# 10: Storage and File System Basics

Last Modified:
6/15/2004 12:10:04 PM
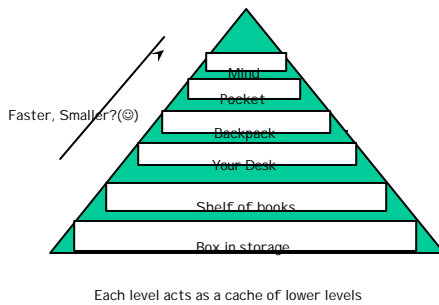
-1

## Storage Hierarchy

Faster, Smaller,
More Expensive

Registers
L1 Cache
L2 Cache
DRAM — Volatile
DISK — Non-Volatile
TAPE

Each level acts as a cache of lower levels

-2

## Example

Faster, Smaller?(☺)

Mind
Pocket
Backpack
Your Desk
Shelf of books
Box in storage

Each level acts as a cache of lower levels

-3

## Secondary Storage

❒ "Secondary" because unlike primary memory does not permit direct execution of instructions or data retrieval via load/store instructions
❒ Usually means hard disks
❒ Tends to be larger, cheaper and slower than primary memory
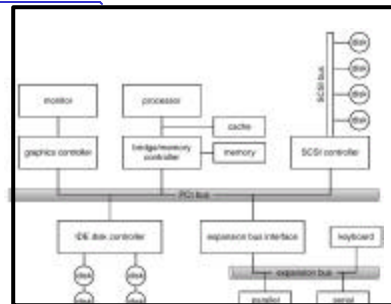❒ Persistent/Non-volatile
  ○ Like "durability" for transactions

-4

## Tertiary Storage Devices

❒ Used primarily as backup and archival storage
❒ Low cost is the defining characteristic
❒ Often consists of *removable media*
  ○ Common examples of removable media are CD-ROMs, tapes, etc.
❒ As disks get cheaper and cheaper, duplicating data on multiple disks becomes more and more attractive as a backup strategy
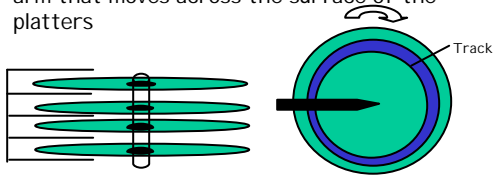
-5

## Typical PC



-6

1

## Disk Basics

❏ Disk drives contain metallic platters spinning around a central spindle

❏ Read/write head assembly is mounted on an arm that moves across the surface of the platters



Track

-7

## Terms

❏ Track = one ring around the surface of one of the platters

❏ Sector = one piece of a track (usually 512 bytes); More sectors in outer tracks

❏ Cylinder = all tracks at the same distance from the center of the platters (I.e. all tracks readable without moving the disk arm)

-8

## Disk Addressing

❏ Early disks were addressed with cylinder #, surface # and sector #

❏ Today disks hide information about their geometry
  o Disks export a logical array of blocks
  o Disk itself maps from logical block address (LBA) to cylinder/surface/sector
  o Allows disk to remap bad sectors (when formatted disk reserves some sectors to use as replacements)
  o Allows disk to hide the non-uniformity of the storage
    • More data on outer tracks, etc.

❏ Disks also have internal caches so that not all requests go to the media
  o On reads take advantage of multiple accesses to the same track
  o On writes, say write is "done" when it is memory inside the disk

-9

## Disk Formatting

❏ Low-level formatting involves dividing the magnetic media into sectors
  o Each sector actually consists of a header, data and a trailer
  o Header and trailer contain information like sector number and error correcting codes (ECC)
  o ECC is additional redundant bits that can often correct for bit errors in the stored value

❏ OS also formats drive
  o 1st divides into partitions – each partition can be treated as a logically separate drive
  o 2cd file system formatting of partitions (more on that later)

-10

## Disk Interfaces

❏ Interface to the disk
  o Request specified with LBA and length
  o Request placed on bus, later reply placed on bus

❏ Device driver hide these details
  o Provide abstraction of synchronous disk read

❏ OS use the disk to provide services
  o Virtual memory

❏ OS exports higher level abstractions
  o File systems

❏ Some applications use the device driver interface to build abstractions of their own (get their own partition)
  o Database systems

-11

## RAID

❏ Expose an array of sectors but implemented as multiple physical disks

❏ Arrangement and relationship of disks

❏ RAID levels

-12

## Disk Performance

❏ Divide the time for an access into stages
  ○ Seek time – time to move the disk arm to the correct cylinder
    • How fast can mechanical arm move? Improves some with smaller disks but not much
  ○ Rotational delay – time waiting for the correct sector to rotate under the read/write head
    • How fast can spindle turn? RPMs go up but slowly
  ○ Transfer time – once head is over the right spot how long to transfer all the data
    • Larger for larger transfers
    • Rate determined by RPMs and by density of the bits on the disk (density going up very quickly!)

