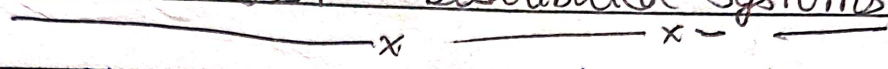


MIT 6.824 Distributed Systems



Infrastructure

- Storage
- Communication
- Computation

Course Contents

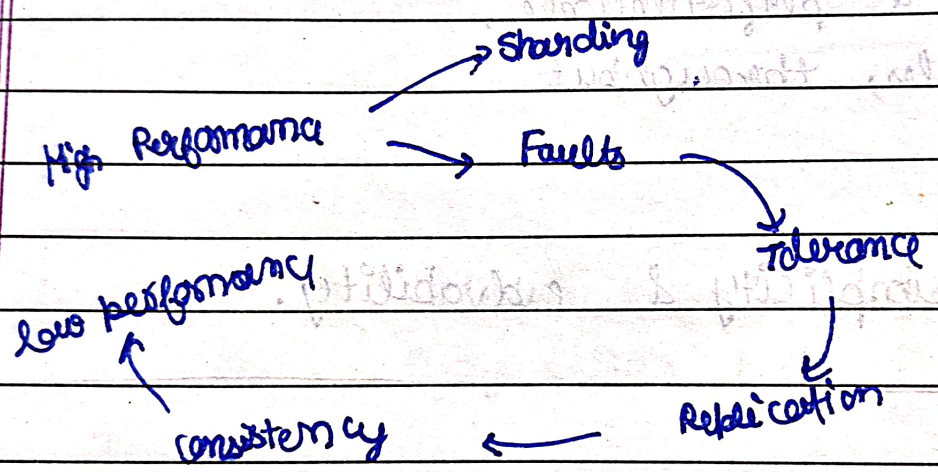
- RPC
- Threads
- Concurrency Control
- Performance
  - ↳ Scalability
  - ↳ Fault Tolerance
  - ↳ Consistency

## ⇒ RPC & Threads.

Why threads?

- I/O concurrency
- Parallelism
- Convenience → ping, background tasks.

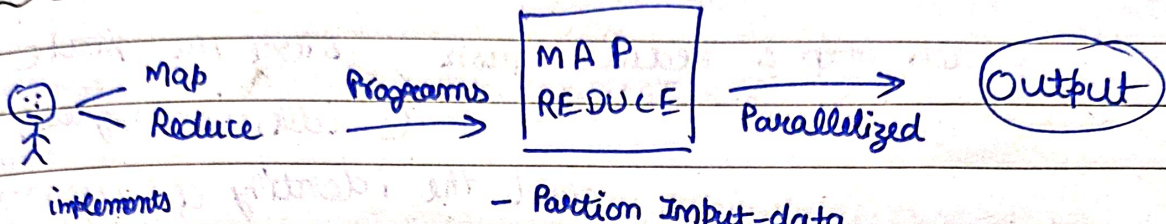
## ⇒ OIFS



(2004)

# MapReduce: Simplified Data Processing on Large Clusters

## ⇒ Abstract



- Partition Input-data
- scheduling program's execution across a set of machines.
- Handles Machine failures.
- Manages Inter-Machine communication.
- High Performance

Input

↓ MAP

Set of Intermediate Key/Value pairs.

↓ REDUCE

Output which shares same key. <sup>intermediate</sup>

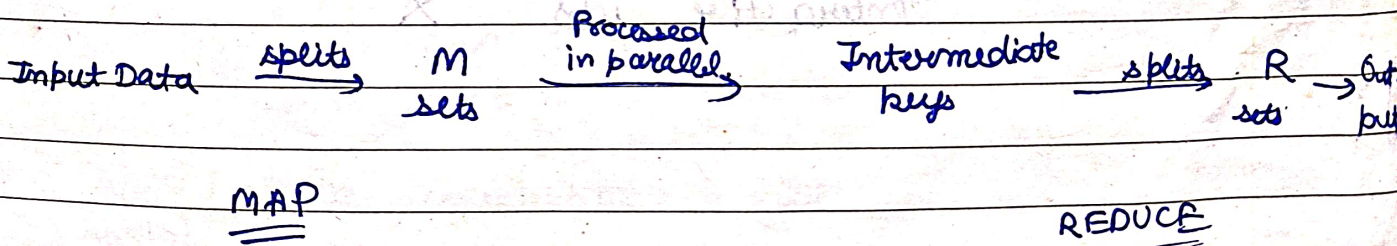
Hence,

map  $(k_1, v_1) \rightarrow \text{list}(k_2, v_2)$

Reduce  $(k_2, \text{list}(v_2)) \rightarrow \text{list}(v_2)$

e.g. Distributed GREP, count of URL access frequency, Inverted Index, Distributed Sort, etc.

## ⇒ Implementation



( users can specify R partitioning functions )

The implementation consists of the following:

① Master Data Structure

for each map & reduce task, store the state (idle, in-progress, or completed) and the identity of worker machines for non-idle tasks.

② Fault Tolerance

Worker (Just Reset)

Master (Can make backup to manage the data structure)

③ tries Locality, & Task Co-locality:

⇒ Refinements

1. I/O types
2. Monitoring Service

Notes  
MapReduce can solve problems that are trivial but big. Might be difficult for hard problems.

i.e. batch processing jobs ✓  
Interactive jobs X

Page \_\_\_\_\_

# The Google File System (2003)

## ⇒ Abstract

scalable distributed file system over an inexpensive commodity hardware.

Traditional File Systems + Google's observations (extensions) distributed for their workload

(Performance + scalability + reliability + availability)

## ⇒ Introduction

Questioning existing design choices

- Component failures

- ↳ They need to be considered as normal, since it is always fails in thousands of machines

- ↳ Problems can occur by (system) OS bugs, apple's bugs, human errors, disk fails, networking, power supplies.

- ↳ Therefore, constant monitoring, error detection, fault tolerance, & automatic recovery.

- Huge file sizes

Moving

- ↳ billions of KB files are unwieldy, since multi

- ↳ GB files are common. I/O operation & block sizes have to be revisited.

- Write / update considerations in files

(Think of OASIS)

- ↳ most files are append-only, over-writing existing data are practically non-existent.

- ↳ once written, the files are only read, and often only sequentially.

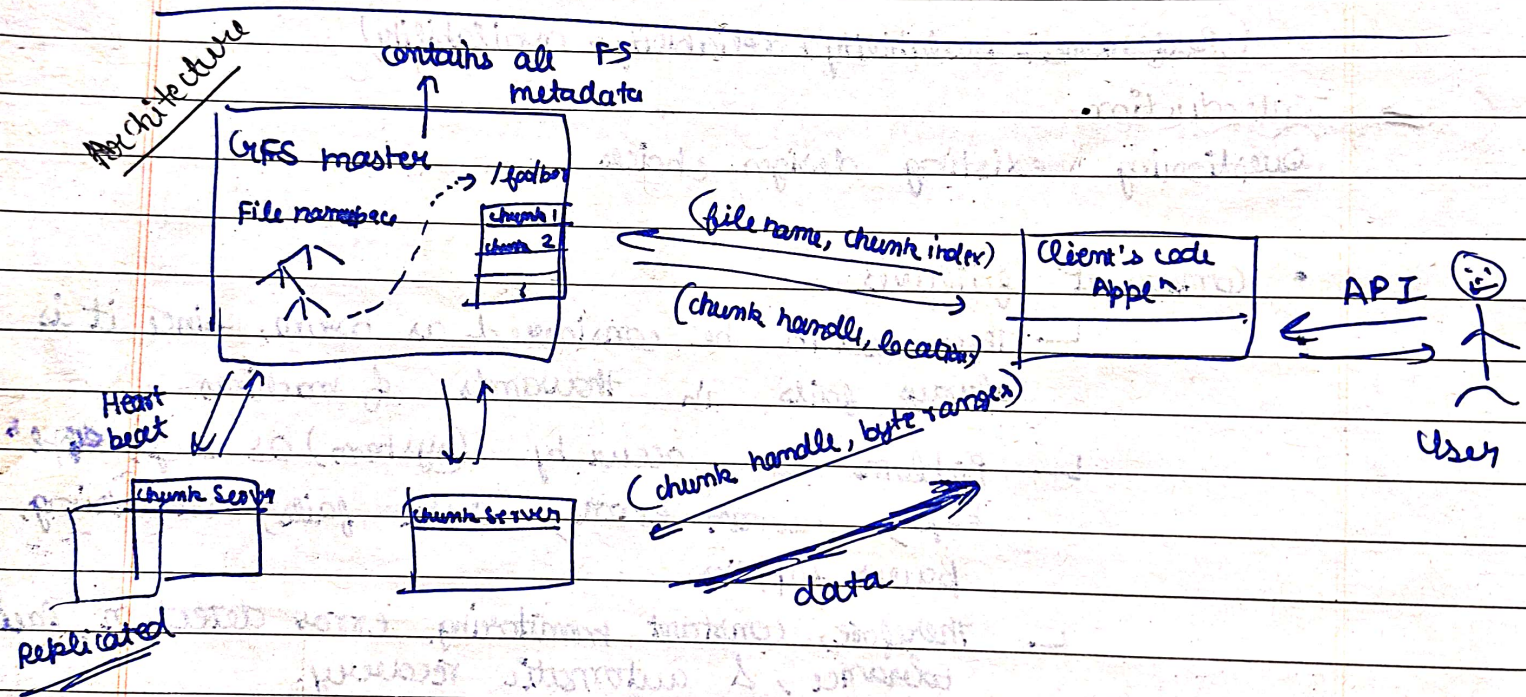
The workloads are often analytics, logs etc.

