

# CS 245: Database System Principles

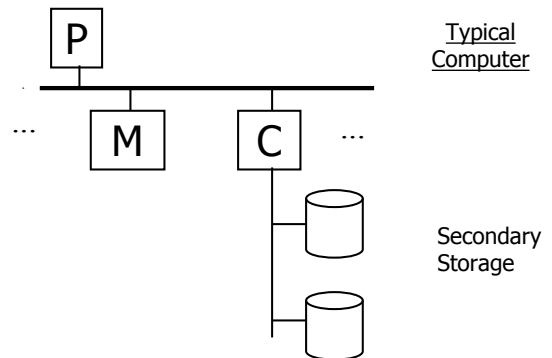
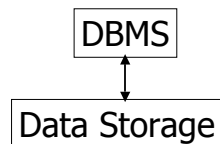
## Notes 02: Hardware

Steven Whang

## Outline

- Hardware: Disks
- Access Times
- Example - Megatron 747
- Optimizations
- Other Topics:
  - Storage costs
  - Using secondary storage
  - Disk failures

### Hardware



### Processor

Fast, slow, reduced instruction set, with cache, pipelined...  
Speed: 100 → 500 → 1000 MIPS

### Memory

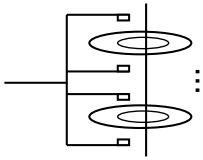
Fast, slow, non-volatile, read-only, ...  
Access time:  $10^{-6}$  →  $10^{-9}$  sec.  
 $1 \mu\text{s}$  →  $1 \text{ ns}$

### Secondary storage

Many flavors:

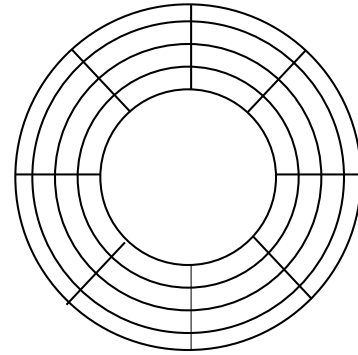
- Disk: Floppy (hard, soft)  
Removable Packs  
Winchester  
Ram disks  
Optical, CD-ROM...  
Arrays
- Tape Reel, cartridge  
Robots

### Focus on: "Typical Disk"



Terms: Platter, Head, Actuator  
Cylinder, Track  
Sector (physical),  
Block (logical), Gap

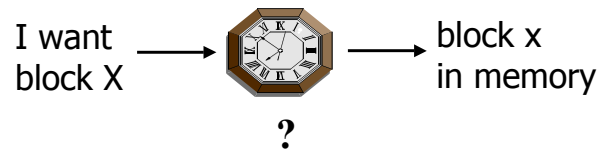
### Top View



### "Typical" Numbers

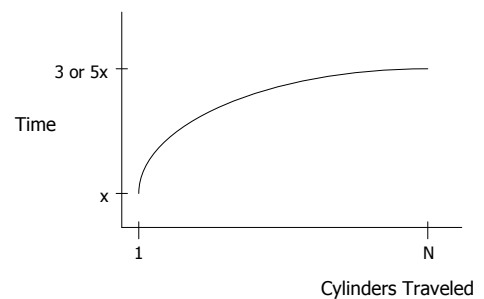
Diameter: 1 inch → 15 inches  
Cylinders: 100 → 2000  
Surfaces: 1 (CDs) →  
(Tracks/cyl) 2 (floppies) → 30  
Sector Size: 512B → 50K  
Capacity: 360 KB (old floppy)  
→ 2 TB

### Disk Access Time



Time = Seek Time +  
Rotational Delay +  
Transfer Time +  
Other

### Seek Time

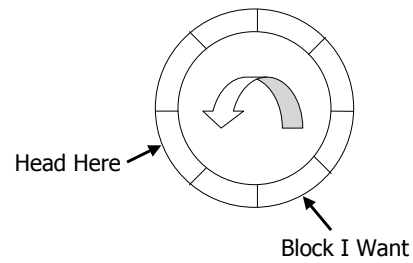


## Average Random Seek Time

$$S = \frac{\sum_{i=1}^N \sum_{\substack{j=1 \\ j \neq i}}^N \text{SEEKTIME}(i \rightarrow j)}{N(N-1)}$$

"Typical" S: 10 ms → 40 ms

## Rotational Delay



## Average Rotational Delay

R = 1/2 revolution

"typical" R = 8.33 ms (3600 RPM)

## Transfer Rate: t

- "typical" t: 1 → 3 MB/second
- transfer time:  $\frac{\text{block size}}{t}$

## Other Delays


- CPU time to issue I/O
- Contention for controller
- Contention for bus, memory

"Typical" Value: 0

- So far: Random Block Access
- What about: Reading "Next" block?

