Hardware

Read Sections 11.2, 11.3, 11.7 of Garcia-Molina et al.

Slides derived from those by Hector Garcia-Molina Some images by Wikipedia

1

<u>Outline</u>

- Hardware: Disks
- Access Times
- Example Megatron 747
- Reliability
- RAID

Hardware





<u>Typical</u> <u>Computer</u>

Secondary Storage

Processor

Fast, slow, reduced instruction set, with cache, pipelined... Speed: $1000 \rightarrow 10000$ MIPS

<u>Memory</u>

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

<u>Secondary storage</u> Hard Disks <u>Tertiary storage</u> **Optical disks:** CD-ROM • DVD-ROM... Tape Cartridges Robots

Focus on: "Typical Disk"



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

Disk Architecture



Top View



9

<u>"Typical" Numbers</u> Diameter: 1 inch \rightarrow 15 inches (1 inch=2.54 cm)2.5 cm \rightarrow 38.1 cm) Cylinders: $10000 \rightarrow 50000$ Surfaces: 2 -> 30 (Tracks/cyl) Sector Size: 512B \rightarrow 50KB Capacity: 72 GB \rightarrow 2TB

Diameter

- Form factors:
 - 8 inches
 - 5.25 inches
 - 3,5 inches
 - 2,5 inches
 - 1,8 inches
 - 1 inch



エム

Disk Access Time

I want _____ block X



___block x in memory

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



Average Random Seek Time



N(N-1)

"Typical" S: 3 ms \rightarrow 10 ms

Rotational Delay



17

Average Rotational Delay

R = 1/2 revolution

"typical" R = 4.17 ms (7200 RPM) R=3 ms (10000 RPM) R=2 ms (15000 RPM)

Transfer Rate: t

- "typical" t: 60 MB/second
- transfer time: <u>block size</u>

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?

Time to get = <u>Block Size</u> + Negligible block t

- skip gap

Cost for <u>Writing</u> similar to <u>Reading</u>

.... unless we want to verify! need to add (full) rotation + <u>Block size</u>

t

To <u>Modify</u> a Block?

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

Block Address:

- Physical Device
- Cylinder (Track) #
- Surface #
- Sector

Complication: Bad Blocks

Messy to handle

 May map via software to integer sequence

)→ Actual Block Addresse

An Example

Megatron 747 Disk

- 8 platters, 16 surfaces

er

- 2¹⁴=16,384 tracks per surface (16,384 cylinders)
- 2⁷=128 sectors per track
- 2¹²=4096 bytes per sector
- Capacity
 - $Disk=2^{4*}2^{14*}2^{7*}2^{12}=2^{37}=128GB$
 - Single track=2⁷*2¹²=512KB

Megatron 747 Disk

- Rotation speed: 7200 RPM
- Average seek time: 8.5 ms

Layout

- Radius: 1.75 inches
- The tracks occupy the outer inch
- The inner 0.75 inch is unoccupied
- Track density in the radial direction: 16,384 tracks per inch
- 10% overhead between blocks

Density of bits

- Outermost track
 - Length= $3.5\pi \approx 11$ inches
 - One track = 512KB = 4Mbits
 - 90% of 11 inches holds 4Mbits
 - Density=420,000 bits per inch
- Innermost track
 - 90% of 4.71 inches holds 4Mbits
 - Density~1Mbit per inch

Density of bits

- To avoid such a high difference of density, the disk stores more sectors on the outer track than on the inner tracks
 - 96 sectors per track in the inner third
 - 128 in the middle third
 - 160 in the outer third
- The density varies from 742,000 bits per inch to 530,000 bits per inch

7200 RPM \rightarrow 120 revolutions / sec \rightarrow 1 rev. = 8.33 msec.

One track:



ime over useful data:(8.33)(0.9)=7.5 ms. ime over gaps: (8.33)(0.1) = 0.833 ms. ransfer time 1 sector = 7.5/128=0.059 ms. rans. time 1 sector+gap=8.33/128=0.065ms

$\frac{\text{Burst Bandwith}}{4 \text{ KB in } 0.059 \text{ ms.}}$ BB = 4/0.059 = 68 KB/ms.

or

BB =68 KB/ms x 1000 ms/1sec x 1MB/1024KB = 68,000/1024 = 66.4 MB/sec

<u>Sustained bandwith</u> (over track) 512 KB in 8.33 ms.

SB = 512/8.33 = 61.5 KB/ms

or

$SB = 61.5 \times 1000/1024 = 60 MB/sec.$

T_1 = Time to read one random block

 $T_1 = seek + rotational delay + TT$

= 8.5 + (8.33/2) + 0.059 = 12.72 ms.

Suppose OS deals with 16 KB blocks



$T_4 = 8.5 + (8.33/2) + 0.059*1 + (0.065) * 3 = 12.92 \text{ ms}$ [Compare to $T_1 = 12.72 \text{ ms}$]



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

Block Size Selection?

