

Designing Data Intensive Applications

CHAPTER 1

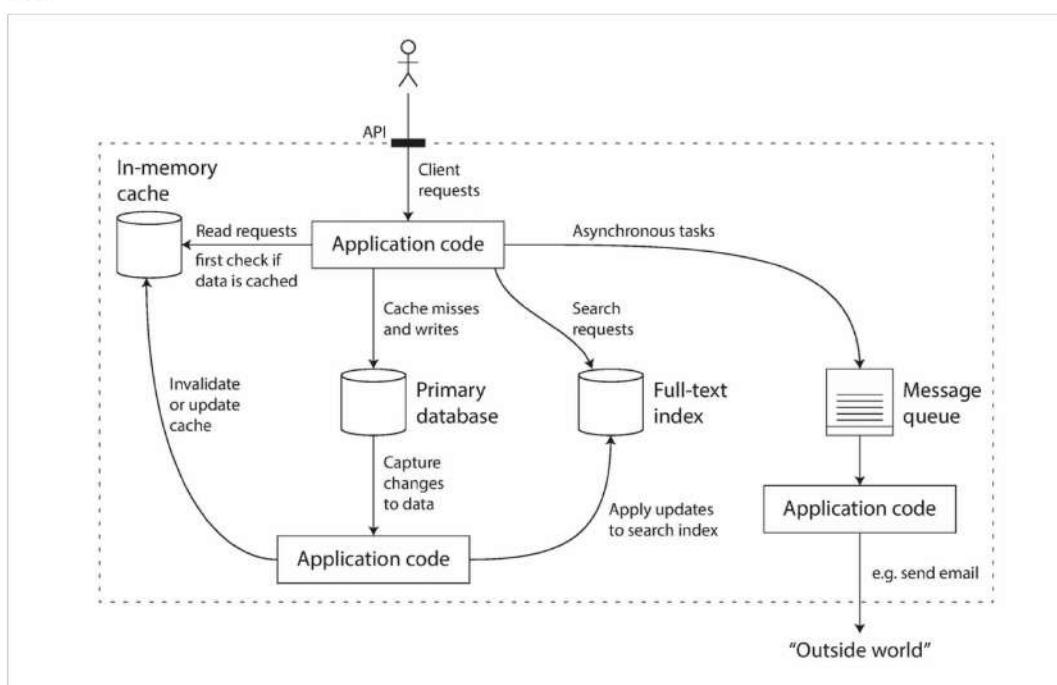
Reliable, Scalable and Maintainable Applications

Reliability → Tolerating Hardware and Software failures
Human Error.

Scalability → Measuring load and performance
Latency percentiles, throughput

Maintainability → Operability, Simplicity and evolvability

One Possible Architecture



How will you ensure below when designing a service →

- ① Ensure data is correct even if things go wrong?
- ② Provide good performance to clients even when parts of system are degraded?
- ③ Handle sudden increase/burst in TPS?
- ④ How will you design an API for system?

Points to consider →

- ① There is no single solution to solve a problem.
It could be possible that ElasticSearch is best search indexing solution available but you went ahead with Solr due to legacy dependencies.

Design can depend on →

- ① Skills and experience of people.
- ② legacy system dependencies.
- ③ time scale for delivery.
- ④ organization's tolerance for different kind of risks.
- ⑤ Regulatory constraints.

Three Important pillars of Software System →

[A] Reliability

System should continue to work correctly even with hardware / software faults or Human errors.

with expected traffic

What can go wrong in system?

These are called faults.

→ fault is different from failures.

↳ watch Error vs fault vs failure video on my youtube channel.

↳ fault means server was unable to serve your request. (probably a specific component of service)

↳ failure means whole system stopped providing required service to user.

↳ To make systems fault Tolerant, you can design experiments to test system limits.

↳ Watch "Chaos Engineering" video on channel.

Hardware Faults →

any issues in system Hardware - CPU, disk, RAM, router
...etc.

Example → On storage cluster of 10000 disks, avg 1 disk die/day.

To solve this, you should have backups in place

↑
if one component dies, redundant component
can take its place.

Software Errors →

Any issue in software which is running on hardware.

Example -

- ① Bug in software
- ② You can try to make your system fault tolerant, but if there are issues in third party software or application we, the system can get into failure state.
- ③ The upstream service is unresponsive.
- ④ Cascading failures - fault in one component triggers fault in another component

↑ Watch "cascading failures" video on YT channel.