If we do things right (e.g., Double Buffer, Stagger Blocks...)

Time to get block =  $\frac{\text{Block Size}}{t} + \text{Negligible}$

- 
- skip gap
  - switch track
  - once in a while, next cylinder

**Rule of Thumb**

Random I/O: Expensive  
Sequential I/O: Much less

- Ex: 1 KB Block
  - » Random I/O: ~ 20 ms.
  - » Sequential I/O: ~ 1 ms.

Cost for Writing similar to Reading

.... unless we want to verify!  
need to add (full) rotation +  $\frac{\text{Block size}}{t}$

- To Modify a Block?

To Modify Block:

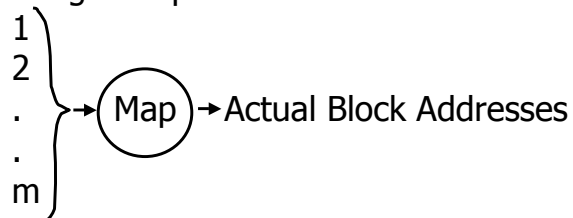
- (a) Read Block
- (b) Modify in Memory
- (c) Write Block
- [(d) Verify?]

Block Address:

- Physical Device
- Cylinder #
- Surface #
- Sector

Complication: Bad Blocks

- Messy to handle
- May map via software to integer sequence



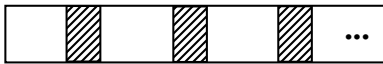
## An Example Megatron 747 Disk (old)

- 3.5 in diameter
- 3600 RPM
- 1 surface
- 16 MB usable capacity ( $16 \times 2^{20}$ )
- 128 cylinders
- seek time: average = 25 ms.  
adjacent cyl = 5 ms.

- 1 KB blocks = sectors
- 10% overhead between blocks
- capacity = 16 MB =  $(2^{20})16 = 2^{24}$
- # cylinders =  $128 = 2^7$
- bytes/cyl =  $2^{24}/2^7 = 2^{17} = 128 \text{ KB}$
- blocks/cyl =  $128 \text{ KB} / 1 \text{ KB} = 128$

3600 RPM  $\rightarrow$  60 revolutions / sec  
 $\rightarrow$  1 rev. = 16.66 msec.

One track:



Time over useful data:  $(16.66)(0.9) = 14.99 \text{ ms}$ .  
Time over gaps:  $(16.66)(0.1) = 1.66 \text{ ms}$ .  
Transfer time 1 block =  $14.99/128 = 0.117 \text{ ms}$ .  
Trans. time 1 block+gap =  $16.66/128 = 0.13 \text{ ms}$ .

## Burst Bandwidth

1 KB in 0.117 ms.

$$\text{BB} = 1/0.117 = 8.54 \text{ KB/ms.}$$

or

$$\begin{aligned} \text{BB} &= 8.54 \text{ KB/ms} \times 1000 \text{ ms/1sec} \times 1 \text{ MB}/1024 \text{ KB} \\ &= 8540/1024 = 8.33 \text{ MB/sec} \end{aligned}$$

Sustained bandwidth (over track)  
128 KB in 16.66 ms.

$$\text{SB} = 128/16.66 = 7.68 \text{ KB/ms}$$

or

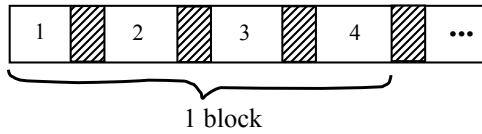
$$\text{SB} = 7.68 \times 1000/1024 = 7.50 \text{ MB/sec.}$$

$T_1$  = Time to read one random block

$T_1$  = seek + rotational delay + TT

$$= 25 + (16.66/2) + .117 = 33.45 \text{ ms.}$$

Suppose OS deals with 4 KB blocks



$$T_4 = 25 + (16.66/2) + (.117) \times 1 + (.130) \times 3 = 33.83 \text{ ms}$$

[Compare to  $T_1 = 33.45 \text{ ms}$ ]

$T_T$  = Time to read a full track  
(start at any block)

$$T_T = 25 + (0.130/2) + 16.66^* = 41.73 \text{ ms}$$

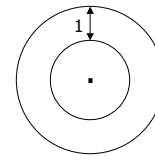
↑  
to get to first block

\* Actually, a bit less; do not have to read last gap.

