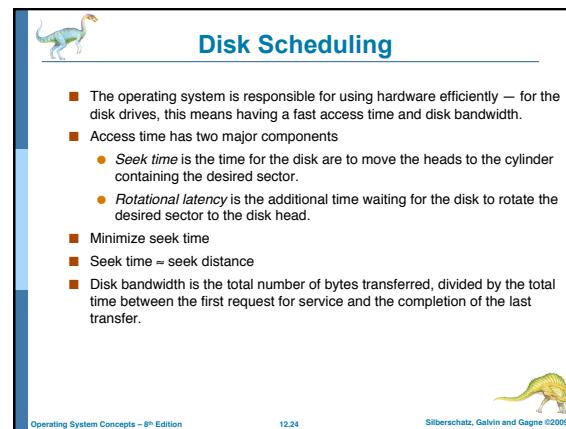
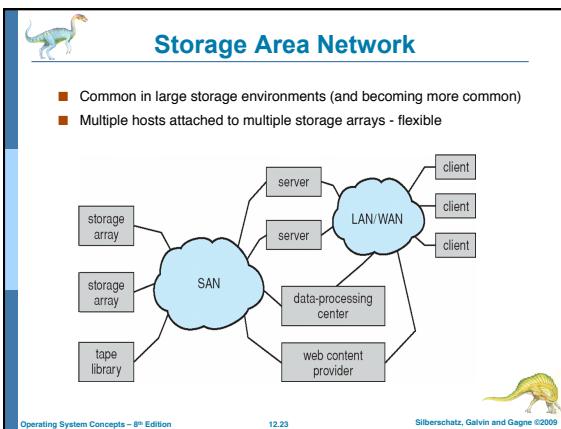
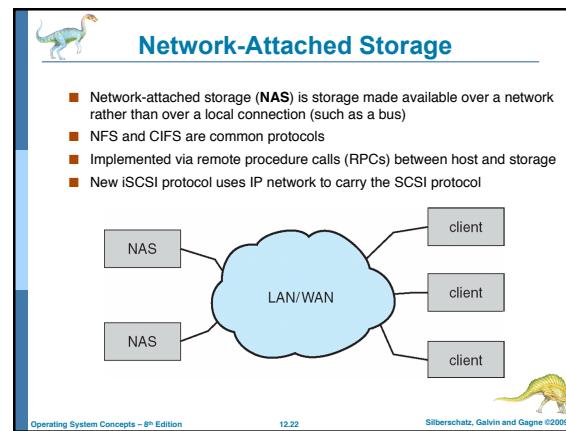
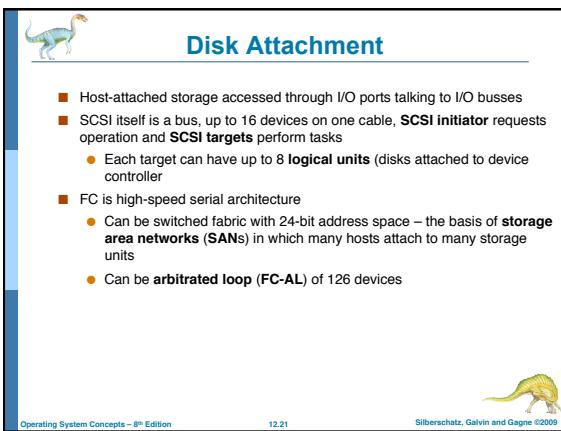
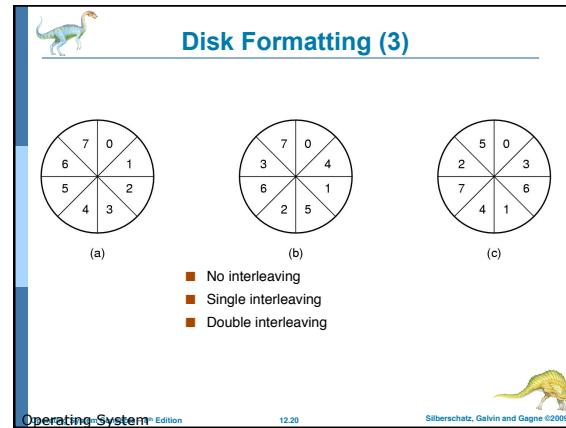
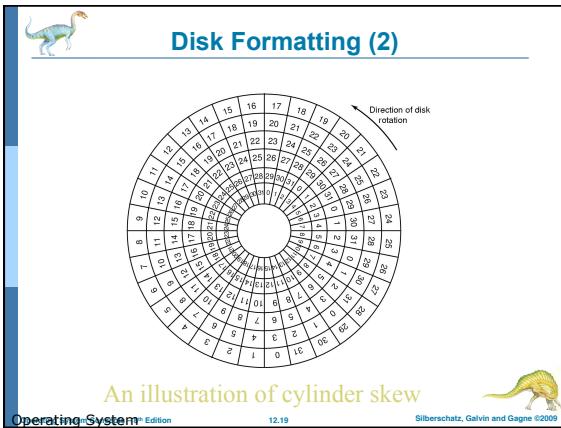
Operating System 12.17 Silberschatz, Galvin and Gagne ©2009

