# Redundancy, Fault Tolerance & RAID

Rishabh Iyer

12-11-2020

# Faults, Errors and Failures

o Fault

  ❖ Underlying defect

  ❖ May or may not cause problems

o Error

  ❖ Active fault

o Failure

  ❖ Error at an interface between modules

# Standard Jargon (Recap)

o MTTF

o MTTR

o MTBF

o Availability

o Bathtub curve

# Reacting to errors within a module

o Do nothing

o Fail-fast

o Fail-safe

o Fail soft

# RAID

o Motivation:

❖ Fault Tolerance

❖ Optionally throughput

o Idea:

❖ Use cheap commodity disks

❖ Split the array in **reliability groups**

❖ Use extra **check disks**

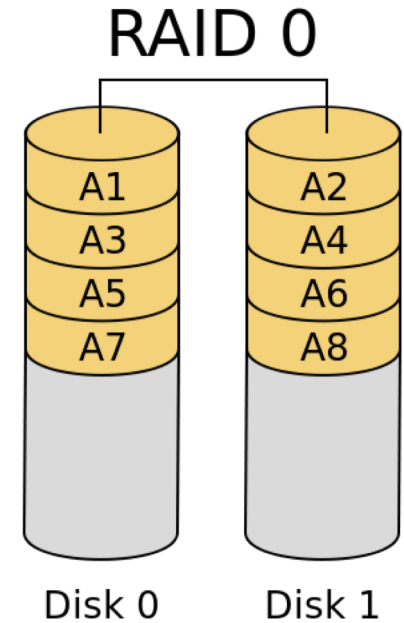# Why do we need redundancy?

o A single disk has an acceptable MTTF

$$MTTF\ of\ a\ Disk\ Array = \frac{MTTF\ of\ a\ Single\ Disk}{Number\ of\ Disks\ in\ the\ Array}$$

# MTTF of RAID?

○ $MTTF_{Group} = \dfrac{MTTF_{Disk}}{G+C} * \dfrac{1}{\substack{Probability\ of\ another\ failure\ in\ the\ group \\ before\ repairing\ the\ first\ one}}$

○ $P(another\ failure) = \dfrac{MTTR}{MTTF_{Disk}/_{(G+C-1)}}$
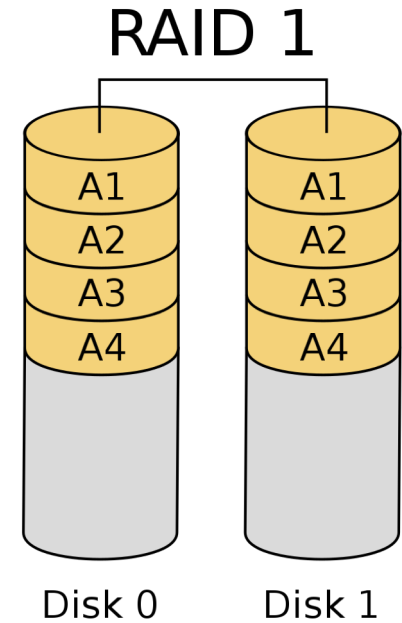
○ $MTTF_{RAID} = \dfrac{MTTF_{Group}}{n_G}$

# RAID 0

o More read/write throughput
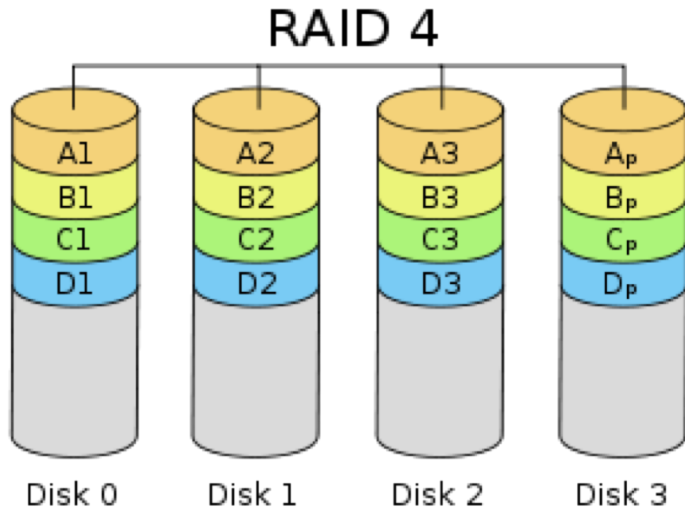
o No redundancy