Co-designing appl<sup>n</sup> & the FS API

Loading atomic append operation

⇒ Design Overview

Operations Supported → Create, delete, open, close, read & write  
+  
Snapshot & record append



FS metadata includes:

- Namespace
- Access Control Information
- Mapping from files to chunks
- Current location of chunks

Operations:

- Chunk-lease mgng.
- garbage collection
- Chunk Migration

(2004)  
Clusters

→ Chaos

→ Single master

- Minimize its involvement in read & writes

→ chunk Size

- 64 MB
- Reduces Metadata on the server, which allows to store it in memory.
- Reduces Request information.

→ Metadata

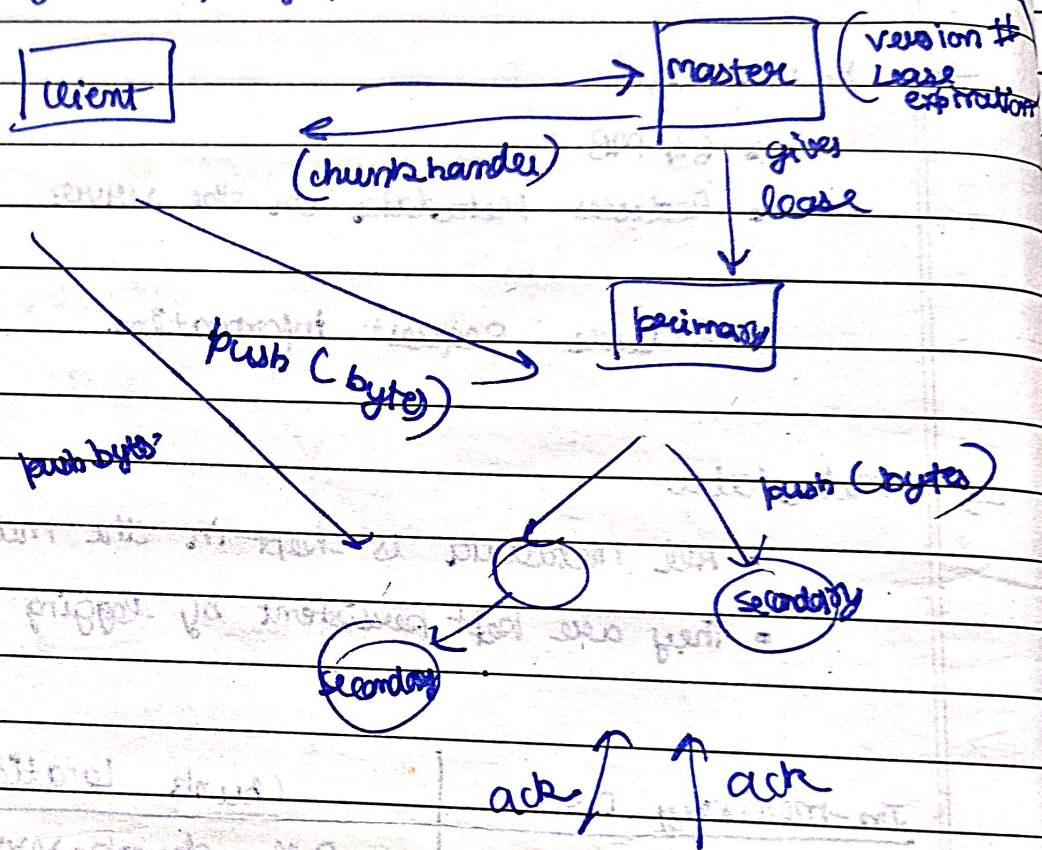
- All metadata is kept in the master's memory
- They are kept persistent by logging mutations to an operation log.

In-memory DS	Chunk Locations	Checkpoints + Operation Log
1. Scanning periodically	1. Polls chunk-servers at startup for locations.	1. Record of critical metadata changes.
2. Fast & efficient operations	2. Can update using heart-beat messages.	2. Persistent & replicated to remote machines.
	3. Also by periodically asking.	3. Also, master creates a checkpoint whenever log grows beyond a certain size.

→ Consistency Model

⇒ Operations  
record appends

append(filename, bytes)



After ack, the client send write request & BOOM

after all "Yes", "success" → done

else → not done

⊛ This still cause in-consistency, some may append & some not.

# Design of Practical System for Fault-Tolerant Machines (2010)

Abstract (works for fail-stop faults)  
Fault-Tolerant VMs via: primary/backup replication approach.

Practical Issues + Design Choices + Implementation & Alternatives

## Introduction

To replicate, we have two approaches:

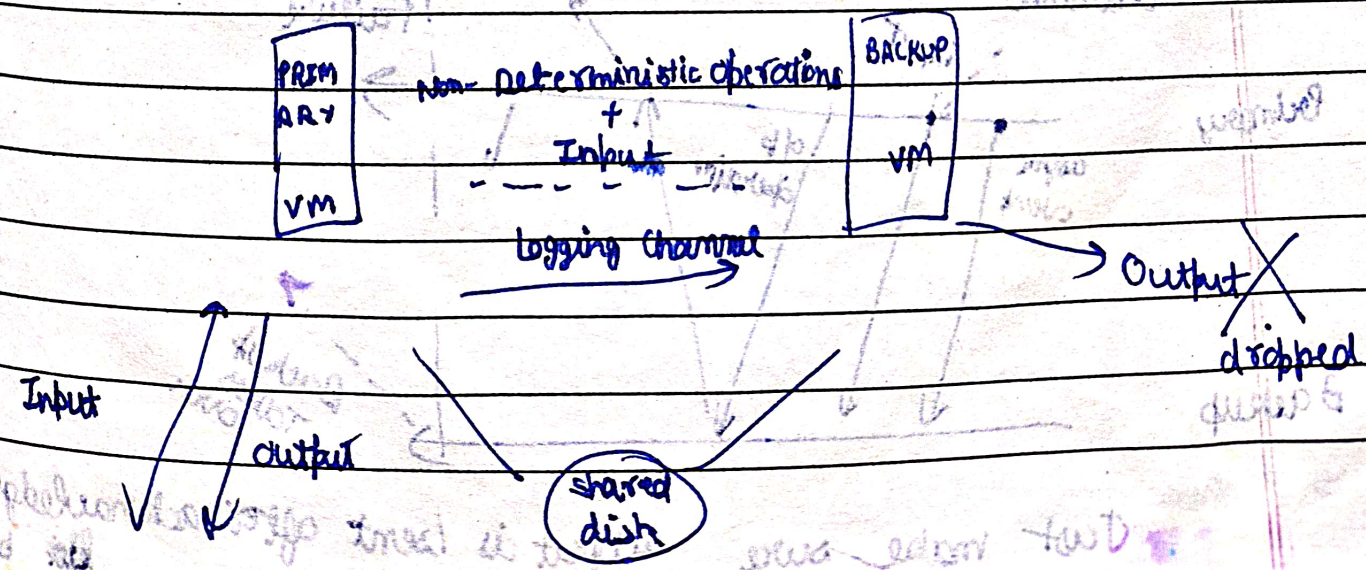
- i) constantly send (all changes to the state of primary i.e. CPU, memory, I/O) to backup continuously.

Bandwidth ↑

Just send some inputs requests to both the server.

Non-deterministic operations such as clock, system interrupts can create trouble.

A VM running on top of hypervisor can implement this approach.



## 1. Deterministic Replay Implementation

There are broad set of inputs, including network packets, disk reads, and input from the keyboard & mouse.  
e.g. Virtual Interrupts, Reading clock cycle. (Non-deterministic)

### Challenges

- Correctly capture all the input & non-determinism
- Correctly applying the inputs & non-determinism
- Doing so in a manner without degrading performance.

→ All inputs to primary VM are recorded & all possible non-determinism associated with the VM execution in a stream of log entries written to a log file.

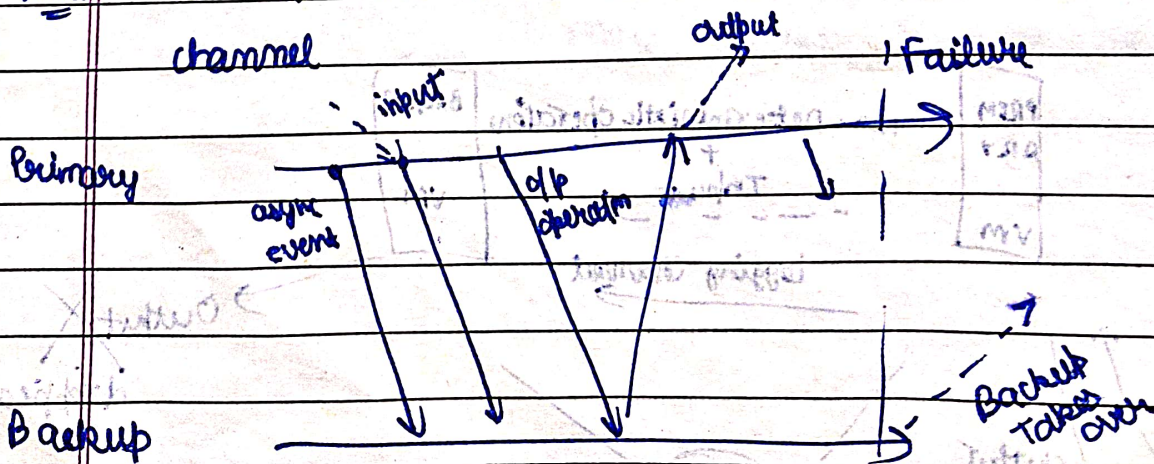
→ For non-deterministic operations, sufficient information is logged as well.