**Disk Formatting (1)**

Preamble	Data	ECC
----------	------	-----

A disk sector

Operating System 12.18 Silberschatz, Galvin and Gagne ©2009



## 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

Operating System Concepts – 8<sup>th</sup> Edition 12.25 Silberschatz, Galvin and Gagne ©2009

## FCFS

Illustration shows total head movement of 640 cylinders.

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

0 14 37 53 65 67 98 122 124 183 199

Operating System Concepts – 8<sup>th</sup> Edition 12.26 Silberschatz, Galvin and Gagne ©2009

## SSTF

- 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.
- Illustration shows total head movement of 236 cylinders.

Operating System Concepts – 8<sup>th</sup> Edition 12.27 Silberschatz, Galvin and Gagne ©2009

## SSTF (Cont.)

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

0 14 37 53 65 67 98 122 124 183 199

Operating System Concepts – 8<sup>th</sup> Edition 12.28 Silberschatz, Galvin and Gagne ©2009

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

Operating System Concepts – 8<sup>th</sup> Edition 12.29 Silberschatz, Galvin and Gagne ©2009

## SCAN (Cont.)

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

0 14 37 53 65 67 98 122 124 183 199

Operating System Concepts – 8<sup>th</sup> Edition 12.30 Silberschatz, Galvin and Gagne ©2009



## C-SCAN

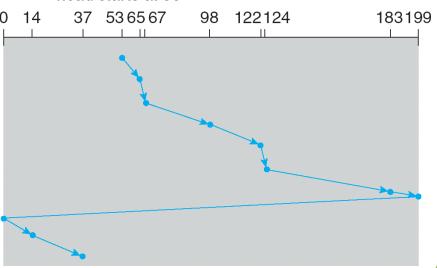
- Provides a more uniform wait time than SCAN.
- The head moves from one end of the disk to the other, servicing requests as it goes. When it reaches the other end, however, it immediately returns to the beginning of the disk, without servicing any requests on the return trip.
- Treats the cylinders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

Operating System Concepts – 8<sup>th</sup> Edition 12.31 Silberschatz, Galvin and Gagne ©2009



## C-SCAN (Cont.)

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



Operating System Concepts – 8<sup>th</sup> Edition 12.32 Silberschatz, Galvin and Gagne ©2009



## C-LOOK

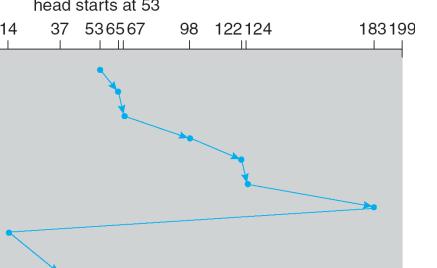
- Version of C-SCAN
- Arm only goes as far as the last request in each direction, then reverses direction immediately, without first going all the way to the end of the disk.

Operating System Concepts – 8<sup>th</sup> Edition 12.33 Silberschatz, Galvin and Gagne ©2009



## C-LOOK (Cont.)

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



Operating System Concepts – 8<sup>th</sup> Edition 12.34 Silberschatz, Galvin and Gagne ©2009



## Question

- Suppose that a disk drive has 5,000 cylinders, numbered 0 to 4999. The drive is currently serving a request at cylinder 143 and the previous request was at cylinder 125. The queue of pending requests, in FIFO order, is 86, 1470, 913, 1774, 948, 1509, 1022, 1750, 130
- Starting from the current head position, what is the total distance that the disk arm moves to satisfy all the pending requests for each of the following scheduling algorithms:

  - FCFS
  - SSTF
  - SCAN
  - LOOK
  - C-SCAN

Operating System Concepts – 8<sup>th</sup> Edition 12.35 Silberschatz, Galvin and Gagne ©2009



## Answer

- FCFS order is: 143, 86, 1470, 913, 1774, 948, 1509, 1022, 1750, 130. Total distance is 7081.
- SSTF order is: 143, 130, 86, 913, 948, 1022, 1470, 1509, 1750, 1744. Total distance is 1745.
- SCAN order is: 143, 913, 948, 1022, 1470, 1509, 1750, 1774, 4999, 130, 86. Total distance is 9769.