o No fault tolerance

RAID 0



Disk 0    Disk 1

# RAID 1

o Increased read throughput

o Almost the same write throughput

o Half of available storage is used

RAID 1

| Disk 0 | Disk 1 |
|--------|--------|
| A1 | A1 |
| A2 | A2 |
| A3 | A3 |
| A4 | A4 |

# RAID 4



RAID 4

A1 A2 A3 Ap
B1 B2 B3 Bp
C1 C2 C3 Cp
D1 D2 D3 Dp

Disk 0  Disk 1  Disk 2  Disk 3



**RAID 4: Small Writes**

Small Write Algorithm

1 Logical Write = 2 Physical Reads + 2 Physical Writes

D0'    D0    D1    D2    D3    P

new data    old data (1. Read)    old parity (2. Read)

+ XOR

+ XOR

(3. Write)    (4. Write)
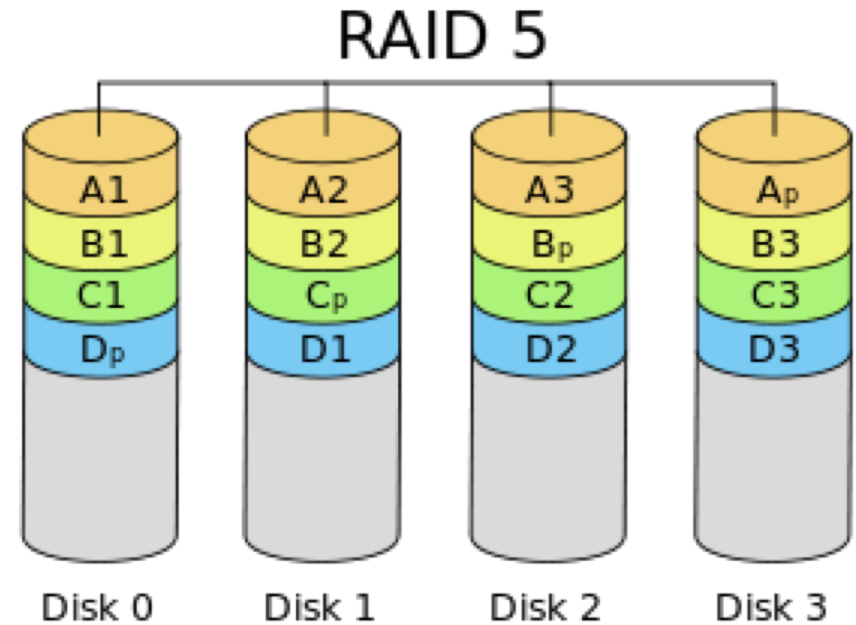
D0'    D1    D2    D3    P'

CSCE430/830    Disk Storage Systems: RAID

o Parallel Reads

o Check disk is a write bottleneck

# RAID 5

o Solves the RAID 4 write bottleneck

o Widely used solution

o Let's see that in action



RAID 5

| Disk 0 | Disk 1 | Disk 2 | Disk 3 |
|--------|--------|--------|--------|
| A1 | A2 | A3 | $A_p$ |
| B1 | B2 | $B_p$ | B3 |
| C1 | $C_p$ | C2 | C3 |
| $D_p$ | D1 | D2 | D3 |