For events & interrupts, efficient event recording & event delivery mechanism is implemented.

## 2. FT Protocol

To ensure no data is lost when primary fails.

Note: The logs are not written to disk, instead sent via logging channel



Just make sure output is sent after acknowledging from backup.

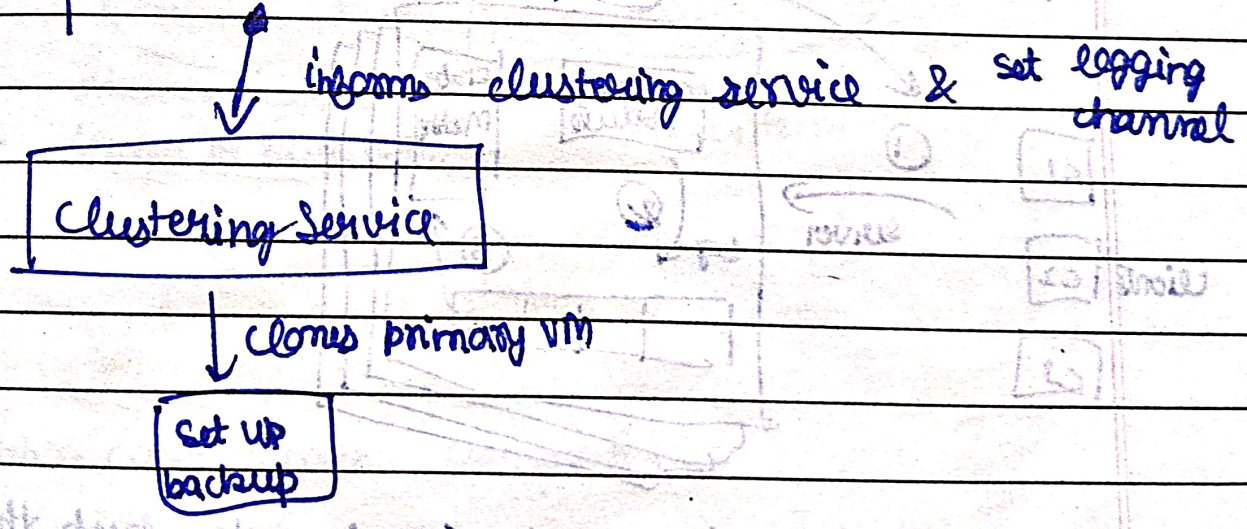
### 3. Detecting & Responding to Failure

Heartbeat + Monitoring  
Log Traffic

To avoid split brain, it uses an atomic set & lock operation on shared storage. Thus, avoids it.

If backup dies, primary will leave recording mode

If primary dies, backup waits until log entry is consumed  
get a lock on shared storage.  
gets another backup VM up



Raft : An understandable Consensus (2014) Algorithm

Raft → Consensus Algorithm

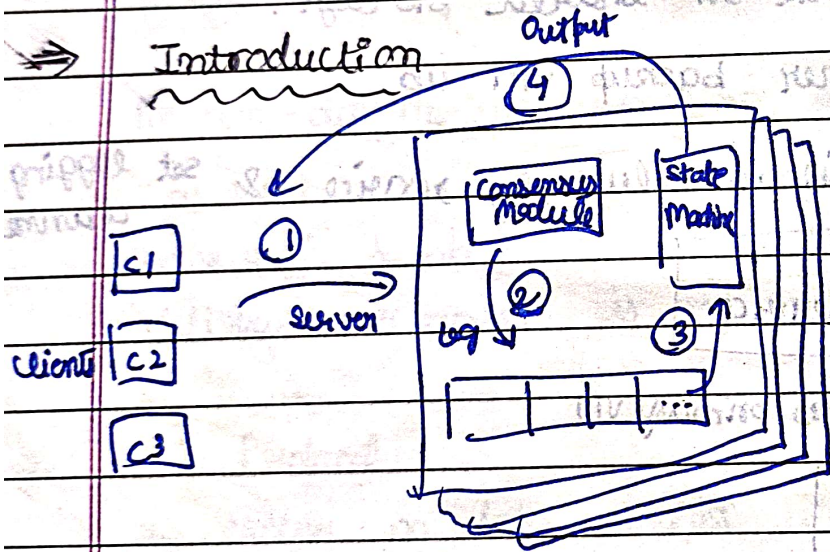
⇒ Abstract (works only on majority)

An easier version of Paxos  
It separates the key elements of the consensus, such as

- leader election
- log Replication
- safety

Understandability  
+  
decomposition  
+  
State Reduction

⇒ Introduction



The consensus algorithm's job is to keep the replicated log consistent.

They (Consensus algo.s) have the following properties

- They ensure safety (never incorrect results) under all non-Byzantine conditions, including network delays, partitions, and packet loss.
- They are fully functional as long as majority of servers are operational.
- They don't depend on timing to ensure consistency of the logs.

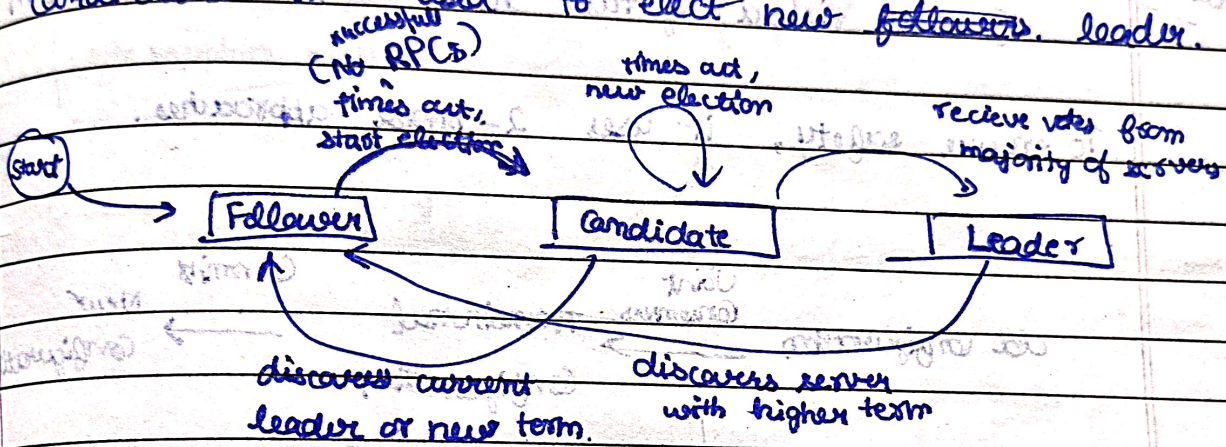
# Algorithm

Notes

A lot of servers (typically 5)  
 Each server can be:  
 - Leader  
 - Follower  
 - Candidate

In normal state, 1 is leader and rest are followers.  
 leader does lot of things

- Followers respond to requests from leaders & candidates
- Candidates are used to elect new followers, leader.



## Server States

Rep. servers communicate using RPCs

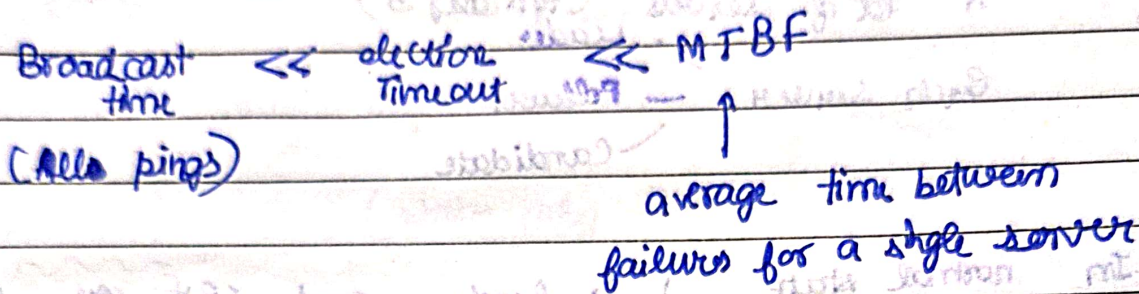
- ↳ Request vote RPCs
- ↳ Append entries RPCs

\* leader only commits an entry once it has been replicated to majority of servers.

\* During election, those who will not get elected, who have lower commit index.

Notes: Every term has at least one leader

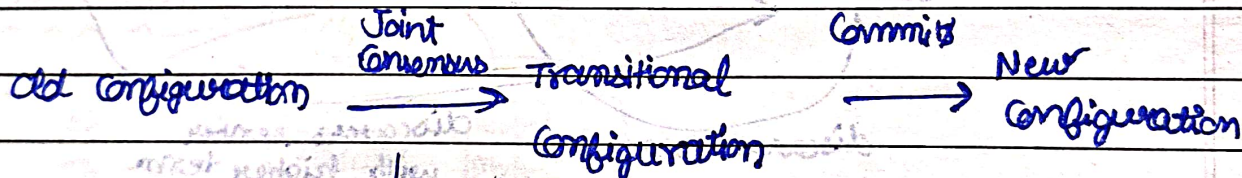
Another requirement for RAFT is timing.