Suggestive Resolutions

- ① Try to close on all unknowns / edge cases during system design.
- ② Testing
- ③ Process Isolation
 - ↳ example tabs in Chrome browser.
- ④ Allow process to crash and restart
- ⑤ Monitoring + chaos engineering + Alerting + Logging

Human Error → Humans are unreliable 😊

- ① Experiments in production should be well monitored.
- ② Right amount of abstraction.
- ③ Thorough Testing.
- ④ Monitoring + Alerting. + Best training practices.

So what do you think?

Is Reliability Important?

Scalability → what happens if your system's load is increased by 10%, will it be reliable?

→ It is used to describe system's capability to handle increased load on system.

↑
What is load? → TPS for web server, read write ratio on database, hit rate on cache..

Example Twitter → How a tweet published by user reach his/her followers?

↳ followers asks for new tweet?

↳ user pushes tweet to all its followers?

→ Watch similar analogy in

"Instagram NewsFeed" video on YouTube.

Here in case of twitter, distribution of followers per user is key load parameter to think around scalability.

Performance → Is system working as per our expectations?

→ What is TPS that can be handled by server for given CPU and memory?

→ What is latency to serve API response?

Response Time determination → p50, p90, p99, etc.

lets say response time = 1 second, then p90 means 90 out of 100 requests take <1sec and other 10 requests take 1sec or more.

How to handle load? Load is increased but performance is not impacted.

- Vertical Scaling
- Horizontal Scaling

Automatic scaling without human interventions → Auto Scaling

→ There is no single solution that can be developed for all large scale systems, it varies use-case by usecase.

Maintainability

Maintain existing software - fixing bugs etc {On call}

① Operability → making life easy for operations.

- ↳ Monitoring health of systems
- ↳ Debugging system issues
- ↳ Software updates
- ↳ System Automation
- ↳ Oncall Rubrocks and System dashboards

② Simplicity → managing complexity.

- ↳ How different components are coupled?
- ↳ hacks introduced in system as short term solve.
- ↳ Complex systems makes testing a cumbersome task.
 - ↳ possibility of introducing a bug.

How can you achieve it?

→ Abstraction

example → Java abstracts out memory cleanup logic.

③ Evolvability → making changes easy

System requirements change and it will requires modification in existing system , how easy is this process ?

Subscribe YouTube channel for more such content.

Channel - MS Deep Singh

Happy learning 😊

Designing Data Intensive Applications

CHAPTER 2

Data models and Query Languages

Relational Model - Data organized in relations ~ tables in SQL.

why NOSQL?

- ① More scalability as compared to SQL databases
- ② open source instead of paid software
- ③ Query operations that are not supported by SQL.

The way we write code → most applications use OOP paradigm to write code and you'll require extra translation layer between objects in code and data models in SQL tables.

users table				
user_id	first_name	last_name	summary	
251	Bill	Gates	Co-chair of ... blogger.	
	region_id	industry_id	photo_id	
	us:91	131	57817532	
regions table				
id	region_name			
us:7	Greater Boston Area			
us:91	Greater Seattle Area			
industries table				
id	industry_name			
43	Financial Services			
48	Construction			
131	Philanthropy			
positions table				
id	user_id	job_title	organization	
458	251	Co-chair	Bill & Melinda Gates F...	
457	251	Co-founder, Chairman	Microsoft	
education table				
id	user_id	school_name	start	end
807	251	Harvard University	1973	1975
806	251	Lakeside School, Seattle	NULL	NULL
contact_info table				
id	user_id	type	url	
155	251	blog	http://thegatesnotes.com	
156	251	twitter	http://twitter.com/BillGates	

Image - Wikimedia Commons

```
{ user_id : 251,  
  first_name : Bill  
  last_name : Gates  
  positions : [  
    { jobtitle : Co-Chair }  
    { jobtitle : Co-founder }  
  ]  
  ...  
}
```

Fetch profile for Bill Gates

- ① Relational → perform multiple queries in tables or perform complex join
- ② Document → all information at single place.