❏ Getting good performance from a drive (seeing impact of a "faster" drive" means avoiding seek and rotational delay)
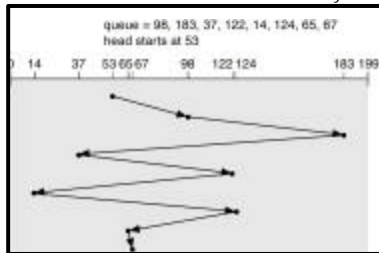
-13

## Avoiding Seek and Rotational Delay

❏ To take advantage of higher transfer rate, OS must transfer larger and larger chunks of data at a time and avoid seek and rotational delay
  ○ Size and placement of virtual memory pages?
  ○ Size and placement of FS blocks?

❏ OS tries to avoid seek and rotational delay by placing things on disk together that will be accessed together

❏ Can also avoid seek and rotational delay by queuing up multiple disk requests and servicing them in an order that minimizes head movement (disk scheduling)
  ○ Like with CPU scheduling, there are many disk scheduling algorithms

-14

## First Come First Serve (FCFS)

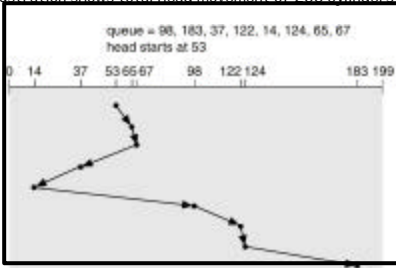Illustration shows total head movement of 640 cylinders



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

0  14  37  53 65 67  98  122 124  183 199

-15

## Shortest Seek Time First (SSTJ)

❏ 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.

-16

## SSTF

Illustration shows total head movement of 236 cylinders



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

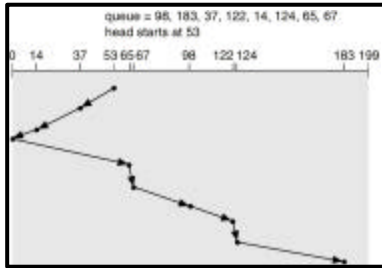0  14  37  53 65 67  98  122 124  183 199

-17

## 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 scheduling*

-18

## SCAN (Cont.)

Illustration shows total head movement of 208 cylinders

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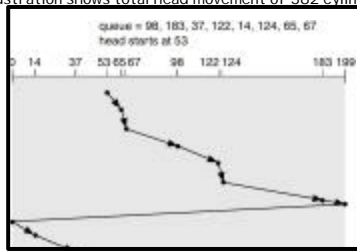head starts at 53

0  14  37  53 65 67  98  122 124  183 199



-19

## C-SCAN

❑ Provides a more uniform wait time than SCAN (with scan those in middle wait less)
❑ 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.

-20

## C-SCAN (Cont.)

Illustration shows total head movement of 382 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



Misleading because seek time not a linear function of number of cylinders

-21

## 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.

-22

## C-LOOK (Cont.)

Illustration shows total head movement of 322 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



-23

## Selecting a Disk-Scheduling Algorithm

❑ SSTF is common and has a natural appeal
  ○ Starvation not observed to be a problem in practice
❑ 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.
❑ Either SSTF or C-LOOK is a reasonable choice for the default algorithm.

-24

4

## Tracking Technology Trends

❑ Exact comparison between technologies changes all the time
   ○ How much slower is disk than main memory?
   ○ Variation even in disks and various memory technologies
❑ Tracking these things takes a fair amount of work