Operating System Concepts – 8<sup>th</sup> Edition 12.36 Silberschatz, Galvin and Gagne ©2009

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

Operating System Concepts – 8<sup>th</sup> Edition 12.37 Silberschatz, Galvin and Gagne ©2009

## Disk Management



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

**Boot block** initializes system.

- The bootstrap is stored in ROM.
- Bootstrap loader program.

**Methods such as sector sparing** used to handle bad blocks.

Operating System Concepts – 8<sup>th</sup> Edition 12.38 Silberschatz, Galvin and Gagne ©2009

## Booting from a Disk in Windows 2000

The diagram illustrates the Windows 2000 boot partition structure. It shows four partitions labeled partition 1, partition 2, partition 3, and partition 4. The MBR (Master Boot Record) is located at the beginning of the disk. A dashed line connects the MBR to a detailed view of the boot partition. This view shows the boot code and the partition table. The partition table contains entries for the four partitions, with partition 4 being the boot partition.

Operating System Concepts – 8<sup>th</sup> Edition 12.39 Silberschatz, Galvin and Gagne ©2009

## Swap-Space Management



**Swap-space** — Virtual memory uses disk space as an extension of main memory.  
Swap-space can be carved out of the normal file system, or, more commonly, it can be in a separate disk partition.  
**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.

Operating System Concepts – 8<sup>th</sup> Edition 12.40 Silberschatz, Galvin and Gagne ©2009

## Data Structures for Swapping on Linux Systems

The diagram shows the data structures for swapping on Linux systems. At the top, a horizontal bar represents the swap area, divided into slots. Below this is the swap partition or swap file, which is divided into pages. Arrows point from the swap area slots down to the swap partition, indicating which pages are currently swapped out. At the bottom is the swap map, which is a table of bits (0 or 1) corresponding to the slots in the swap area. The bit values 1, 0, 3, 0, 1 are shown under the first five slots respectively.

Operating System Concepts – 8<sup>th</sup> Edition 12.41 Silberschatz, Galvin and Gagne ©2009

## RAID Structure



**RAID** — multiple disk drives provides **reliability** via **redundancy**.  
RAID is arranged into six different levels.

Operating System Concepts – 8<sup>th</sup> Edition 12.42 Silberschatz, Galvin and Gagne ©2009

## RAID (cont)

- Several improvements in disk-use techniques involve the use of multiple disks working cooperatively.
- 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* keeps duplicate of each disk.
  - Block interleaved parity* uses much less redundancy.

Operating System Concepts – 8<sup>th</sup> Edition 12.43 Silberschatz, Galvin and Gagne ©2009

## 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.  
 (e) RAID 4: block-interleaved parity.  
 (f) RAID 5: block-interleaved distributed parity.  
 (g) RAID 6: P + Q redundancy.

Operating System Concepts – 8<sup>th</sup> Edition 12.44 Silberschatz, Galvin and Gagne ©2009

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

a) RAID 0 + 1 with a single disk failure.  
 b) RAID 1 + 0 with a single disk failure.

Operating System Concepts – 8<sup>th</sup> Edition 12.45 Silberschatz, Galvin and Gagne ©2009

## Question

A RAID can fail if two or more of its drives crash within a short time interval. Suppose the probability of one drive crashing in a given hour is  $p$ . What is the probability of a  $k$ -drive RAID failing in a given hour?

Operating System Concepts – 8<sup>th</sup> Edition 12.46 Silberschatz, Galvin and Gagne ©2009

## Answer

- Probability of zero failures is  $(1-p)^k$
- Probability of one failure is  $kp(1-p)^{k-1}$
- Probability of RAID failure is  $1-(1-p)^k \cdot kp(1-p)^{k-1}$
- There is a flawed assumption here: What is it?

