

Project Announcement

- Grading:
 - Part A: 3/4 of total grade
 - Part B: 1/4 of total grade
- Actual grade formula:
 - $3/4 * (\text{Max}(A1, 0.7 * A2)) + 1/4 * (\text{Max}(B_self, 0.9 * B_with_staff_code))$
 - A1 is submission for P1a on or before the part A deadline.
 - A2 is version of P1a submitted "any time" before the part B deadline.
 - B_self is part B running against your own implementation of LSP
 - 0.9*B_with_staff_code is part B running against the reference LSP

15-440 Distributed Systems

Lecture 13 – RAID

Thanks to Greg Ganger and Remzi Arapaci-Dusseau for slides

Replacement Rates

HPC1		COM1		COM2	
Component	%	Component	%	Component	%
Hard drive	30.6	Power supply	34.8	Hard drive	49.1
Memory	28.5	Memory	20.1	Motherboard	23.4
Misc/Unk	14.4	Hard drive	18.1	Power supply	10.1
CPU	12.4	Case	11.4	RAID card	4.1
motherboard	4.9	Fan	8	Memory	3.4
Controller	2.9	CPU	2	SCSI cable	2.2
QSW	1.7	SCSI Board	0.6	Fan	2.2
Power supply	1.6	NIC Card	1.2	CPU	2.2
MLB	1	LV Pwr Board	0.6	CD-ROM	0.6
SCSI BP	0.3	CPU heatsink	0.6	Raid Controller	0.6

Outline

- Using multiple disks
 - Why have multiple disks?
 - problem and approaches
- RAID levels and performance
- Estimating availability

Motivation: Why use multiple disks?

- Capacity
 - More disks allows us to store more data
- Performance
 - Access multiple disks in parallel
 - Each disk can be working on independent read or write
 - Overlap seek and rotational positioning time for all
- Reliability
 - Recover from disk (or single sector) failures
 - Will need to store multiple copies of data to recover
- So, what is the simplest arrangement?

Just a bunch of disks (JBOD)

- Yes, it's a goofy name
 - industry really does sell "JBOD enclosures"

Disk Subsystem Load Balancing

- I/O requests are almost never evenly distributed
 - Some data is requested more than other data
 - Depends on the apps, usage, time, ...
- What is the right data-to-disk assignment policy?
 - Common approach: Fixed data placement
 - Your data is on disk X, period!
 - For good reasons too: you bought it or you're paying more...
 - Fancy: Dynamic data placement
 - If some of your files are accessed a lot, the admin(or even system) may separate the "hot" files across multiple disks
 - In this scenario, entire files systems (or even files) are manually moved by the system admin to specific disks
 - Alternative: Disk striping
 - Stripe all of the data across all of the disks

7

Disk Striping

- Interleave data across multiple disks
 - Large file streaming can enjoy parallel transfers
 - High throughput requests can enjoy thorough load balancing
 - If blocks of hot files equally likely on all disks (really?)

File Foo:

Stripe:

stripe unit or block

8

Disk striping details

- How disk striping works
 - Break up total space into fixed-size stripe units
 - Distribute the stripe units among disks in round-robin
 - Compute location of block #B as follows
 - $disk\# = B \% N$ ($\% = \text{modulo}, N = \text{\#ofdisks}$)
 - $LBN\# = B / N$ (computes the LBN on given disk)

9

Now, What If A Disk Fails?