~~Cluster Membership Change~~

Can be risked of two leaders for the same term.

To ensure safety, it uses 2-phase approaches.



i) Log entries are replicated to all servers in both configurations.

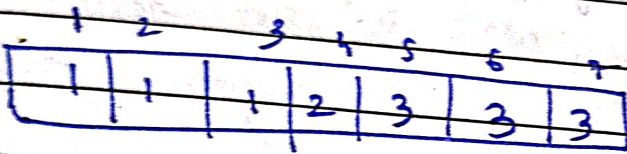
ii) Any server from either configuration

iii) Agreement (for elections & entry committed)

requires separate majority from both the old & new configurations.

Log Compaction

Snapshot discarded. current system state & log up to that point is



log index before.

snapshot

last included index: 5  
last included term: 3  
state machine state:  
=



after

← committed →

YAHOO

# Zookeeper (2010)

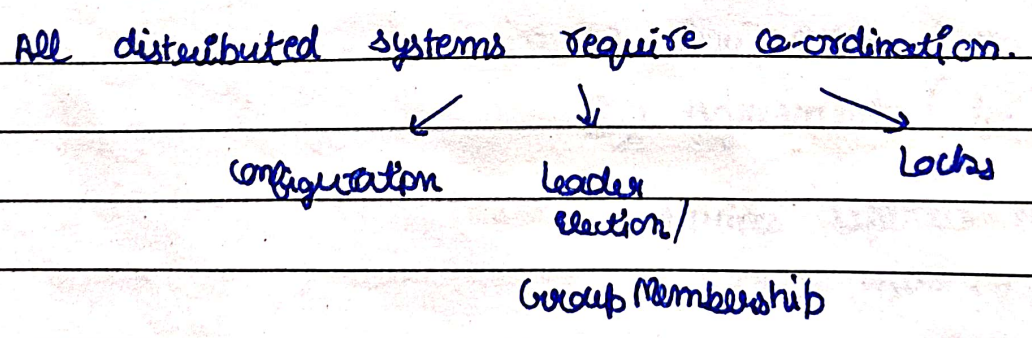
zookeeper → <sup>Service to</sup> Co-ordinating process of distributed systems / applications

## ⇒ Abstract

↓  
Elements of Group Messaging + shared registers + distributed locks service  
in a replicated centralized service.

- It has wait-free aspects with an event-driven mechanism to provide powerful co-ordination service.
- FIFO requests & linearizability for all requests.

## ⇒ Introduction



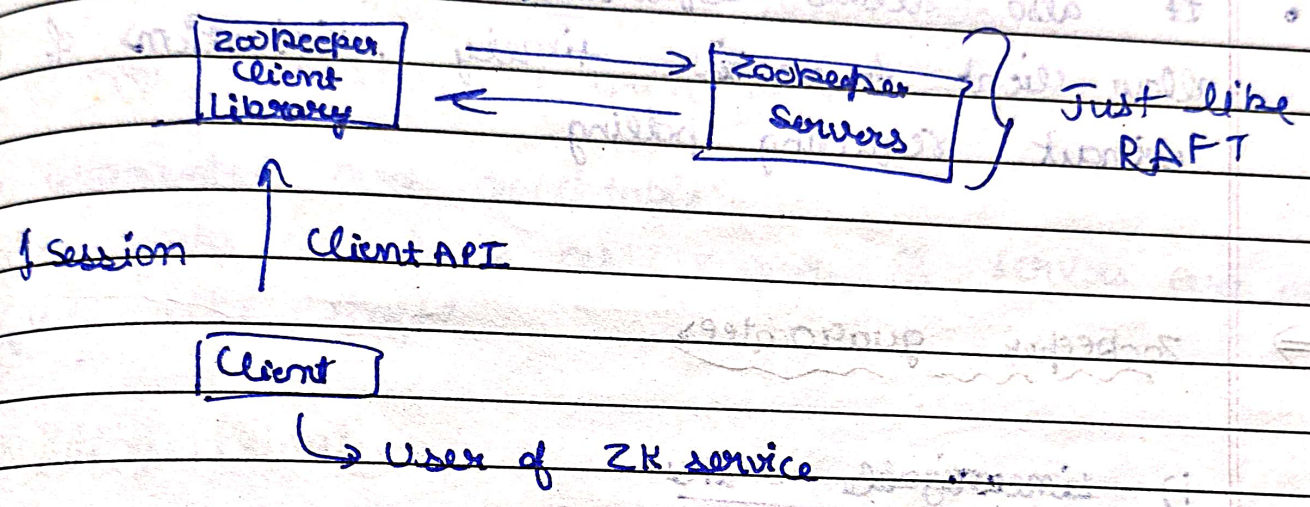
There can be several services for each of different co-ordination services but...

- Zookeeper implemented a co-ordination kernel and

exposed an API that enables appl developers to implement their own primitives.

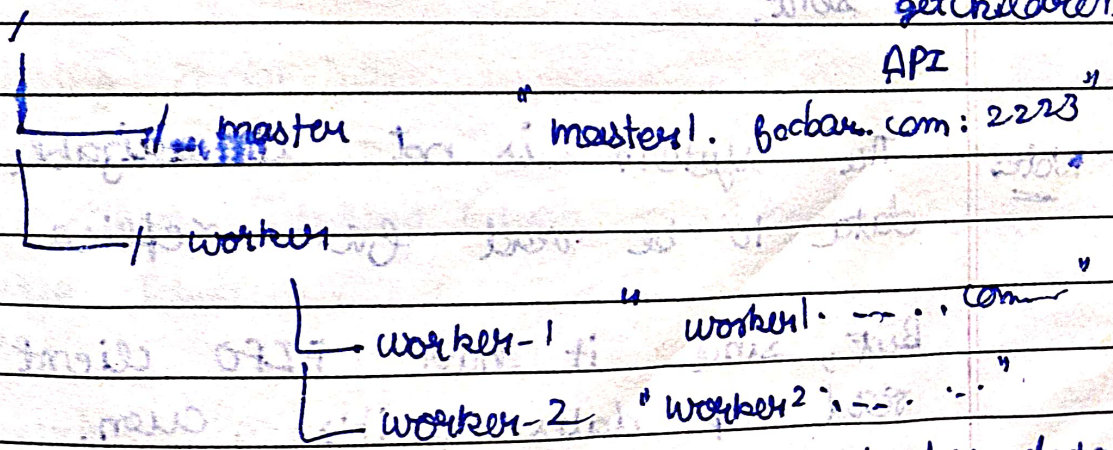
It also avoided locks, and used word-free data objects organized hierarchically as in file systems.

Zookeeper Service



zookeeper instead implementing a more general purpose system, it provides a file system.

Replace concept of "file" with "znode" and provides create, delete, exists, setData, getData, getChildren.



znode is an in-memory data node of zookeeper data.

There are two types of znodes, clients can create:

- i) Regular  
Persistent, creation & deletion explicit
- ii) ephemeral  
let system remove automatically

It also allows sequential flag & watches to allow client to receive timely notifications of changes without requiring polling.

⇒ Zookeeper guarantees

i) Linearizable writes

All ~~any~~ request that update the state of zookeeper are serializable.

ii) FIFO client order

All requests from a given client that were sent are executed the same way.

Notes

The system is not linearizable, since it allows stale data to be read from replicas.

But, since it have FIFO client order, client will read up dates of it's own.

It also allows service to scale linearly

$\frac{1}{N} \rightarrow N \times$  performance  $\rightarrow$  scale better  $\text{☺}$

$\uparrow$   
only for read-heavy appl<sup>n</sup> though.

This can be explained via

Leader config. update + Notification

"sync" operation (write before read)

So, guarantees are acceptable.

This works as long as majority of servers are alive  
(just like Raft)

### Building Primitives

Zookeeper locks cool IMO.

It can do

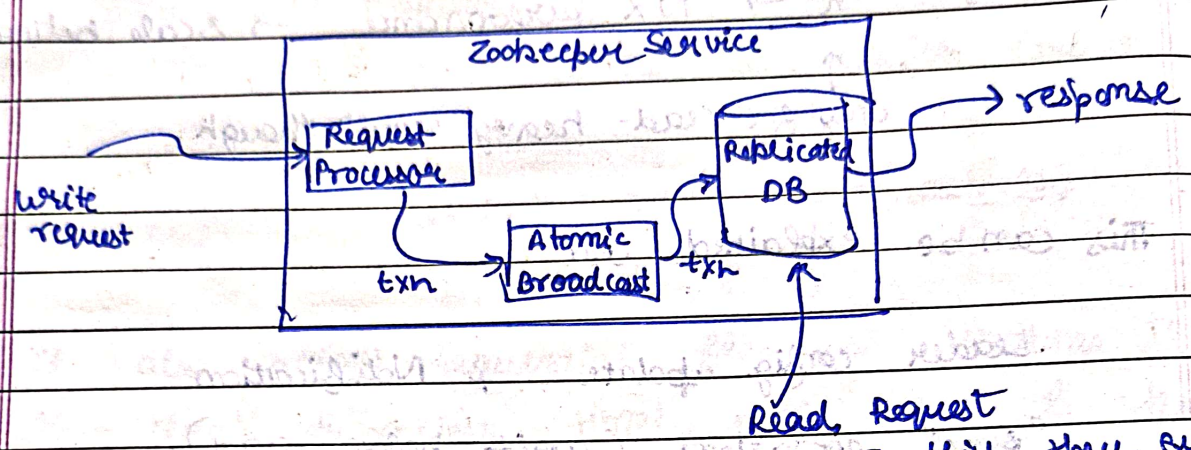
- i) Group Membership
- ii) Efficient Lock Mechanisms
- iii) Read/Write Locks

### Applications using Zookeeper

- i) ~~zoo~~ Yahoo Fetching (web crawling) Service
- ii) Message Broker
- iii) Kattq

& much more...