• Big Block \rightarrow Amortize I/O Cost



Big Block ⇒ Read in more useless stuff!
 and takes longer to read

Reliability

- Measured by the Mean Time to Failure (MTTF):
 - Length of time by which 50% of a population of disks will have failed catastrophically (head crash, no longer readable)
 - For modern disks, the MTTF is 10 years
 - This means that, on average, after 10 years it will crash
 - We can assume that every year 5% of the disks fail (uniform distribution assumption)
 - Probability that a disk fails in one year $P_F = 5\% = 1/20$



MTTF

- Expected value of the failure year:
- MTTF=E(year)=
 =0.05*1+....+0.05*20=
 =0.05*20*(20+1)/2=21/2 ≈ 10

Disk Arrays

- Redundant Arrays of Inexpensive Disks (RAID)
- Two aims: increase speed and reliability

- Uses "block level striping"
 - Blocks that are consecutive for the OS are distributed evenly across different disks
 - RAID 0
- A1 A2 consecutive blocks: A1-A8
 A3 A4
 A5 A6
 A7 A8

- Improves reading and writing speed
 - With two disks, two blocks can be read at the same time
 - A request for block "A1" would be serviced by disk 1. A simultaneous request for block A3 would have to wait, but a request for A2 could be serviced concurrently
- Reduces reliability: if one disk fails, the data is lost.

- P(data loss)=P(disk1 fails or disk2 fails)=
- =P(disk1 fails)+P(disk2 fails)-P(disk1 fails and disk2 fails)=
- $= P_{F} + P_{F} P_{F} + P_{F} = 2P_{F} P_{F}^{2} =$
- =2*0.05-0.0025=0.0975

- Number of years=1/0.0975 \approx 10
- MTTF =E(year)=
 ≈ 0.0975*10*(10+1)/2 ≈ 11/2 ≈5.5

- Creates an exact copy (or mirror) of a set of data on two or more disks.
- Typically, a RAID 1 array contains two disks
- Improved
 - Reading speed: two blocks can be read at the same time
 - Reliability: if one disk crashes, we can use the other
- Writing speed remains the same

RAID1A1A1A2A2A3A3A4A4

- Two disks with MTTF of 10 years
- What is the MTTF resulting in data loss?
- Data loss happens when one disk fails and the other fails as well while we are replacing the first.
- Supposing it takes 3 hours to replace the first disk. This is 1/2920 of a year
- P(fails rep)=1/2920=3.42E-04

- The probability that the second disk fails while replacing the first is P(fails1 and fails2 rep)= =5E-2*5E-2*3.42E-04=8.55E-07
- P(data loss)=P(fails1 and fails2 rep or fails2 and fails1 rep)=

- = P(fails1 and fails2 rep) + P(fails2
 and fails1 rep)-P(fails2 and fails1
 rep and fails1 and fails2 rep)=
 ≈2*8.55E-07
 -1 71E 06
- =1.71E-06

- Number of years=1/1.71E-06 ≈ 584795
- MTTF=E(years)=
 =1.71E-06*584795*584796/2=
 =584796/2=292398

 Uses block-level striping with a dedicated parity disk. RAID 4 A1 A2 A3 Ap Consecutive blocks A1-A3,B1-B3, B1 B2 B3 Bp C1 C2 C3 Cp C1-C3, D1-D3 D1 D2 D3 Dp

Parity block

- Bit i of the block in position j on the parity disk is the parity bit of the bits in position i in the blocks in position j in the other disks
- Eg., blocks of one byte, blocks A1-A3
 Disk1 11110000
 Disk2 10101010
 <u>Disk3 00111000</u>
 Disk4 01100010 (parity disk)

- Improves reading time: multiple blocks can be read at the same time
- Improves reliability: if one disk fails, we can reconstruct its content (assuming the others are correct)

- Problem:
 - When writing a block, we need to read and write the parity disk's block
 - This creates a bottleneck

 Uses block-level striping with parity data distributed across all member disks.

RAID 5

A1 A2 A3 Ap

B1 B2 Bp B3 C1 Cp C2 C3

Dp D1 D2 D3

- Reading and reliability as RAID 4
- Writing improved because the parity blocks are not all on one disk

 Uses block-level striping with dual parity data distributed across all member disks.

RAID 6

A1 A2 A3 Ap Aq

B1 B2 Bp Bq B3 C1 Cp Cq C2 C3 Dp Dq D1 D2 D3

- p and q blocks are computed with two different algorithms, e.g.
 - parity and Reed-Solomon
 - orthogonal dual parity
 - diagonal parity

- It is able to recover from the loss of two disks
- Writing improved because the parity blocks are not all on one disk

Nested RAID Levels

• RAID 0+1:



RAID 0+1

- If a disk fails, it can be rebuilt from the corresponding disk in the other RAID 0 batch
- If two disk fails from the same stripe, no recovery

RAID 1+0 o RAID 10



RAID 1+0 o RAID 10

- If a disk fails, it can be rebuilt from the corresponding disk in the other RAID 1 batches
- If two disk fails from the same RAID 1 batch, no recovery

RAID 5+0