- In a JBOD (independent disk) system
 - one or more file systems lost
- In a striped system
 - a part of each file system lost
- Backups can help, but
 - backing up takes time and effort
 - backup doesn't help recover data lost during that day
 - Any data loss is a big deal to a bank or stock exchange

10

Tolerating and masking disk failures

- If a disk fails, it's data is gone
 - may be recoverable, but may not be
- To keep operating in face of failure
 - must have some kind of data redundancy
- Common forms of data redundancy
 - replication
 - erasure-correcting codes
 - error-correcting codes

11

Redundancy via replicas

- Two (or more) copies
 - mirroring, shadowing, duplexing, etc.
- Write both, read either

12

Mirroring & Striping

- Mirror to 2 virtual drives, where each virtual drive is really a set of striped drives
 - Provides reliability of mirroring
 - Provides striping for performance (with write update costs)

13

Implementing Disk Mirroring

- Mirroring can be done in either software or hardware
- Software solutions are available in most OS's
 - Windows2000, Linux, Solaris
- Hardware solutions
 - Could be done in Host Bus Adaptor(s)
 - Could be done in Disk Array Controller

14

Hardware vs. Software RAID

- Hardware RAID
 - Storage box you attach to computer
 - Same interface as single disk, but internally much more
 - Multiple disks
 - More complex controller
 - NVRAM (holding parity blocks)
- Software RAID
 - OS (device driver layer) treats multiple disks like a single disk
 - Software does all extra work
- Interface for both
 - Linear array of bytes, just like a single disk (but larger)

Lower Cost Data Redundancy

- Single failure protecting codes
 - general single-error-correcting code is overkill
 - General code finds error and fixes it
- Disk failures are self-identifying (a.k.a. erasures)
 - Don't have to find the error
- Fact: N-error-detecting code is also N-erasure-correcting
 - Error-detecting codes can't find an error, just know its there
 - But if you independently know where error is, allows repair
- Parity is single-disk-failure-correcting code
 - recall that parity is computed via XOR
 - it's like the low bit of the sum

16

Simplest approach: Parity Disk

- One extra disk
- All writes update parity disk
 - Potential bottleneck

17

Updating and using the parity

18

The parity disk bottleneck

- Reads go only to the data disks
 - But, hopefully load balanced across the disks
- All writes go to the parity disk
 - And, worse, usually result in Read-Modify-Write sequence
 - So, parity disk can easily be a bottleneck

Solution: Striping the Parity

- Removes parity disk bottleneck

The diagram shows five vertical cylinders representing disks. Each disk has four segments labeled A, B, C, and D. The parity blocks are labeled Ap, Bp, Cp, and Dp. In the first disk, A is white, B is white, C is white, and D is white. In the second disk, A is white, B is white, C is white, and D is white. In the third disk, A is white, B is white, C is white, and D is white. In the fourth disk, A is white, B is white, C is white, and D is white. In the fifth disk, Ap is white, B is white, C is white, and D is white.

Outline

- Using multiple disks
 - Why have multiple disks?
 - problem and approaches
- RAID levels and performance
- Estimating availability

RAID Taxonomy

- Redundant Array of Inexpensive Independent Disks
 - Constructed by UC-Berkeley researchers in late 80s (Garth)
- RAID 0 – Coarse-grained Striping with no redundancy
- RAID 1 – Mirroring of independent disks
- RAID 2 – Fine-grained data striping plus Hamming code disks
 - Uses Hamming codes to detect and correct multiple errors
 - Originally implemented when drives didn't always detect errors
 - Not used in real systems
- RAID 3 – Fine-grained data striping plus parity disk
- RAID 4 – Coarse-grained data striping plus parity disk
- RAID 5 – Coarse-grained data striping plus striped parity
- RAID 6 – Coarse-grained data striping plus 2 striped codes

RAID-0: Striping

- Stripe blocks across disks in a “chunk” size
 - How to pick a reasonable chunk size?

The diagram shows four vertical cylinders representing disks. Each disk has four segments labeled 0, 4, 8, 12; 1, 5, 9, 13; 2, 6, 10, 14; and 3, 7, 11, 15.

How to calculate where chunk # lives?
 Disk:
 Offset within disk:

RAID-0: Striping

The diagram shows four vertical cylinders representing disks. Each disk has four segments labeled 0, 4, 8, 12; 1, 5, 9, 13; 2, 6, 10, 14; and 3, 7, 11, 15.