⇒ Zookeeper Implementation



Note: Zookeeper can deliver same messages twice, but I think they are idempotent.

i) Atomic Broadcast

Update Zookeeper State → Leader → Broadcast

ii) Request Processor

setDataTxn if success  
 errorTxn if fail

iii) Replicated DB

In-memory  
 Uses snapshots to recover faster.

iv) Client-Server Interactions

On write request → Notifications → Watchers

Read request → tagged with zxid → FIFO order  
 (may return stale value)

# Bit coin (2009)



P2P version of e-cash allowing online payments without having to go through a financial institution.

It have:

- Network time-stamp transactions
- Hashed into blocks
- based on proof-of-work

## ⇒ Introduction

Motive



→ Having computationally impractical mechanisms to reverse or commit fraud.

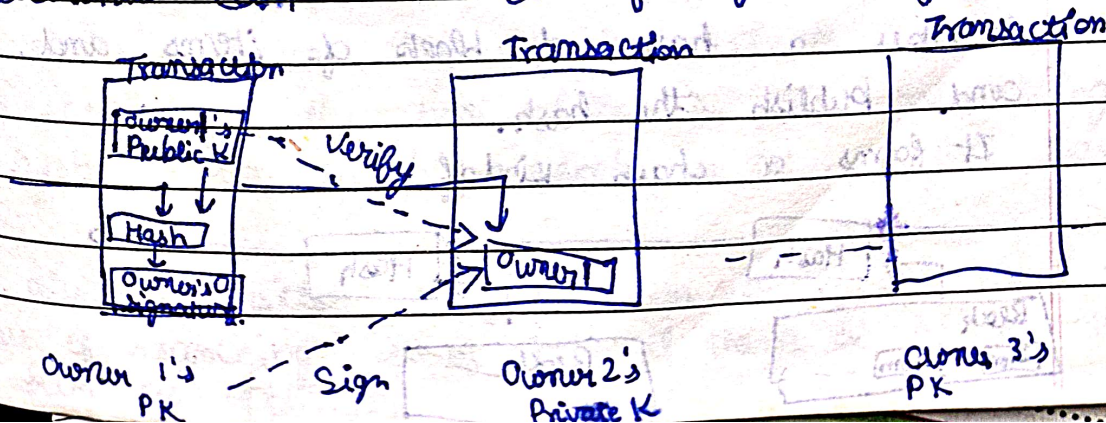
→ and allow routine escrow mechanisms.

both buyer & seller are satisfied

★ This paper solves double-spending problem.

## ⇒ Transactions

Electronic coin = Chain of digital signatures.



It is basically a digital signature stored on HDD.

→ Payer signs the <sup>hash of</sup> [public key + Previous transaction] <sub>↳ next of owner</sub>

→ Payee can ~~also~~ verify the transaction but don't know if payer double-spends the coin.

Approach 1: Use a Bank like system  
X Not decentralized.

Approach 2: Use bitcoin system.

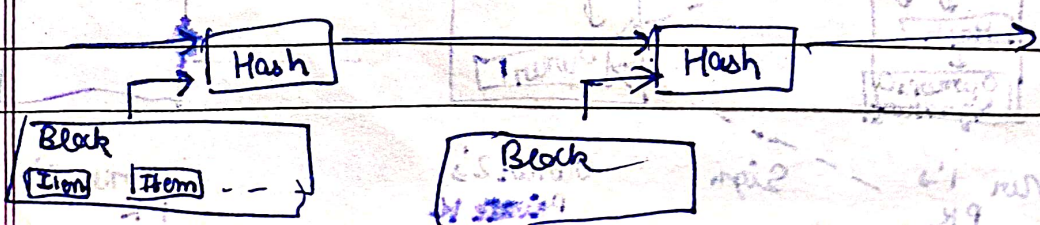
↓  
a system for participants to agree on a single history of the order which they were received.

⇒ Timestamp server

To avoid double spending, it introduces timestamp server:

It takes a hash of block of items and timestamp it and publish the hash.

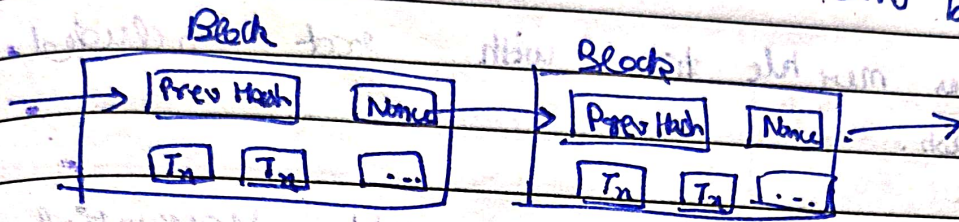
It forms a chain kindoff:



## ⇒ Proof of work

For a distributed timestamp server on a p2p basis, bitcoin uses proof-of-work.

add nonce to prev hash → starts with number of zero bits.



It avoids IP-based-one-vote system and uses computation-power-based government.

## ⇒ Network

- 1) New transactions are broadcasted to all nodes.
- 2) Each node collects new transactions into a block.
- 3) Each node works on finding a difficult proof-of-work for its block.
- 4) When a node finds proof-of-work, it broadcasts the block to all nodes.
- 5) Nodes accept the block only if all transactions in it are valid & not already spent.
- 6) Nodes express their acceptance of the block by working on creating the next block in the chain, using the hash of the accepted block as the previous hash.

Nodes always consider longest chain to be the correct one & in case of branching saves both of them & select longest one as soon as applicable.

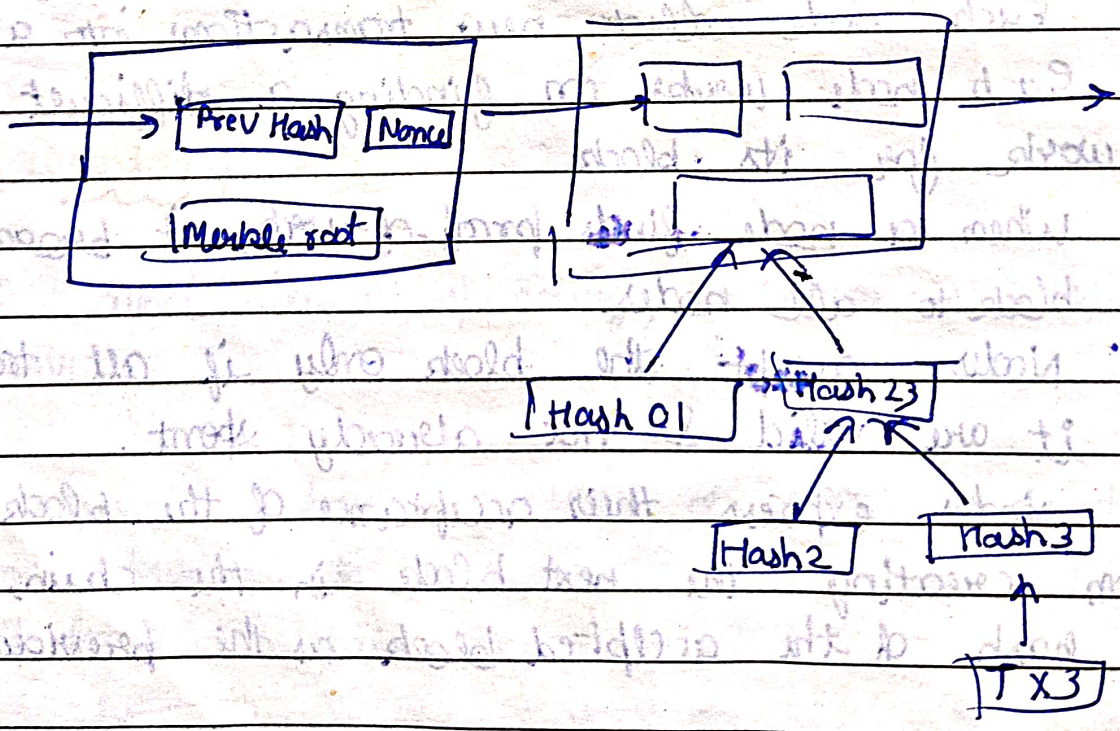
### ⇒ Incentive

Halves every 4 years.

+  
Transaction fees

### ⇒ Reclaiming disk space

- Uses merkle tree with root included in block's hash.
- Allows less space and avoid sequential storing of transactions.
- A user only needs to keep a copy of the block headers of the longest proof of work chain, by querying network nodes.