Comparison of document Model vs Relational Model

① Application code simplicity →

- ↳ It will depend on application, what is the usecase it is solving.
- ↳

Property	Document	Relational
① Data in application ?	→ requirement is document like structure. → poor joins support.	→ requirement is many to many relationships.
② Schema flexibility	→ Don't enforce strict schema. → schema on read - schema is maintained in application code. → Adding new fields is easy task.	→ strict schema → schema on write. → Adding new columns require proper migration strategy ↳ could be a real pain point if data size is huge. → If all records are expected to have same structure, RDBMS is good way to enforce it.

③ Data Locality → document is stored as continuous string. for queries → If app require entire doc. to be read, it is useful.	→ If data is split into multiple tables and entire data needs to be retrieved, there are performance issues. → Oracle has feature "multi-table index cluster tables" which provide locality property in RDBMS.
---	---

Similarity

- ↳ There can always be feature similarity between SQL & NOSQL.
It totally depends on usecase we're trying to solve.
- ↳ NOSQL database has support for joins.
- ↳ SQL databases can store data in document form.

Query Languages for Data

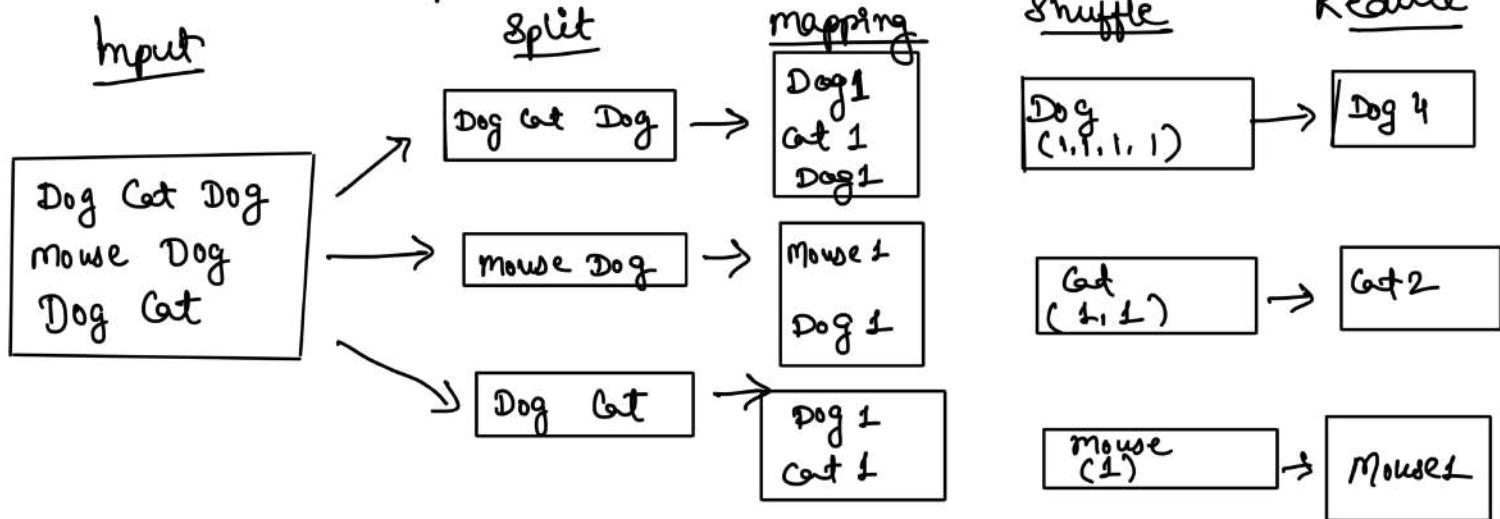
SQL - declarative query language

Declarative Language	Imperative Language
<p>→ You specify pattern of data you need.</p> <p>↳ but you don't specify how these results are achieved.</p> <p>↳ Database system's query optimizer to decide how to gather this data.</p>	<p>→ Tells computer to perform certain operations in certain order.</p>

Map Reduce Querying

→ programming model for processing large amounts of data.

↳ in later chapters → CH 10



Graph Data Model

↳ Document is good if there are mostly 1-many relationships or no relationships between records.