- Evaluate for D disks
- Capacity: How much space is wasted?
- Performance: How much faster than 1 disk?
- Reliability: More or less reliable than 1 disk?

RAID-1: Mirroring

- Motivation: Handle disk failures
- Put copy (mirror or replica) of each chunk on another disk

Capacity:
Reliability:
Performance:

RAID-4: Parity

- Motivation: Improve capacity
- Idea: Allocate parity block to encode info about blocks
 - Parity checks all other blocks in stripe across other disks
- Parity block = XOR over others (gives "even" parity)
 - Example: 0 1 0 → Parity value?
- How do you recover from a failed disk?
 - Example: x 0 0 and parity of 1
 - What is the failed value?

RAID-4: Parity

- Capacity:
- Reliability:
- Performance:
 - Reads
 - Writes: How to update parity block?
 - Two different approaches
 - Small number of disks (or large write):
 - Large number of disks (or small write):
 - Parity disk is the bottleneck

RAID-5: Rotated Parity

Rotate location of parity across all disks

- Capacity:
- Reliability:
- Performance:
 - Reads:
 - Writes:
 - Still requires 4 I/Os per write, but not always to same parity disk

Comparison

	RAID-0	RAID-1	RAID-4	RAID-5
Capacity	N	$N/2$	$N - 1$	$N - 1$
Reliability	0	1 (for sure) $\frac{N}{2}$ (if lucky)	1	1
Throughput				
Sequential Read	$N \cdot S$	$(N/2) \cdot S$	$(N - 1) \cdot S$	$(N - 1) \cdot S$
Sequential Write	$N \cdot S$	$(N/2) \cdot S$	$(N - 1) \cdot S$	$(N - 1) \cdot S$
Random Read	$N \cdot R$	$N \cdot R$	$(N - 1) \cdot R$	$N \cdot R$
Random Write	$N \cdot R$	$(N/2) \cdot R$	$\frac{1}{2} \cdot R$	$\frac{N}{4} \cdot R$
Latency				
Read	D	D	D	D
Write	D	D	$2D$	$2D$

Table 38.7: RAID Capacity, Reliability, and Performance

Advanced Issues

- What happens if more than one fault?
 - Example: One disk fails plus "latent sector error" on another
 - RAID-5 cannot handle two faults
 - Solution: RAID-6 (e.g., RDP) Add multiple parity blocks
- Why is NVRAM useful?
 - Example: What if update 2, don't update P0 before power failure (or crash), and then disk 1 fails?
 - NVRAM solution: Use to store blocks updated in same stripe
 - If power failure, can replay all writes in NVRAM
 - Software RAID solution: Perform parity scrub over entire disk

RAID 6

- P+Q Redundancy
 - Protects against multiple failures using Reed-Solomon codes
 - Uses 2 "parity" disks
 - P is parity
 - Q is a second code
 - It's two equations with two unknowns, just make "biggerbits"
 - Group bits into "nibbles" and add different coefficients to each equation (two independent equations in two unknowns)
- Similar to parity striping
 - De-clusters both sets of parity across all drives
 - For small writes, requires 6 I/Os
 - Read old data, old parity1, old parity2
 - Write new data, new parity1, new parity2

31

The Disk Array Matrix

	Independent	Fine Striping	Course Striping
None	JBOD		RAID0
Replication	Mirroring RAID1		RAID0+1
Parity Disk		RAID3	RAID4
Striped Parity	Gray90		RAID5

32

Outline

- Using multiple disks
 - Why have multiple disks?
 - problem and approaches
- RAID levels and performance
- Estimating availability

33

Sidebar: Availability metric

- Fraction of time that server is able to handle requests
 - Computed from MTBF and MTTR (Mean Time To Repair)

$$\text{Availability} = \frac{\text{MTBF}}{\text{MTBF} + \text{MTTR}}$$

34

How often are failures?