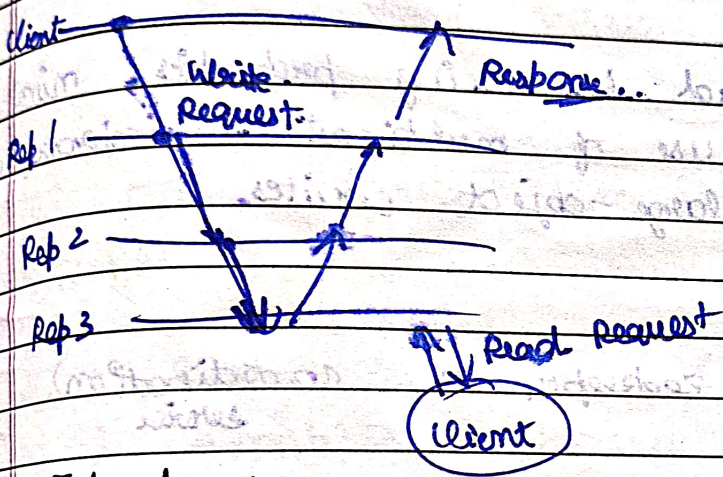


# Object Storage on CRAQ

Chain replication with Appointed Quies.

## Overview

A distributed storage system focussed on improving chain replication & read throughput while maintaining strong consistency.



- It creates not set as reads goes through one node
- It have stronger consistency & cheaper write operations

## Introduction

Object stores supports

↳ Read (Return data block under an object name)

↳ write (change the state of a single object)

These are like key-value DBs and are highly scalable. For this, object namespace is partitioned over many machines & each data object is replicated several times.

Chain replication looks promising but only read throughput pulls it down, that's where CRAQ comes in.

YAY? Basically, it divides read operations over all nodes in a chain:

1. Supporting read operation while preserving strong consistency. (load balancing).
2. Allows eventual consistency among read operations and specify maximum staleness value.
3. Geographical load balancing properties, mini-transactions & use of multicast to improve write performance for large-object writes.

⇒ CRAQ (Uses zookeeper as coordination service)

It works as follows:

1. Each node in the chain can store multiple versions of an object: one clean version and a dirty version per recent write. All versions are initially clean.
2. When node receives a new version of an object, it appends this latest version to its list for the object.
  - If not tail, mark as dirty
  - Otherwise, mark as clean (committed)

When a node receives a message for an object version, the node marks object version as clean, and delete all prior versions of the object.

on read request

- If version (latest) is clean, return value.
- otherwise, contact tail & ask for tail's last committed version number.

This model allows CRAQ to support three forms of consistency

- i) Strong Consistency (As above)
- ii) eventual Consistency (Don't ask tail for latest committed version no.)
- iii) GC with Maximum-Bounded Inconsistency (As the name suggests)

Problems with CRAQ:

- i) If one node is slow, affects all of the nodes in a chain
- ii) We can also create multiple chains instead of CRAQ to have even more simplicity.

# Aurora (2017)

↓  
Design Considerations for High throughput cloud-native relational DBs.

## ⇒ Abstract

A relational DB service for OLTP offered as part of AWS (Transaction Prof)

- Since, today storage & compute is better today than network, aurora addresses this constraint.
- It also uses an efficient asynchronous scheme to achieve consensus without chatty & expensive recovery protocols.

## ⇒ Introduction

Well known problem nowadays is network between DB tier and storage tier.

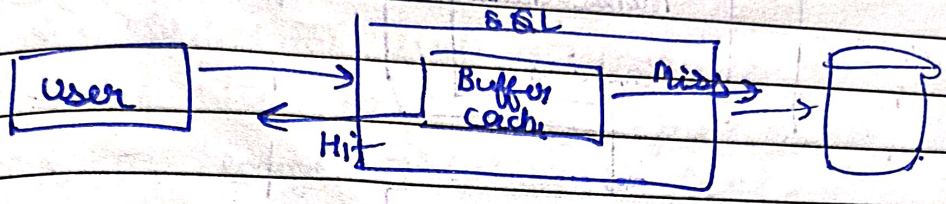
→ Before starting it, let's talk about internals of DB: (SQL)

• SQL server writes modified data page to the disk asynchronously:

- Lazy writing → Inspect & remove infrequently used pages.
- Eager writing → Fixing non-logged operations
- checkpoint → Periodic check to scan the buffer cache & write dirty page into the disk.

• On read request / SQL query. (But SQL server store data in 8 KB pages).

- i) It locates the data & retrieves the pages that contain these information.
- ii) by retrieving it first from the disk to buffer cache
- iii) Finally, returns result to the reader.



Read

• on write request, update query to be precise:

- i) locate page in buffer cache, if not found, retrieve into buffer cache from disk.
- ii) Acquire locks on the page to be updated, update them now.

iii) Now, pages are called dirty page.

iv) This modification is recorded in transaction log and an ack is sent to the user. (changed page is still in buffer cache).

Write-check logging

v) later, by one of the processes, changes are flushed into disk.

Coming back next, this page miss or background tasks can cause context switches & resource contention

→ Amazon Aurora, addresses these issues by leveraging the redo-log across a highly distributed cloud environment

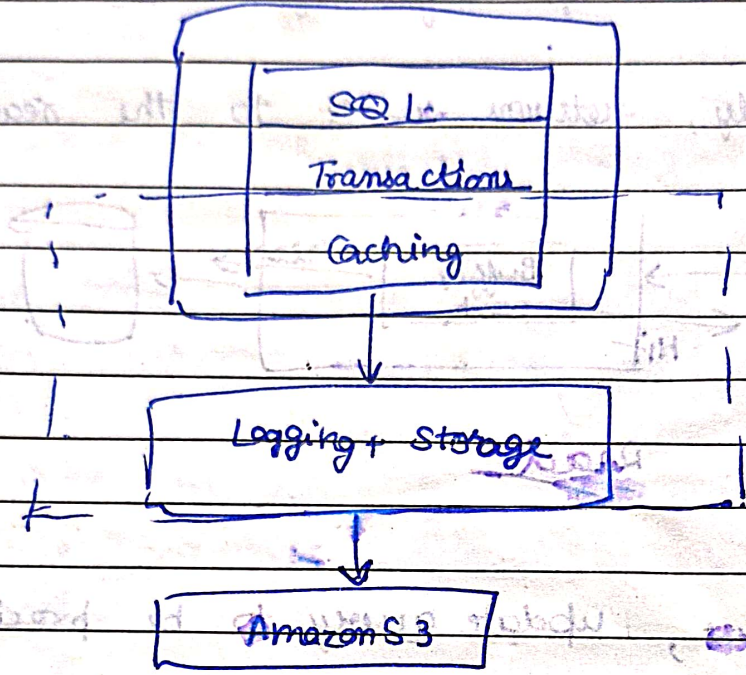
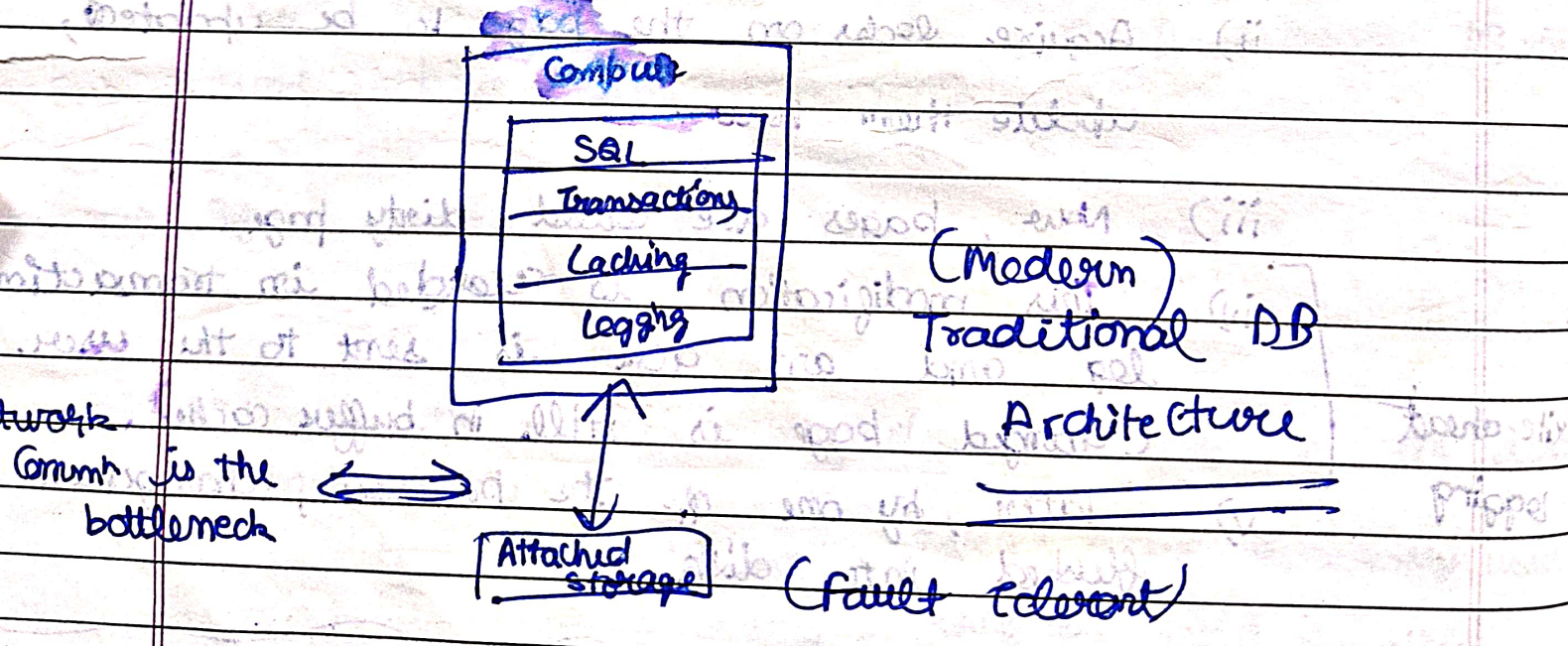


Fig: Move logging & storage of the database engine



→ The Aurora architecture offers 3 major advantages

1. Protection of the DB from performance variance & transient or permanent failures.
2. Reduce network IOPS by only writing redo log records to storage.
3. Expensive operations such as (backup & redo recovery) are offloaded to continuous asynchronous operations across a large distributed fleet.