### The NEW Megatron 747 (Example 11.3, 1st Ed.)

- 16 Surfaces, 3.5 Inch diameter  
– outer 1 inch used
- $2^{14} = 16,384$  Tracks/surface
- 128 Sectors/track
- $2^{12} = 4096$  Bytes/sector

- 128 GB Disk
- If all tracks have 128 sectors
  - Outermost density: 420,000 bits/inch
  - Inner density: 1000,000 bits/inch



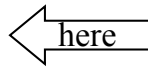
- Outer third of tracks: 160 sectors
- Middle third of tracks: 128
- Inner third of tracks: 96
- Density: 530,000 → 742,000 bits/inch

### Timing for new Megatron 747 (Ex 11.5)

- Time to read 16,384-byte block:
  - MIN: 0.253 ms
  - MAX: 25.96 ms
  - AVE: 10.88 ms

## Outline

- Hardware: Disks
- Access Times
- Example: Megatron 747
- Optimizations
- Other Topics
  - Storage Costs
  - Using Secondary Storage
  - Disk Failures



## Optimizations (in controller or O.S.)

- Disk Scheduling Algorithms
  - e.g., elevator algorithm
- Track (or larger) Buffer
- Pre-fetch
- Arrays
- Mirrored Disks
- On Disk Cache

## Double Buffering

Problem: Have a File

- » Sequence of Blocks B1, B2

Have a Program

- » Process B1
- » Process B2
- » Process B3
- ⋮

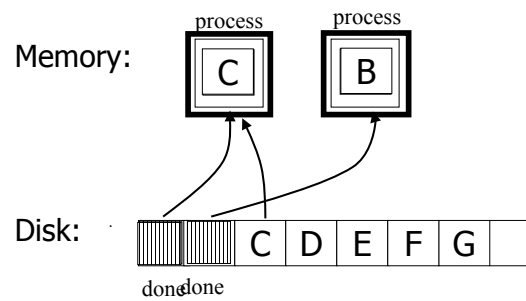
## Single Buffer Solution

- (1) Read B1 → Buffer
- (2) Process Data in Buffer
- (3) Read B2 → Buffer
- (4) Process Data in Buffer ...

Say  $P$  = time to process/block  
 $R$  = time to read in 1 block  
 $n$  = # blocks

Single buffer time =  $n(P+R)$

## Double Buffering



Say  $P \geq R$

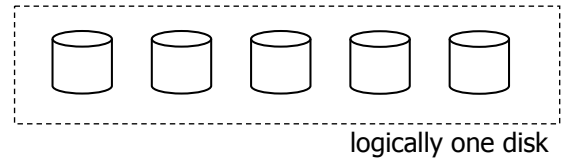
$P$  = Processing time/block  
 $R$  = IO time/block  
 $n$  = # blocks

What is processing time?

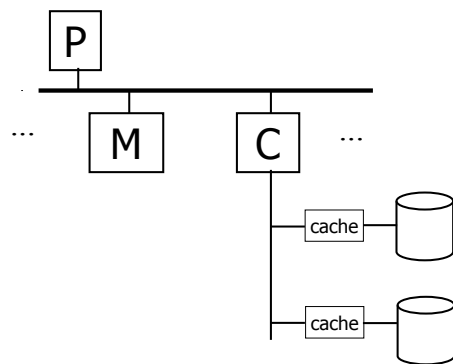
- Double buffering time =  $R + nP$
- Single buffering time =  $n(R+P)$

## Disk Arrays

- RAIDs (various flavors)
- Block Striping
- Mirrored

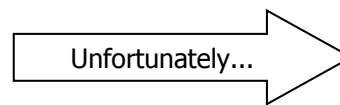


## On Disk Cache



## Block Size Selection?

- Big Block  $\rightarrow$  Amortize I/O Cost

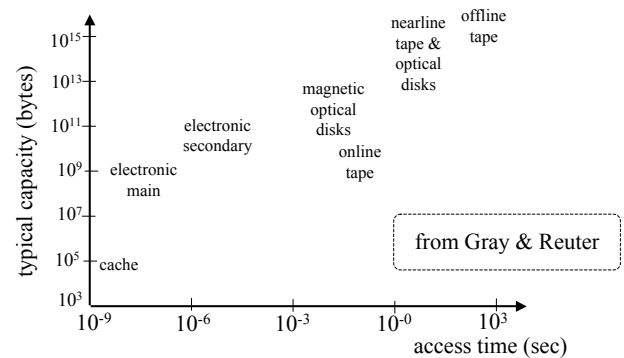


- Big Block  $\Rightarrow$  Read in more useless stuff!  
and takes longer to read

## Trend

- As memory prices drop,  
blocks get bigger ...