Operating System Concepts – 8<sup>th</sup> Edition 12.47 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.48 Silberschatz, Galvin and Gagne ©2009



## 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.

Operating System Concepts – 8<sup>th</sup> Edition

12.49

Silberschatz, Galvin and Gagne ©2009



## 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.

Operating System Concepts – 8<sup>th</sup> Edition

12.50

Silberschatz, Galvin and Gagne ©2009



## 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.

Operating System Concepts – 8<sup>th</sup> Edition

12.51

Silberschatz, Galvin and Gagne ©2009



## 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 as CD-ROM and DVD, come from the factory with the data pre-recorded.

Operating System Concepts – 8<sup>th</sup> Edition

12.52

Silberschatz, Galvin and Gagne ©2009



## Tapes

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

Operating System Concepts – 8<sup>th</sup> Edition

12.53

Silberschatz, Galvin and Gagne ©2009



## Operating System Issues

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

Operating System Concepts – 8<sup>th</sup> Edition

12.54

Silberschatz, Galvin and Gagne ©2009



## 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 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.55 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.56 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.57 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.58 Silberschatz, Galvin and Gagne ©2009

## IU MDSS

- <http://storage.iu.edu/mdss.html>
- One terabyte disk cache
- 175TB storage
- Directory structure always visible to user
- Less frequently accessed items are migrated to tape

Operating System Concepts – 8<sup>th</sup> Edition 12.59 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.60 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.61 Silberschatz, Galvin and Gagne ©2009

## Reliability

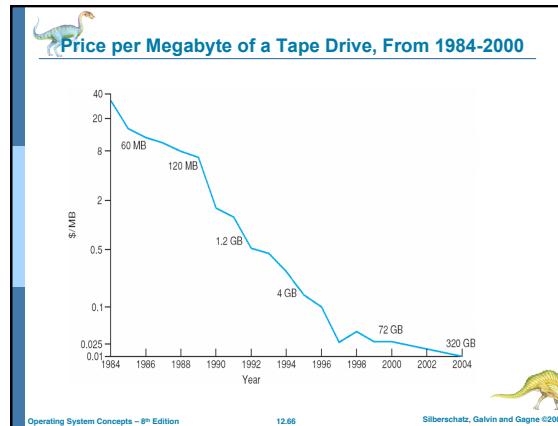
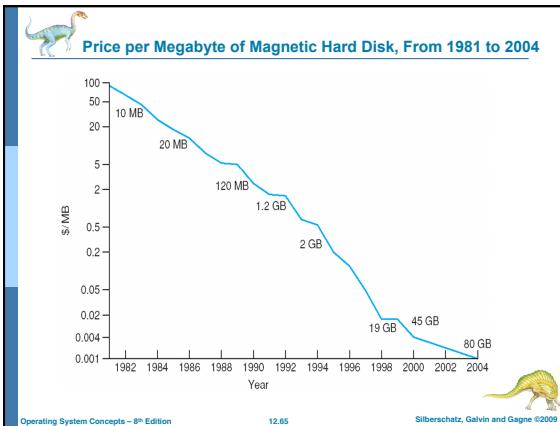
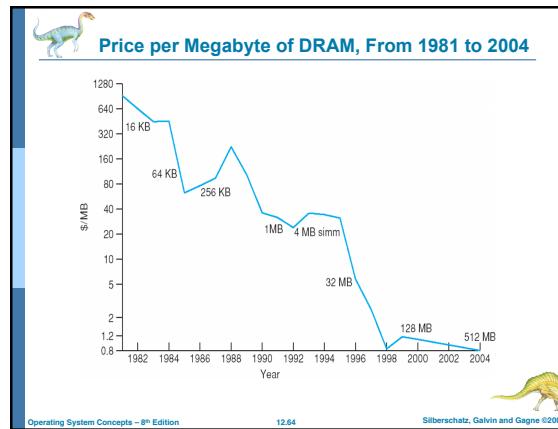
- A fixed disk drive is likely to be more reliable than a removable disk or tape drive.
- An optical cartridge is likely to be more reliable than a magnetic disk or tape.
- A head crash in a fixed hard disk generally destroys the data, whereas the failure of a tape drive or optical disk drive often leaves the data cartridge unharmed.

Operating System Concepts – 8<sup>th</sup> Edition 12.62 Silberschatz, Galvin and Gagne ©2009

## 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.

Operating System Concepts – 8<sup>th</sup> Edition 12.63 Silberschatz, Galvin and Gagne ©2009



## End of Chapter 12



Operating System Concepts - 8<sup>th</sup> Edition,

Silberschatz, Galvin and Gagne ©2009