# RAID 5 Example

|         | DISK0 | DISK1 | DISK2 | DISK3 |
|---------|-------|-------|-------|-------|
| STRIPE0 | 0100  | 0101  | 0010  | ?     |
| STRIPE1 | 0010  | 0000  | ?     | 0100  |
| STRIPE2 | 0011  | ?     | 1010  | 1000  |
| STRIPE3 | ?     | 0001  | 1101  | 1010  |

o Complete the parity entries

# RAID 5 Example

| | DISK0 | DISK1 | DISK2 | DISK3 |
|---|---|---|---|---|
| STRIPE0 | 0100 | 0101 | 0010 | 0011 |
| STRIPE1 | 0010 | 0000 | 0110 | 0100 |
| STRIPE2 | 0011 | 0001 | 1010 | 1000 |
| STRIPE3 | 0110 | 0001 | 1101 | 1010 |

o STRIPE0,DISK3 = 0100 XOR 0101 XOR 0010 = 0011
o STRIPE1,DISK2 = 0010 XOR 0000 XOR 0100 = 0110
o STRIPE2,DISK1 = 0011 XOR 1010 XOR 1000 = 0001
o STRIPE3,DISK0 = 0001 XOR 1101 XOR 1010 = 0110

# RAID 5 Example - Writes

|  | DISK0 | DISK1 | DISK2 | DISK3 |
|---|---|---|---|---|
| STRIPE0 | 0100 | 0101 | 0010 | 0011 |
| STRIPE1 | 0010 | 0000 | 0110 | 0100 |
| STRIPE2 | 0011 | 0001 | 1010 | 1000 |
| STRIPE3 | 0110 | 0001 | 1101 | 1010 |

Modifying STRIPE0 in DISK2 to 1101. Outline the steps. How many reads and writes will you need?

# RAID 5 Example - Writes

| | DISK0 | DISK1 | DISK2 | DISK3 |
|---|---|---|---|---|
| STRIPE0 | 0100 | 0101 | ~~0010~~ 1101 | ~~0011~~ 1100 |
| STRIPE1 | 0010 | 0000 | 0110 | 0100 |
| STRIPE2 | 0011 | 0001 | 1010 | 1000 |
| STRIPE3 | 0110 | 0001 | 1101 | 1010 |

STRIPE0, DISK3 = 0010 XOR 1101 XOR 0011 = 1100

# RAID 5 Example – Disk Failure

| | DISK0 | DISK1 | DISK2 | DISK3 |
|---|---|---|---|---|
| STRIPE0 | 0100 | 0101 | 1110 | 1100 |
| STRIPE1 | 0010 | 0000 | 0110 | 0100 |
| STRIPE2 | 0011 | 0001 | 1010 | 1000 |
| STRIPE3 | 0110 | 0001 | 1111 | 1010 |

Disk 2 died. How can the RAID controller serve a read requests for STRIPE0 for DISK2?
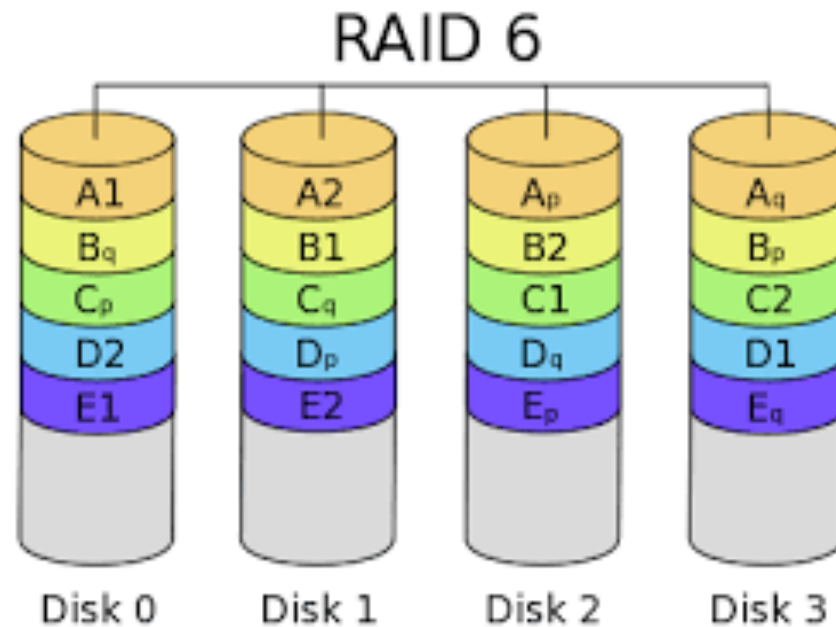
# RAID 5 Example – Disk Failure