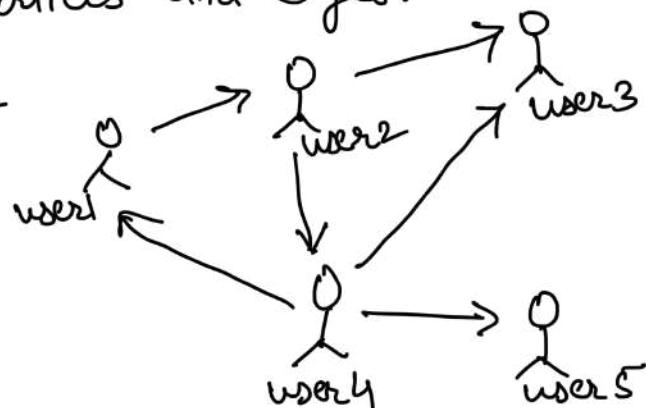
↳ many-many relationships can be very well handled by RDBMS.

But what if these relationships become very complex and nested.

Then graph data model is more suitable.

Graph → vertices and edges.

Example



↳ Graph provides flexibility in data modelling as compared to RDBMS.

Example Neo4j

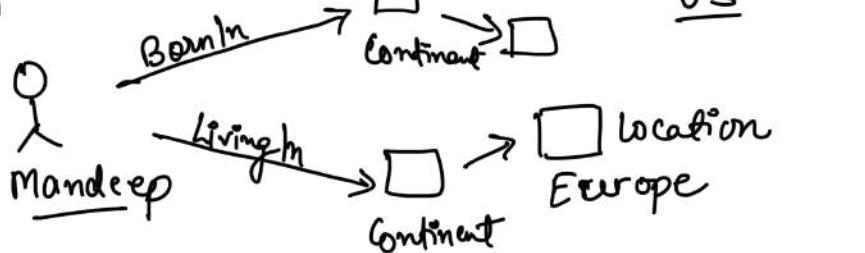
query language - cypher (declarative)

example - find all people who immigrated from US to Europe

↳ BornIn = US

&&

LivingIn = Europe



Graph Databases Compared to Network Model

Network	Graph
<ul style="list-style-type: none">① schema specifies record type could be nested within each record type.② To reach particular record, traverse the access path.③ Children of record are ordered set.④ Imperative queries	<ul style="list-style-type: none">① No such restriction. A vertex can have edge to any other vertex.② You can directly go to any vertex by its unique identifier.③ vertices & edges are not ordered.④ Support for both declarative & imperative queries.

Stay Tuned for follow up chapters.

Subscribe for more. YouTube - Ms Deep Singh

Happy Learning ☺

Designing Data Intensive Applications

CHAPTER - 3

Storage And Retrieval

How will data be stored in database?

How data will be retrieved from database?

Date Structures that Power Your Database

Why do I care what's happening in backend? \Rightarrow It helps in selecting a data storage for your system. There are lot of options these days 😊

Log File \rightarrow Assume you're keeping a key value pair and store in text file in append mode.

For any new write query, you'll just append an entry.

How will read value for given key?

\hookrightarrow Scan entire file for occurrence of key? \rightarrow that will be $O(n)$ operation

\hookrightarrow It's not efficient if data grows to large scale

\hookrightarrow You'll need indexes (watch youtube video on Indexes on channel)
 \uparrow
create indexes for data written.

\hookrightarrow Should I create indexes for all the data? It will make read faster

→ NO, it will make write slower. Along with append operation, you also need to write data to indexes.

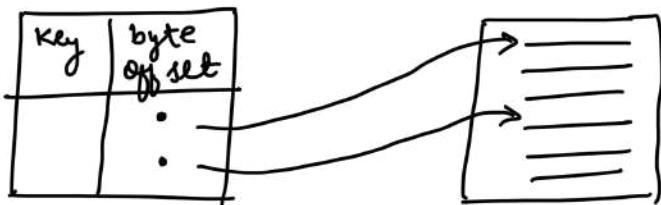
Hash Indexes → How will you index data on disk?

Assumption - Data storage provides only append operation on file

Indexing Strategy →

↳ Create in-memory hashMap.

↳ Every key is mapped to byte offset in file.