- MTBF (Mean Time Between Failures)
 - $\text{MTBF}_{\text{disk}} \sim 1,200,00$ hours (~136 years, <1% per year)
- $\text{MTBF}_{\text{multi-disk system}} = \text{mean time to first disk failure}$
 - which is $\text{MTBF}_{\text{disk}} / (\text{number of disks})$
 - For a striped array of 200 drives
 - $\text{MTBF}_{\text{array}} = 136 \text{ years} / 200 \text{ drives} = 0.65 \text{ years}$

35

Back to Mean Time To Data Loss (MTTDL)

- MTBF (Mean Time Between Failures)
 - $\text{MTBF}_{\text{disk}} \sim 1,200,00$ hours (~136 years, <1% per year)
- $\text{MTBF}_{\text{multi-disk system}} = \text{mean time to first disk failure}$
 - which is $\text{MTBF}_{\text{disk}} / (\text{number of disks})$
 - For a striped array of 200 drives
 - $\text{MTBF}_{\text{array}} = 136 \text{ years} / 200 \text{ drives} = 0.65 \text{ years}$

36

Reliability without rebuild

- 200 data drives with MTBF_{drive}
 - $MTTDL_{array} = MTBF_{drive} / 200$
- Add 200 drives and do mirroring
 - $MTBF_{pair} = (MTBF_{drive} / 2) + MTBF_{drive} = 1.5 * MTBF_{drive}$
 - $MTTDL_{array} = MTBF_{pair} / 200 = MTBF_{drive} / 133$
- Add 50 drives, each with parity across 4 data disks
 - $MTBF_{set} = (MTBF_{drive} / 5) + (MTBF_{drive} / 4) = 0.45 * MTBF_{drive}$
 - $MTTDL_{array} = MTBF_{set} / 50 = MTBF_{drive} / 111$

37

Rebuild: restoring redundancy after failure

- After a drive failure
 - data is still available for access
 - but, a second failure is BAD
- So, should reconstruct the data onto a new drive
 - on-line spares are common features of high-end disk arrays
 - reduce time to start rebuild
 - must balance rebuild rate with foreground performance impact
 - a performance vs. reliability trade-offs
- How data is reconstructed
 - Mirroring: just read good copy
 - Parity: read all remaining drives (including parity) and compute

38

Reliability consequences of adding rebuild

- No data loss, if fast enough
 - That is, if first failure fixed before second one happens
- New math is...
 - $MTTDL_{array} = MTBF_{firstdrive} * (1 / \text{prob of } 2^{nd} \text{ failure before repair})$
 - ... which is $MTTR_{drive} / MTBF_{seconddrive}$
- For mirroring
 - $MTBF_{pair} = (MTBF_{drive} / 2) * (MTBF_{drive} / MTTR_{drive})$
- For 5-disk parity-protected arrays
 - $MTBF_{set} = (MTBF_{drive} / 5) * ((MTBF_{drive} / 4) / MTTR_{drive})$

39

Three modes of operation

- Normal mode
 - everything working; maximum efficiency
- Degraded mode
 - some disk unavailable
 - must use degraded mode operations
- Rebuild mode
 - reconstructing lost disk's contents onto spare
 - degraded mode operations plus competition with rebuild

40

Mechanics of rebuild

- Background process
 - use degraded mode read to reconstruct data
 - then, write it to replacement disk
- Implementation issues
 - Interference with foreground activity and controlling rate
 - Rebuild is important for reliability
 - Foreground activity is important for performance
 - Using the rebuilt disk
 - For rebuilt part, reads can use replacement disk
 - Must balance performance benefit with rebuild interference

41

Conclusions

- RAID turns multiple disks into a larger, faster, more reliable disk
- RAID-0: Striping
 - Good when performance and capacity really matter, but reliability doesn't
- RAID-1: Mirroring
 - Good when reliability and write performance matter, but capacity (cost) doesn't
- RAID-5: Rotating Parity
 - Good when capacity and cost matter or workload is read-mostly
 - Good compromise choice