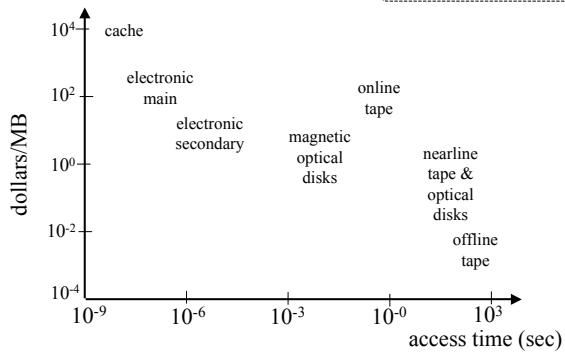
## Storage Cost





## Storage Cost

from Gray & Reuter



CS 245

Notes 2

49

## Using secondary storage effectively

(Sec. 11.4)

- Example: Sorting data on disk
- Conclusion:
  - I/O costs dominate
  - Design algorithms to reduce I/O
- Also: How big should blocks be?

CS 245

Notes 2

50

## Five Minute Rule

- THE 5 MINUTE RULE FOR TRADING MEMORY FOR DISC ACCESSES

Jim Gray & Franco Putzolu  
May 1985

- The Five Minute Rule, Ten Years Later  
Goetz Graefe & Jim Gray  
December 1997

CS 245

Notes 2

51

## Five Minute Rule

- Say a page is accessed every X seconds
- CD = cost if we keep that page on disk
  - \$D = cost of disk unit
  - I = numbers IOs that unit can perform
  - In X seconds, unit can do XI IOs
  - So  $CD = \$D / XI$

CS 245

Notes 2

52

## Five Minute Rule

- Say a page is accessed every X seconds
- CM = cost if we keep that page on RAM
  - \$M = cost of 1 MB of RAM
  - P = numbers of pages in 1 MB RAM
  - So  $CM = \$M / P$

CS 245

Notes 2

53

## Five Minute Rule

- Say a page is accessed every X seconds
- If CD is smaller than CM,
  - keep page on disk
  - else keep in memory
- Break even point when  $CD = CM$ , or

$$X = \frac{\$D P}{I \$M}$$

CS 245

Notes 2

54

## Using '97 Numbers

- P = 128 pages/MB (8KB pages)
- I = 64 accesses/sec/disk
- \$D = 2000 dollars/disk (9GB + controller)
- \$M = 15 dollars/MB of DRAM
  
- X = 266 seconds (about 5 minutes)  
(did not change much from 85 to 97)

## Disk Failures

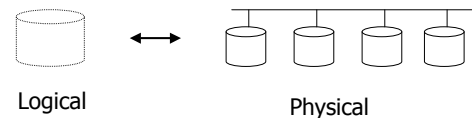
- Partial → Total
- Intermittent → Permanent

## Coping with Disk Failures

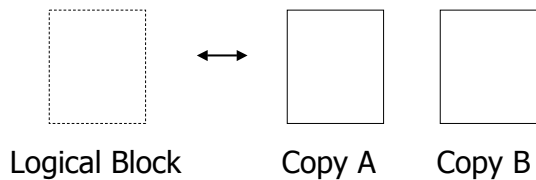
- Detection
  - e.g. Checksum
  
- Correction
  - ⇒ Redundancy

## At what level do we cope?

- Single Disk
  - e.g., Error Correcting Codes
- Disk Array

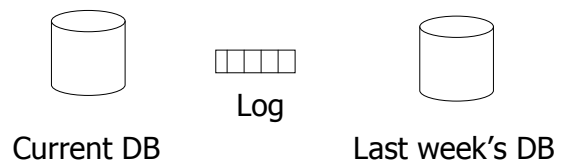


→ Operating System  
e.g., Stable Storage



→ Database System

- e.g.,



## Summary

- Secondary storage, mainly disks
- I/O times
- I/Os should be avoided,  
especially random ones.....

## Outline

- Hardware: Disks
- Access Times
- Example: Megatron 747
- Optimizations
- Other Topics
  - Storage Costs
  - Using Secondary Storage
  - Disk Failures