↳ gft is used by Bitcask (Storage engine in Riak)

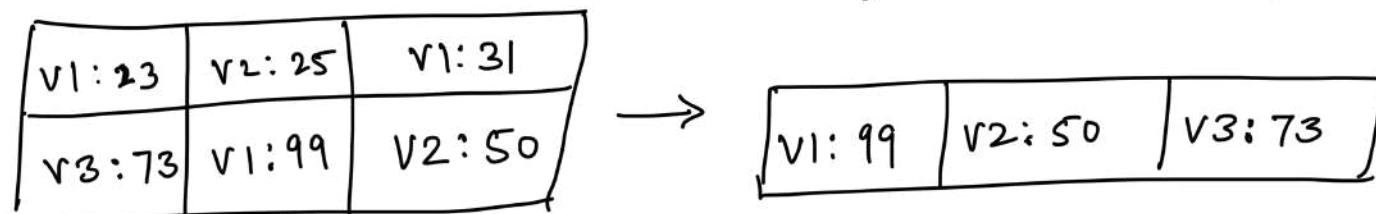
↳ In-memory HashMap helps if there are limited no of keys.

↳ Example - maintain a counter how many times video is played.

↳ How will you ensure you're not running out of ^{since you're only appending} disk space

↳ break the log into segments of certain size.

↳ ② perform compaction - remove duplicate keys in log and only keep recent update.



↳ Each segment has its own in-memory hash Table.

↳ find value for key : check most recent segment's hash Map.

↳ if not found, check 2nd most recent.

Implementation Details

- ① file format - Binary → encodes the length of string in bytes, followed by raw string.
- ② Record Deletion - append special deletion record to data file.
↳ called Tombstone {watch Discard msg storage video}
- ③ Crash Recovery - Maps in memory will be lost in events of crash.
↳ keep snapshot of segment's hash map on disk.
- ④ Partially written records - use checksums, helps in identifying corrupted parts and delete them.
- ⑤ Concurrency - Data files are append only → immutable.
↳ can be read concurrently.

Why not update the key instead of appending?

- ① Much faster than random writes, especially on magnetic spinning disks.
- ② Concurrency and crash recovery is simple.

Limitations

- ① Hash Table must fit in memory.
- ② Range queries are not efficient.

SSTables and LSM-Trees

optimizing the limitations mentioned above ↗

↳ change the format of how segment files are stored

↳ make sequence of key-value pair "sorted by key"
↑

This is called Sorted String Table (SSTable)

↳ Each key should appear only once in merged segment.
(handled in compaction process)

A) Merging of segments → algorithm similar to merge sort.

① Read input file side by side.

② look for first key in each file.

③ Copy lowest key to output file. Repeat again.

④ If same key in multiple segments → pick value from most recent segment.

S1	V1: 5	V2: 3	V3: 3	V4: 7
----	-------	-------	-------	-------

S2	V5: 9	V6: 11	V2: 5	V3: 9
----	-------	--------	-------	-------

S3	V1: 11	V2: 9	V7: 23	V3: 21
----	--------	-------	--------	--------

↓

V1: 11	V2: 9	V3: 21	V4: 7	V5: 9	V6: 11	V7: 23
--------	-------	--------	-------	-------	--------	--------

Search for a key

- ↳ ① Inmemory index is not required for maintaining offset for key.
- ② Maintain offset for some of keys (sparse index) and they traverse between these offsets to find particular key.
- ③ Group the records for keys you're maintaining in Inmemory index.
↳ saves disk space and reduce I/O.

Construction & maintenance of SSTables

for any write query -

- ① Add it to in-memory balanced tree structure. — called memtable.
- ② Write memtable to disk after certain threshold (say 1MB) — called SSTable
- ③ For read query, Search in memtable, then most recent segment of SST.
- ④ Run merge & compaction time to time.

Performance of Database built on LSM Tree

- ① Reads can be slow as you traverse from one segment to another.
 ↳ Bloom filters are used → data structure helpful in approximating contents of set.
 Example → will tell if key is present in DB or not saving traversing DB.
- ② Compaction & merging → how to determine order & timing.