| | DISK0 | DISK1 | DISK2 | DISK3 |
|---|---|---|---|---|
| STRIPE0 | 0100 | 0101 | 110 | 1100 |
| STRIPE1 | 0010 | 0000 | 0110 | 0100 |
| STRIPE2 | 0011 | 0001 | 1010 | 1000 |
| STRIPE3 | 0110 | 0001 | 1 1 | 1010 |

STRIPE0,DISK2 = 0100 XOR 0101 XOR 1100 = 1101

# RAID 6

o Double Parity

o Can tolerate 2 failures in a group

# Changing times

o Is RAID still the go-to for fault tolerant storage?

&#10070; Scale

&#10070; Evolving speeds/bottlenecks

# RAID for SSDs?

o SSDs have higher MTTF

  ❖ Check Flash Translation Layer

o No performance argument (for RAID 1)

  ❖ High throughput SSDs

  ❖ Internal parallelism

o TRIM command for garbage collection

# RAID vs E2E Argument

o The E2E argument argues against providing reliability in lower layers of the stack because ultimately it is the endpoints that know what they want. RAID provides a lot of reliability at the hardware device level. Is it a violation of the E2E argument? Why are both still considered good design guidelines/ well designed systems?

# RAID vs E2E Argument

o RAID is a layer above the hardware device, software RAID implementations exist

o Sometimes errors are better masked at a lower layer, because it understands the error better

o All applications have the same interface to disk, so which layer you put the reliability in does not matter.

# Summary

o Redundancy is the de-facto method used to achieve fault tolerance

o RAID:

   ❖ Required increased throughput from an array of disks

   ❖ Used redundancy to ensure improved reliability

   ❖ Decreasing performance, storage overheads with increasing RAID levels

# The RAMCloud Storage System

Lei Yan
(Slides partially adopted from Marios Kogias)

# RAMCloud

Serving large data (100s of TB) with low latency (5 - 10us)

- Store data in DRAM
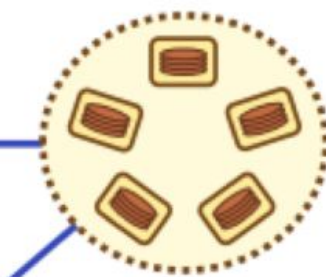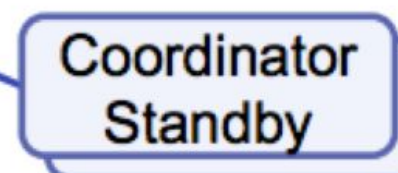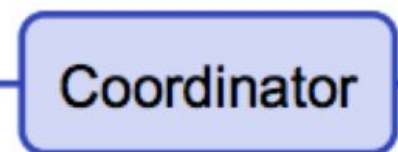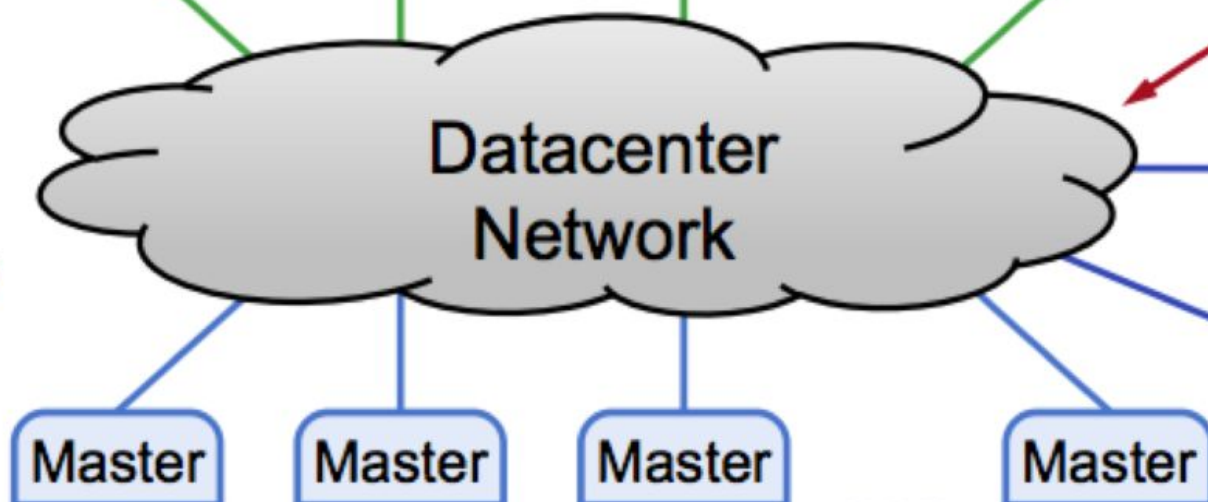- Distributed