Sub 16 Durability

Option 1: Replication

3 ways with 1 copy per AZ  
use read/write quorums of 2/3

Can it survive?

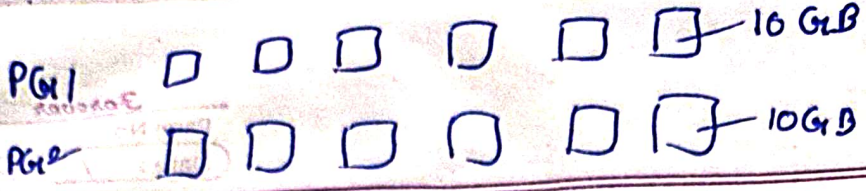
- i) AZ Failure ✓ (2/3 still live)
- ii) AZ+1 Failure X (lose 2/3)

Option 2: Replicate 6 ways with 2 copy per AZ

write quorum of 4/6, read quorum 3/6

Can it survive

- i) AZ failure (4/6 still live)
- ii) AZ+1 failure (still have 3 copies)



with  
Now, the 6 copies, how to get the largest DB.

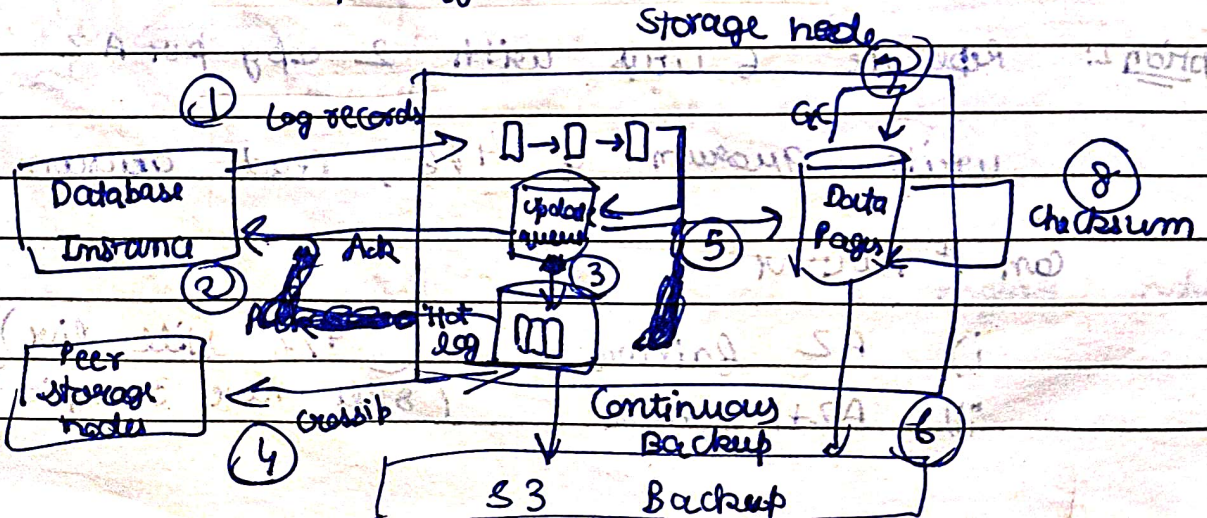
i) Partition volume into  $n$  fixed-size segments and replicate each segment 6 ways into a protection group (PG)

ii) If segments are too small, failures are more likely otherwise, repairs take too long.

ii) so, they choose 10 GB as segment size. since on a 10 Gbps network line, it takes only 10 seconds to repair. (so, a total of 1 minute)

Sub 2: offloading the lower quarters of traditional DB to this tier. (LOG IS THE DB)

In traditional DB, replication leads to write amplification, so  $\times$  overhead, with it's different architecture, offloads work by:



Normal operations

For decision, it uses LSN (latest) to

Log Sequential Numbers



Writes

VDL = Volume durable LSN (Trusted data for each replica)

VCL = Volume Complete LSN (Guarantee Availability of all prior records)

$LSN \leq VDL + LSN \text{ Allocation Limit}$   
(Throttling purposes)

i) Database determine the PCs

ii) Sends LSN to all replicas in PC

iii) On quorum of 4 ack, send response to client

Reads (No need for quorum though)

i) Same like traditional DB

ii) Establish a read-point & determine segment to get data from

OR

Basic read-quorum & get the latest version data

Frangipani (2000)



A scalable distributed FS

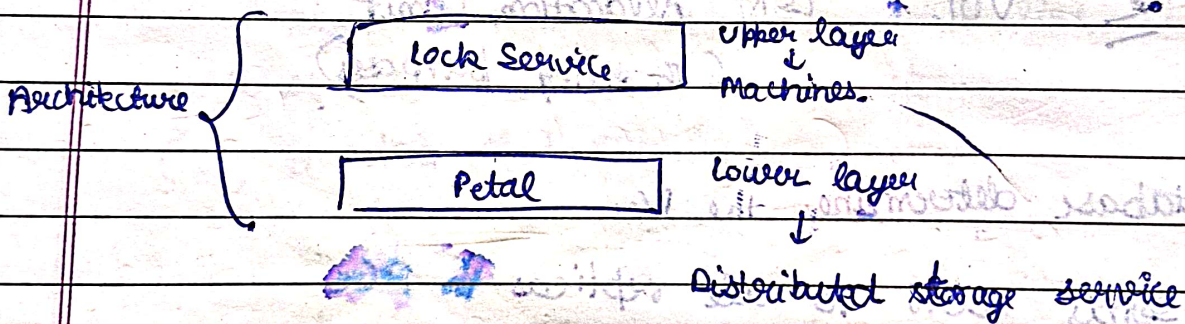
⇒ Abstract

An ideal Distributed FS would provide

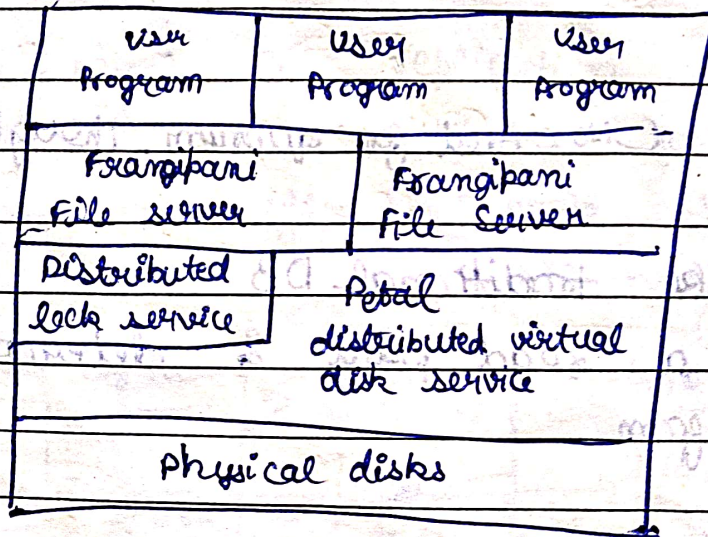
- i) Shared access to the same set of files
- ii) Scalable
- iii) Availability is high
- iv) Minimum Human administration.

It has a fairly simple set of features.

Frangipani approximates this ideal, by a two layer architecture:

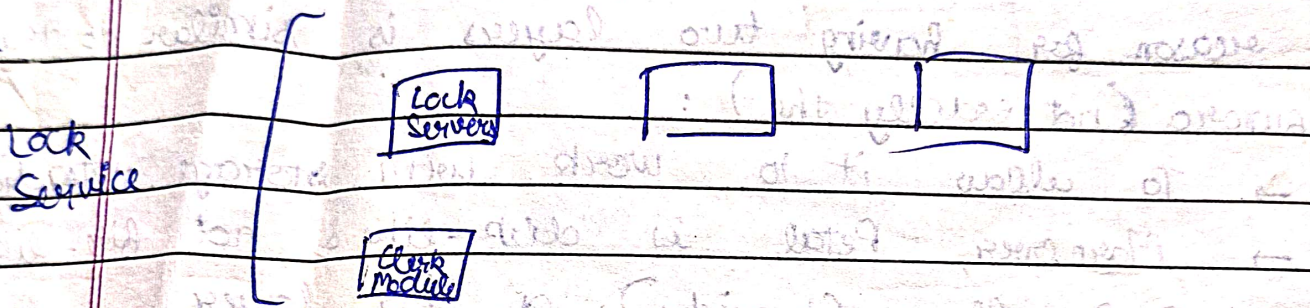


⇒ Introduction



## ⇒ Lock Service

- Multiple readers / single-writer
- All locks are sticky, i.e., a client will retain a lock until some other client needs a conflicting one.
- Based on lease, 30s on start, renew later on.