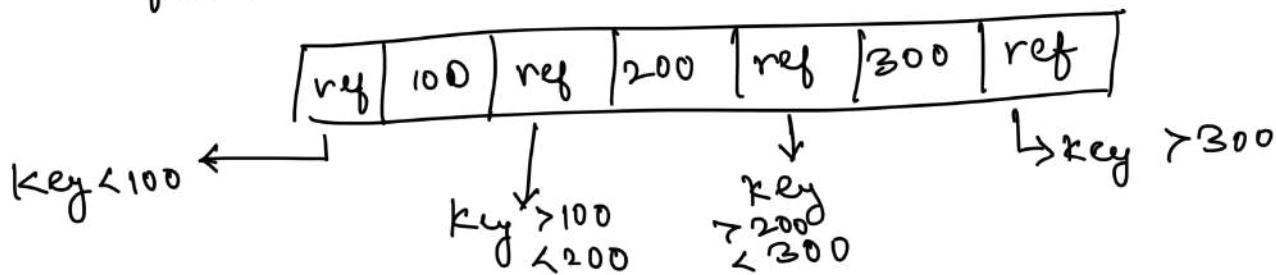
Size-Tiered	Leveled Compaction
HBase. Cassandra uses both.	LevelDB, RocksDB
smaller SSTables are successively merged into larger SSTables.	key range is split up into small SSTable and older data moved to separate levels.

B Tree → It is not "Binary" Tree.

↳ key value pairs sorted by key. → on disk

↳ BTree break database into fixed size blocks.

↳ Each Block can be identified by address, using which one page refer to another.



Handling Crashes in Btree → what if hardware crashed at time of any update?

It maintains additional data structure → WAL (Write ahead Log)

↳ append only file - every modification to Btree is written here before actual update.

Write amplification → one write to DB resulting in multiple writes on disk over course of DB's lifetime.

L S M Tree

- ① Typically faster for writes
- ② It can sustain higher write throughput due to less write amplification.
- ③ Can be compressed better
- ④ Compaction process can affect ongoing read/write.
 - ↳ more observable at high throughput.

BTree

- ① Typically faster for reads.
- ② More write overhead, update existing page.
- ③ Leaves some space unused due to fragmentation.
- ④ Each key is present a single place.
 - ↳ offer strong transactional support

Full Text Search and fuzzy indexes

↳ Lucene example Elasticsearch

ex. to search for mis-spelled words or synonyms.

In-memory Databases

memcached → intended for caching purpose, it's ok to lose data if machine is restarted.

↳ Inmemory databases for durability → dump periodic snapshots to disk.

↳ Redis and couchbase provide weak durability by writing to disk asynchronously.

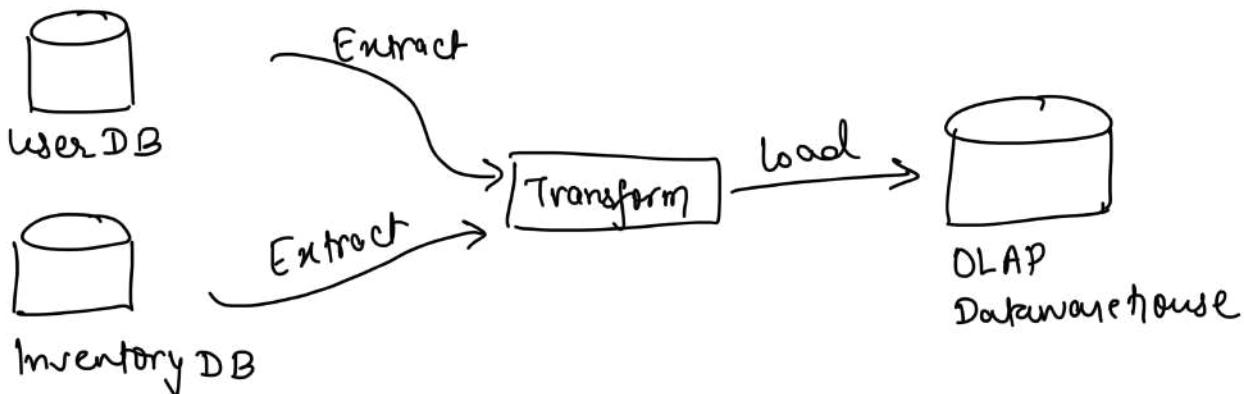
Your use-case → Transaction Processing or Analytics ?

OLTP

- ① Read - small number of records per query fetched by key
- ② Write - Random access, low latency writes
- ③ GB to TB

OLAP

- ① Aggregate over large set of data.
- ② ETL
- ③ TB to PB



OLTP databases

Column Oriented Data Storage

- ① The table can be ≥ 100 columns wide but we might be accessing only 4 to 5 columns at a time.
- ② To get the data, you need to load all rows matching criteria to these columns.
- ③ Column Oriented \rightarrow Store values from each column together instead of rows.
- ④ Column compression is easier \rightarrow there are chances that values in column are same for different rows.
- ⑤ Example - Cassandra, HBase
 - \hookrightarrow not strongly column oriented.
 - \rightarrow they store all columns for row together along with rowkey.
 - \rightarrow don't use column compression.