-25

## Drive Specs (6/2004)



-26

## More details!



-27

## Memory Types and Prices (6/2004)
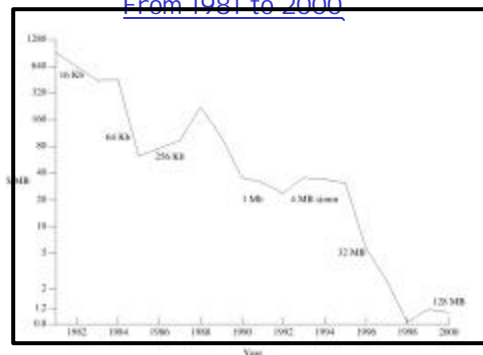


-28

## Two random points (2002)

❑ Memory: 128 MB, PC 133, SDRAM, $45
   ○ $0.35/MB
   ○ ~8 nanosecond access time
❑ Disk: 20 GB, Ultra ATA/100, $109
   ○ $0.005/MB (1/2 penny per MB!!)
   ○ 9.5 ms average seek (what is average? Seek time increases with number of tracks moved but not linearly)
   ○ 4.16 ms average latency (1/2 rotation at 7200 RPM?)
   ○ 100 MB/sec burst transfer (25-41 MB/sec sustained transfer)
❑ Disk/Memory Ratios
   ○ Price: 1/70
   ○ Size: 160/1
   ○ Speed (Access time): 13 ms/8ns = 1625000/1
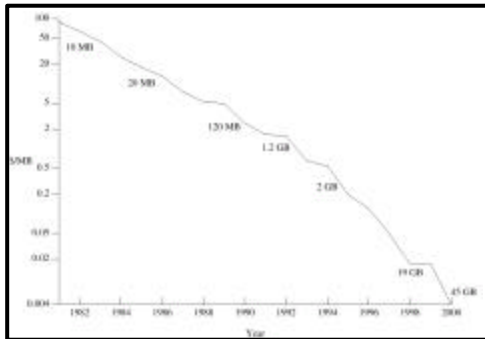   ○ Speed (Transfer rate): 40 MB/s / 1.1 GB/s = 1/30

-29

## Price per Megabyte of DRAM, From 1981 to 2000



-30

5

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

## OS adapts to performance trends?

❒ For the OS to make the right choices if needs to be aware of the trade-offs
  ○ Is the speed comparison between registers, DRAM and disk like the difference between your mind, your pocket and your book shelf *OR* is more like the difference between your pocket, the bookstore and Pluto?
  ○ How much computation/meta-data storage is reasonable to do to avoid a disk access?
  ○ Should we use DRAM as a file cache or to store more memory page for processes?
❒ "Right" answer changes with new generations of technology and OS source lives much longer than that?
❒ Can OS measure performance and be coded to react to measurements?

## File Systems

❒ Today talked a bit about disk internals
❒ Despite complex internals, disks export a simple array of sectors

❒ Next, how do we go from that to a file system?
❒ What do we exactly do we expect from a file system?

## Outtakes

## Drive interfaces

❒ SCSI
  ○ Fast, wide,....
❒ UltraATA
  ○ ATA, IDE
❒ SerialATA

## Next time

❒ Consider other tertiary storage device options instead of start on file systems
❒ Floppy drives
  ○ Thin flexible disk coated with magnetic material enclosed in a protective plastic case
❒ CDROMS
  ○ Spiral towards center not concentric circles like hard drives
❒ CD-Rs and CD-RWs
❒ WORM Disks
❒ Magneto-Optical disks
❒ Tapes
❒ Introduce reliability aspect of different media

# RAID

# Fall 2002: Current Drive Specs

# Fall 2002: More details!

# Fall 2002: Memory Types and Prices

# Fall 2002: Two random points

❏ Memory: 128 MB, PC 133, SDRAM, $45
  ○ $0.35/MB
  ○ ~8 nanosecond access time
❏ Disk: 20 GB, Ultra ATA/100, $109
  ○ $0.005/MB (1/2 penny per MB!!)
  ○ 9.5 ms average seek (what is average? Seek time increases with number of tracks moved but not linearly)
  ○ 4.16 ms average latency (1/2 rotation at 7200 RPM?)

# 100 MB/sec burst transfer (25-41 MB/sec sustained transfer)

❏ Disk/Memory Ratios
  ○ Price: 1/70
  ○ Size: 160/1
  ○ Speed (Access time): 13 ms/8ns = 1625000/1
  ○ Speed (Transfer rate): 40 MB/s / 1.1 GB/s = 1/30

## Disk Scheduling

❒ FCFS
  ○ Service requests in order they arrive
  ○ No possibility of data inconsistency
❒ SSTF (Shortest seek time first)
  ○ Do closest request first
  ○ Unfairly favors middle tracks
❒ SCAN (elevator scheduling)
  ○ Service requests in one direction until done, then reverse
❒ C-SCAN
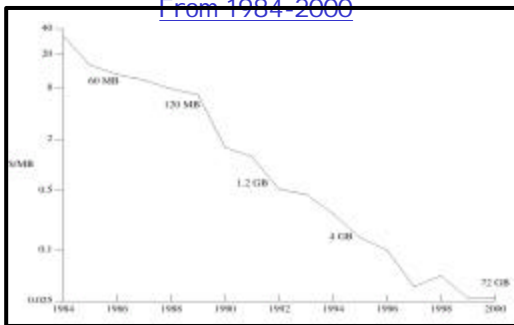  ○ Like SCAN, but when done do not reverse, return to the beginning

-43

## 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.

-44

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



-45