## frangipani

- Uses fault-tolerant, distributed failure detection mechanism to detect the crash of lock servers.
- i) When a frangipani server crashes, the locks that it owns can't be released until recovery have been performed.
- ii) So, lock service asks clerk on another frangipani machine to perform recovery & then release all locks belonging to crashed server.

- FGP file servers can be scaled up & uses locks for co-ordinating their actions.
- Petal & the lock service are also distributed for scalability, fault tolerance. They provide large virtual disks (replicated).
- FGP file server module uses Petal device driver for every comm<sup>n</sup> required.

- The reason for having two layers is similar to that of Aurora (not exactly tho):
  - To allow it to work with storage abstractions
  - Moreover, Petal is disk-like & not file-like, FGP allows (provides) a FS layer.

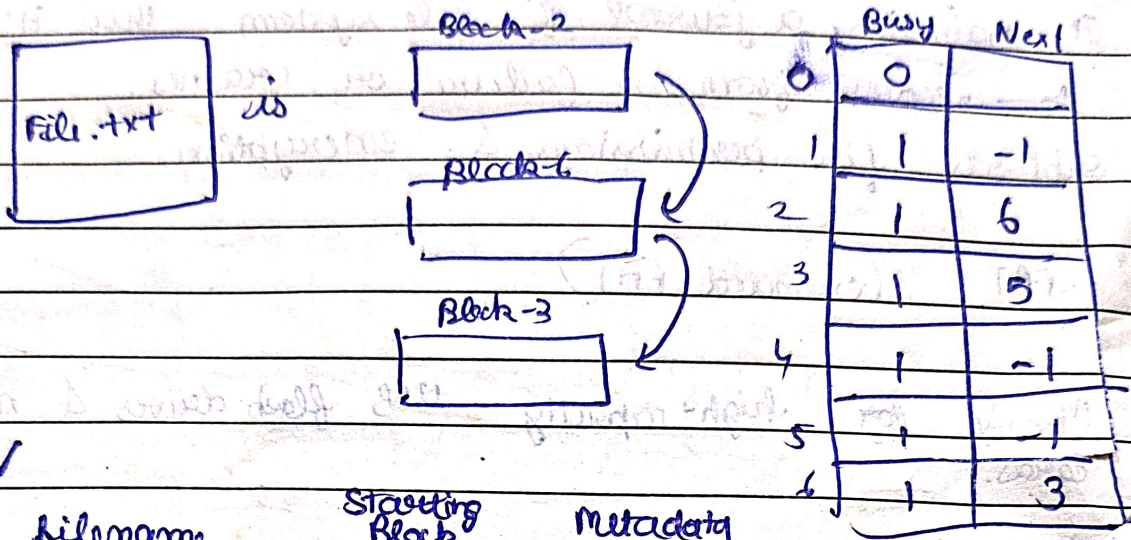
### ⇒ Disk Structure

FGP uses the large, sparse disk address space of Petal.

- A Petal virtual disk has  $2^{64}$  bytes of address space.
- If file is written, petal commits physical disk space to virtual address.
- Petal also provides a "decommit" primitive that frees the physical space.
- Each commit size is 64 KB of chunk

→ Before starting it, let's talk about FS:

i) FAT (File Allocation Table)



filename	Starting Block	Metadata
file.txt	2	---

(Directory Table format)

Note: Root file directory has a fixed address, so we know how to get started.

FAT 12, 16, 32 are (represent) the bits in each table element of FAT.

So,	Max file size	Max Volume size.
FAT12		
FAT16		
FAT32		

ii) NTFS (New technology file system)

- File size limit = 16 EB (exabytes)
- It maintains a journal of file system, thus it's able to recover from a failure or crashes.
- supports file permissions & encryption.

iii) exFAT (extended FAT)

- Mainly for high-capacity USB flash drives & memory cards.

iv) ext4

- Linux default FS
- No native windows or macOS support.

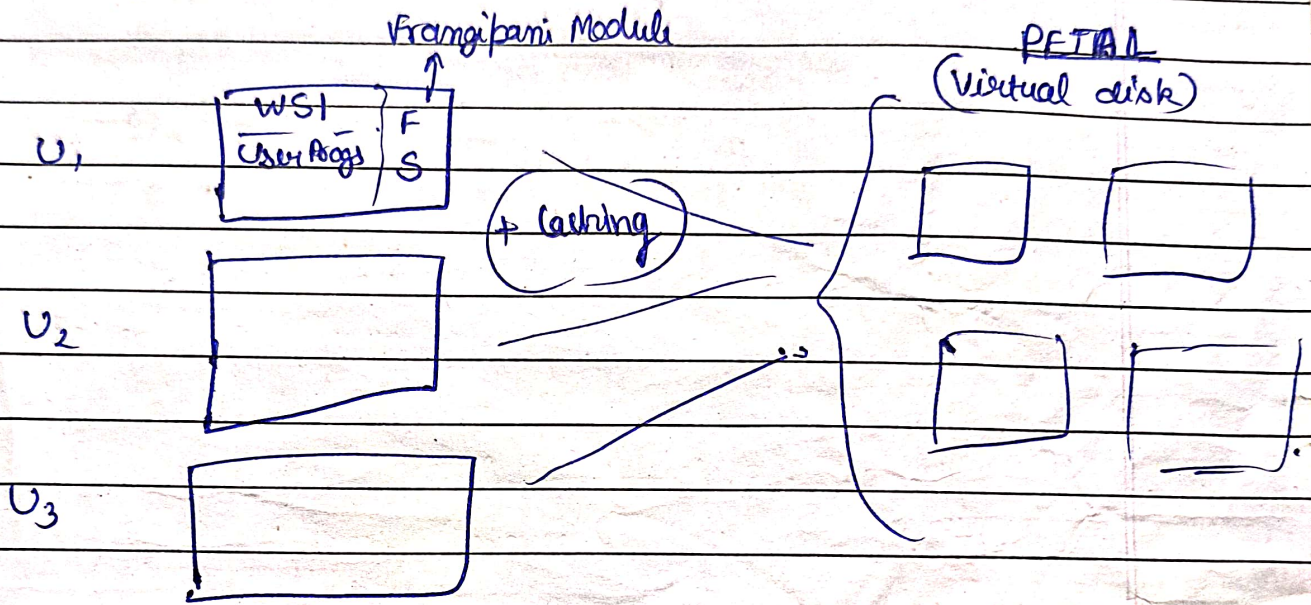
on the other hand, ~~windows~~ (linux) ubuntu can read-write NTFS

v) HFS, HFS+ & APFS ← Apple file system

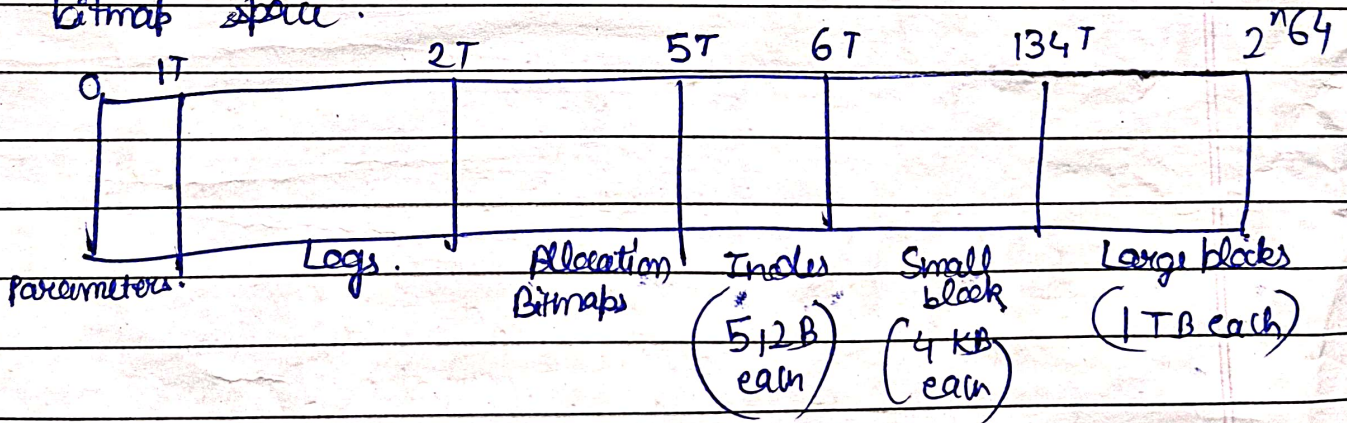
- Hierarchical file system for macOS
- No native windows or linux support.

Summarizing

- i) Cache-coherence
- ii) Public Client logs
- iii) Framigipani version numbers



\* Each server has its own log & its own blocks of allocation bitmap space.



So, to have caching & a decentralized platform, we need to work upon

- a) Cache coherence
- b) Atomicity
- c) Crash Recovery.