Subscribe for more such content.

YouTube channel- MsDeep Singh

Happy Learning 😊

Designing Data Intensive Applications

CHAPTER 4

Encoding And Evolution

As application grows over time, new changes will be introduced in application.

- ↳ How will you rollout these changes without breaking the application and continue serving customers?
 - ↳ Rolling upgrade or staged rollout.
 - ↳ don't deploy at once. Deploy on one node then validate, Deploy on other nodes and validate and so on.
- ① Backward compatible → New code can work with older data.
- ② Forward compatible → Old code can work with newer data.

How programs work with data?

- ① In memory → data to be accessed by CPU on machine via objects lists, hashmaps, etc.
- ② Send data over network → encode it to some different format example JSON.

Translation from ① → ② is called encoding/serialization/marshalling
reverse is called Decoding.

How to encode?

- ① Language specific support
 Java - `java.io.Serializable`
 Python - `Pickle`

They should not be avoided for any extensive purposes

- ↳ Encoded in Java will cause problems to decode in Python.
- ↳ Lack of forward/backward compatibility.

② JSON, XML and Binary Variants

- ↳ Support across programming languages.

JSON/XML/CSV

- ↳ Can't distinguish between number and string.
- ↳ No support for binary strings. ↳ Workaround using Base64.
- ↳ CSV don't have any schema, application define meaning of each row/column.

Binary Encoding

- ↳ very helpful as data size grows.
- ↳ Example for JSON - BSON, BJSON, etc.

Thrift and Protocol Buffers

- ↳ Both require schema for data to be encoded.

Example Thrift

```
struct Person {
    1: required string username,
    2: optional i64 favouriteNo
}
```

↳ `optional` & runtime check and don't reflect in encoding.

→ Thrift has two binary encoding formats

↳ Binary Protocol

↳ Compact Protocol. → takes less space.

Schema Evolution

- ↳ Schema evolves over time. How will you ensure encoding/decoding & backward/forward compatible.

- ↳ Every field in encoded data is identified by tag number.
- ↳ If field is not set, it is skipped from record.
- ↳ Field names can be changed, but field tags should not be changed as they are used for encoding data.
- ↳ Each new field added should be given a new tag number.

Backward compatible

- ↳ ignore the field(datatype) helps in skipping byte offset, which are newly added on read.

Forward compatible

- ↳ As we ensure unique tag number for each field, it will not cause any issues.
- ↳ It should not be made required to avoid runtime errors.

Data Type evolution in Schema

- ↳ value can loose precision or get truncated.

Avro → binary encoding format

↳ use schema to specify structure for data to be encoded.
↳ no tag numbers.
record Person {
 string username;
 array<string> interests;
}

- ↳ most compact as compared to others.

- ↳ Avro defines schema in terms of writer's schema (who writes) and reader's schema (who reads)

- ↳ It works fine as far as they are compatible.

↳ to add/remove value, field should have default value.

↳ maintain version number as schema evolves.

Dynamically generated schema in Avro

↳ You don't have to worry about tag numbers unlike Thrift.
↳ the schema can be generated from your object schema.

Dataflow

① Dataflow through Databases

↳ database will change over time e.g. add new columns to RDBMS

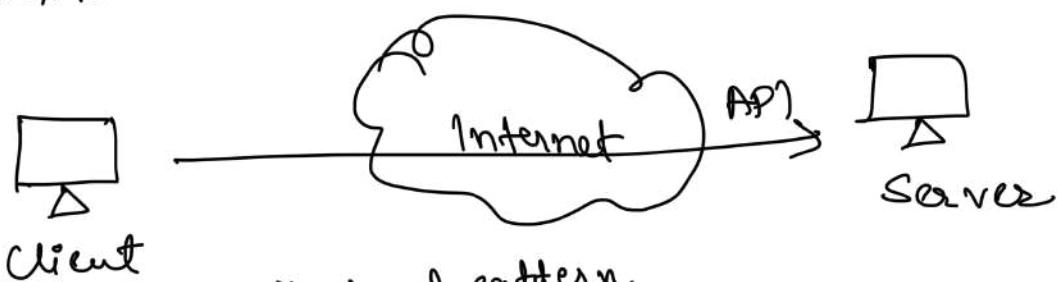
↳ You can specify a default value for newly introduced columns.

