

RAID

- Redundant Array of Inexpensive Disks
 - “A Case for Redundant Array of Inexpensive Disks” Patterson, Gibson, and Katz 1988
 -old but crucial....

28
1

K. Selcuk Candan (CSE510)

281

Observations

- CPU speed
 - MIPS = $2^{\text{Year}-1984}$ (Joy's Law, 1985)
- Memory
 - Amdahl's constant: each CPU instruction per second requires one byte of Main Memory
 - 1GHz CPU speed....needs ~ 1GB memory
- Chip capacity
 - #Trans/Chip = $2^{\text{Year}-1964}$ (Moore's Law, 1975)

28
2

K. Selcuk Candan (CSE510)

282

Observation

- Secondary storage is largely mechanical!!!!
 - Density: $10^{(\text{Year}-1971)/10}$ (Frank, 1987)
 - Speed??

28
3

K. Selcuk Candan (CSE510)

283

Observation

- Secondary storage is largely mechanical!!!!
 - Density: $10^{(\text{Year}-1971)/10}$ (Frank, 1987)
 - Speed??
 - Seek & rotation
 - From '71 to '81
 - seek time improved 2 times
 - rotation 0 times (now rotation is faster...but not by much)

28
4

K. Selcuk Candan (CSE510)

284

Observation

- Secondary storage is largely mechanical!!!!
 - Density: $10^{(\text{Year}-1971)/10}$ (Frank, 1987)
 - Speed??
 - Seek & rotation
 - From '71 to '81
 - seek time improved 2 times
 - rotation 0 times (now rotation is faster...but not by much)
 - ...caching helps
 - there are “Main Memory DBMSs”, e.g., TimesTen
 -not a practical solution for all applications

28
5

K. Selcuk Candan (CSE510)

285

Solution.....

- Find a way to increase disk speed...
 - But how??

28
8

K. Selcuk Candan (CSE510)

288

Solution.....

- Find a way to increase disk speed...
 - But how??
- Use parallel I/O!!!
 - But how??

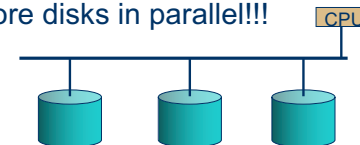
28
9

K. Selcuk Candan (CSE510)

289

Solution.....

- Find a way to increase disk speed...
 - But how??
- Use parallel I/O!!!
 - But how??
- Price of disks keep decreasing.....so, use more disks in parallel!!!



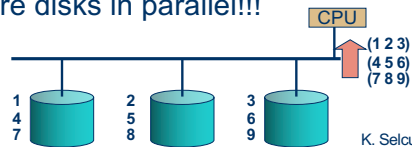
29
0

K. Selcuk Candan (CSE510)

290

Solution.....

- Find a way to increase disk speed...
 - But how??
- Use parallel I/O!!!
 - But how??
- Price of disks keep decreasing.....so, use more disks in parallel!!!



K. Selcuk Candan (CSE510)

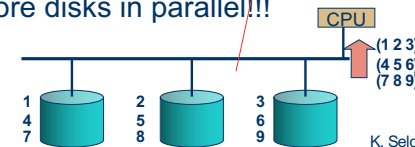
29
1

291

Solution.....

- Find a way to increase disk speed...
 - But how??
- Use parallel I/O!!!
 - But how??
- Price of disks keep decreasing.....so, use more disks in parallel!!!

Good news: Disk don't need to be perfectly synchronized



K. Selcuk Candan (CSE510)

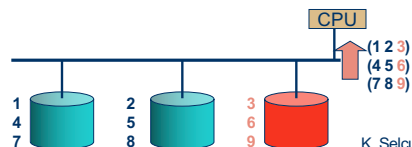
29
2

292

Problem...

- Reliability
 - as you add more disks Mean Time to Failure (MTTF) of the system drops:

$$MTTF(n) = MTTF/n$$



K. Selcuk Candan (CSE510)

29
3

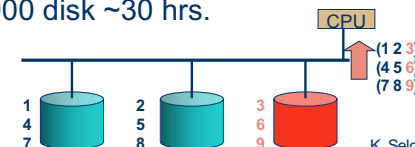
293

Problem...