# RAMCloud Architecture



**1000 – 100,000 Application Servers**

Appl. / Library  Appl. / Library  Appl. / Library  ...  Appl. / Library

High-speed networking:
- 5 µs round-trip
- Full bisection bandwidth

Datacenter Network

Coordinator

Coordinator Standby

External Storage (ZooKeeper)

Commodity Servers

Master  Master  Master  Master

64-256 GB per server

**1000 – 10,000 Storage Servers**

# Durability and High Availability

- Durability

  No loss of committed updates

  At least as durable as disk-based systems (assuming every server's data is replicated in the disks of another two servers)

- High Availability

  < 5 sec down time

# Failures

- Server Crashes
- Data center power Failures
- Malicious acts, e.g., hackers deleting data

# Design Exercise

**Goals:**

- **Durability:** Probability of losing data at a given time point should be as low as disk-based systems

- **High availability:** < 5 sec down time

- **Low latency:** 5-10us

- **High throughput:** Highest possible while achieving the above goals

**Assumptions:**

- 100+ Servers

- 100MB/s disk bandwidth

- 1GB/s network bandwidth per server

- 50GB DRAM per server

- Probability of server crash = Probability of disk failure

- Probability that a disk and one replica fail < Probability of power failure

# Durability

First try: Replicate data in DRAM

- Does it work?

Second try: Replicate data in Disk

- Does it work?

# Fast Recovery

First try: Replicate entire data on one primary backup, restore entire data on a single new server

- What is the bottleneck?

    Disk accesses

- Recovery time?

    100MB/s for 50GB -> 500 sec

# Fast Recovery

Second try: Partition replicas among 100 primary backup servers, restore data on a single new server

- What is the bottleneck?

  Network bandwidth

- Recovery time?

  1GB/s for 50GB -> 50 sec

# Fast Recovery

Last try: Partition replicas among 100 primary backup servers, restore data on the backups

- Recovery time?

100 * 100MB/s for 50GB -> 5 sec

# High Throughput

First try: Store replicas on disk as hash table

- Random write access to disk, low write throughput

Second try: Store replicas as logs on disk, every update to the data is appended to the log

- Sequential write access to disk

# Low Latency

When should storage servers reply to the client that the write succeeds?

*How to handle other failures within the storage servers, e.g., memory corruption, after we detecting them?*

*How to handle these failures with the techniques available in the system we just designed after we detecting them?*

- Hint: Recall George's one-hammer-for-all approach of handling failures with rebooting

- Failure promotion: promote the failures to a crash, and reusing techniques used for handling crashes.

- Reduce system complexity: BL's "make it simple"

- Assumption: The failures happens rare, so the cost of crash recovery is amortized.

# Full RAMCloud paper:

https://dl.acm.org/doi/pdf/10.1145/2806887

RAMCloud full paper is a good read, gives a lot of details of the system, covers many aspects of system design, scalability, fault tolerance and so on...
A lot of analysis of various trade-offs and how they affect the design of the system. Definitely have a look if you are interested in systems