- Reliability
 - as you add more disks Mean Time to Failure (MTTF) of the system drops:

$$MTTF(n) = MTTF/n$$

- 1 disk ~30K hrs.
- 1000 disk ~30 hrs.



K. Selcuk Candan (CSE510)

29
4

294

Solution

- Use redundancy to increase reliability
- For every G disks in the array use C disks for error detection and correction
 - D: total # of disks with data (exc. check disks)
 - G: number of data disks in a group (exc. check disks)
 - In a single-group RAID G=D
 - C number of check disks in a group
 - ng: number of groups

29
5

K. Selcuk Candan (CSE510)

295

So what do we pay??

$$\text{overhead} = \frac{C}{G}$$

29
8

K. Selcuk Candan (CSE510)

298

Different applications

- Streaming media
 - High data rate needed
 - G disks are used in synchrony
- Databases (transactions)
 - Fast small reads/writes needed
 - G disks work independently

29
9

K. Selcuk Candan (CSE510)

299

RAID 1 (Mirror Disks)

- G=1, C = 1
- Every write to a disk also does to the check disk
 - Overhead: 100%
 - Large reads: 2D/S
 - Large writes: D/S
 - Large R-M-W = 4D/3S
 - Small reads: 2D
 - Small writes: D
 - Small R-M-W = 4D/3

S: synchronization slowdown

30
0

K. Selcuk Candan (CSE510)

300

RAID 2 (Hamming code for ECC)

- Hamming code: corrects 1bit, detects 2bit errors

D7	D6	D5	D4	D3	D2	D1	D0
----	----	----	----	----	----	----	----

30
1

K. Selcuk Candan (CSE510)

301

RAID 2 (Hamming code for ECC)

- Hamming code: corrects 1bit, detects 2bit errors

D7	D6	D5	D4	D3	D2	D1	D0
----	----	----	----	----	----	----	----

C3	C2	C1	C0
----	----	----	----

30
2

K. Selcuk Candan (CSE510)

302

RAID 2 (Hamming code for ECC)

- Hamming code: corrects 1bit, detects 2bit errors

D7	D6	D5	D4	D3	D2	D1	D0
----	----	----	----	----	----	----	----

C3	C2	C1	C0
----	----	----	----

$$C0 = D7 \oplus D5 \oplus D3 \oplus D1$$

$$C1 = D7 \oplus D6 \oplus D3 \oplus D2$$

$$C2 = D7 \oplus D6 \oplus D5 \oplus D4$$

$$C3 = D7 \oplus D6 \oplus D5 \oplus D4 \oplus D3 \oplus D2 \oplus D1 \oplus D0$$

30
3

K. Selcuk Candan (CSE510)

303

RAID 2

- $G=25, C=5$
- Disks are bit interleaved!!!!

- Overhead: 20%
- Large reads: D/S
- Large writes: D/S
- Large R-M-W = D/SG
- Small reads: D/SG
- Small writes: D/2SG
- Small R-M-W = D/SG

small write penalty (you need to read first to verify)!!!

30
4

K. Selcuk Candan (CSE510)

304

RAID 2

On average better than RAID 1, because more disks can be used as data disks

- $G=25, C = 5$
- Disks are bit interleaved!!!!
 - Overhead: 20%
 - Large reads: D/S
 - Large writes: D/S
 - Large R-M-W = D/S
 - Small reads: D/SG
 - Small writes: $D/2SG$
 - Small R-M-W = D/SG

small write penalty (you need to read first to verify)!!!

30
5

K. Selcuk Candan (CSE510)

305

RAID 3

- Single (!) check disk per group
 - Can't find which disk failed...but, usually the disk controller will signal it anyway
- ..i.e., only one parity disk per group!!

30
6

K. Selcuk Candan (CSE510)

306

RAID 3

- $G=25, C = 1$
- Disks are bit interleaved!!!!
 - Overhead: 4%
 - Large reads: D/S
 - Large writes: D/S
 - Large R-M-W = D/S
 - Small reads: D/SG
 - Small writes: $D/2SG$
 - Small R-M-W = D/SG

30
7

K. Selcuk Candan (CSE510)

307

RAID 3

- $G=25, C = 1$
- Disks are bit interleaved!!!!
 - Overhead: 4%
 - Large reads: D/S
 - Large writes: D/S
 - Large R-M-W = D/S
 - Small reads: D/SG
 - Small writes: $D/2SG$
 - Small R-M-W = D/SG

SO...WHERE IS THE GAIN??

30
8

K. Selcuk Candan (CSE510)

308

RAID 3

- $G=25, C = 1$
- Disks are bit interleaved!!!!
 - Overhead: 4%
 - Large reads: D/S
 - Large writes: D/S
 - Large R-M-W = D/S
 - Small reads: D/SG
 - Small writes: $D/2SG$
 - Small R-M-W = D/SG

SO...WHERE IS THE GAIN??

More disks can be utilized for storing data...
....so, on the average per disk performance increases!!!

30
9

K. Selcuk Candan (CSE510)

309

RAID 4

- How to fix small accesses??

31
0

K. Selcuk Candan (CSE510)

310

RAID 4

- How to fix small accesses??
 - Can we achieve independent reads writes?

31
1

K. Selcuk Candan (CSE510)

311

RAID 4

- How to fix small accesses??
 - Can we achieve independent reads writes?
- We need to do more than one I/O per group

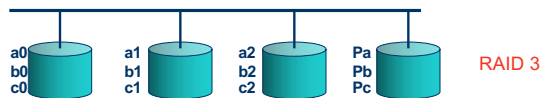
31
2

K. Selcuk Candan (CSE510)

312

RAID 4

- How to fix small accesses??
 - Can we achieve independent reads/writes?
- We need to do more than one I/O per group



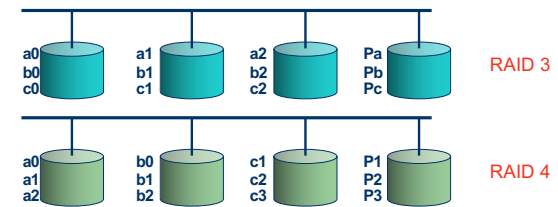
31
3

K. Selcuk Candan (CSE510)

313

RAID 4 (sector interleave)

- How to fix small accesses??
 - Can we achieve independent reads/writes?
- We need to do more than one I/O per group



31
4

K. Selcuk Candan (CSE510)

314

RAID 4

- $G=25, C=1$
- Disks are sector interleaved!!!!
 - Overhead: 4%
 - Large reads: D/S
 - Large writes: D/S
 - Large R-M-W = D/S
 - Small reads: **D**
 - Small writes: D/2G
 - Small R-M-W = D/G

31
5

K. Selcuk Candan (CSE510)

315

RAID 5

- Check disk is bottleneck for small writes
 - How can we remove it?

31
6

K. Selcuk Candan (CSE510)

316

RAID 5 (no single check disk)

- Check disk is bottleneck for small writes
 - How can we remove it?
- Distribute data onto check disks, vice versa!

31
7

K. Selcuk Candan (CSE510)

317

RAID 5 (no single check disk)

- Check disk is bottleneck for small writes
 - How can we remove it?
- Distribute data onto check disks, vice versa!

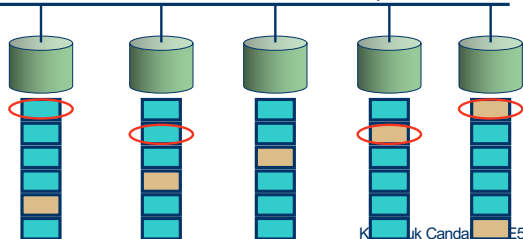
31
8

K. Selcuk Candan (CSE510)

318

RAID 5 (no single check disk)

- Check disk is bottleneck for small writes
 - How can we remove it?
- Distribute data onto check disks, vice versa!



31
9

K. Selcuk Candan (CSE510)

319

RAID 5

- $G=25, C=1$
- Disks are sector interleaved!!!!
 - Overhead: 4%
 - Large reads: D/S
 - Large writes: D/S
 - Large R-M-W = D/S
 - Small reads: $(1+C/G)D$
 - Small writes: $(1+C/G)D/4$
 - Small R-M-W = $(1+C/G)D/2$

32
0

K. Selcuk Candan (CSE510)

320