↳ Handle specifically in application code as per requirement.

→ It is good choice as compared to migrating entire db to new schema.

② Data Flow through Services - REST & RPC

↳ communication over network.



architectural pattern

REST → Design philosophy that builds upon principles of HTTP

SOAP → XML based protocol for network API calls.

e.g. GET, PUT

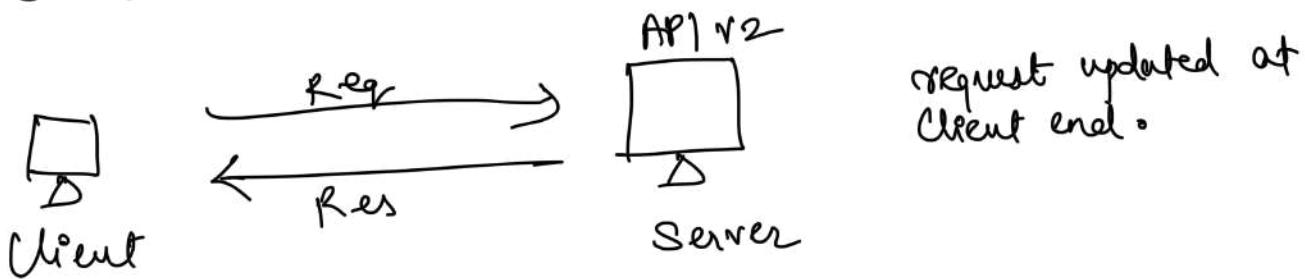
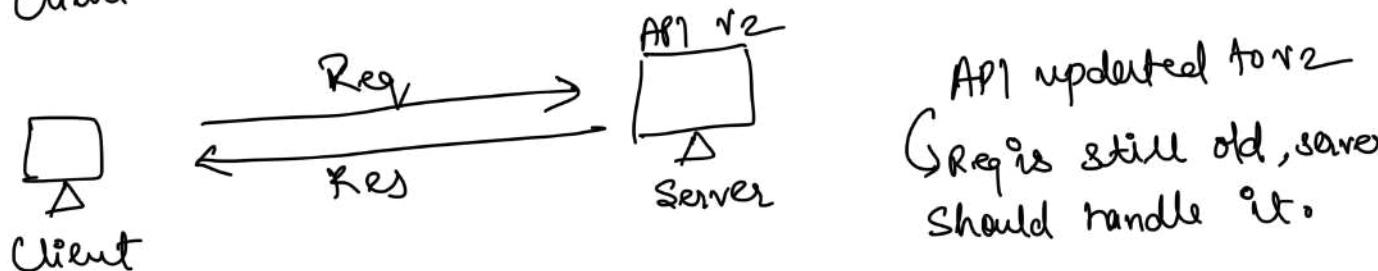
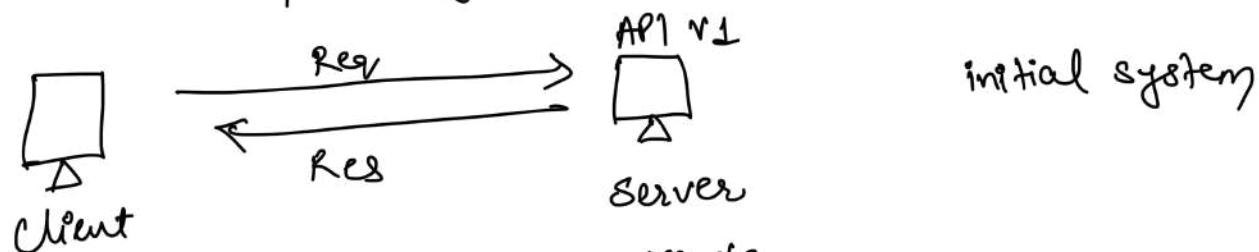
RPC → Remote Procedure Call

↳ make request to remote network service.

One example - gRPC → support streams, a call can consist of series of request and response over time.

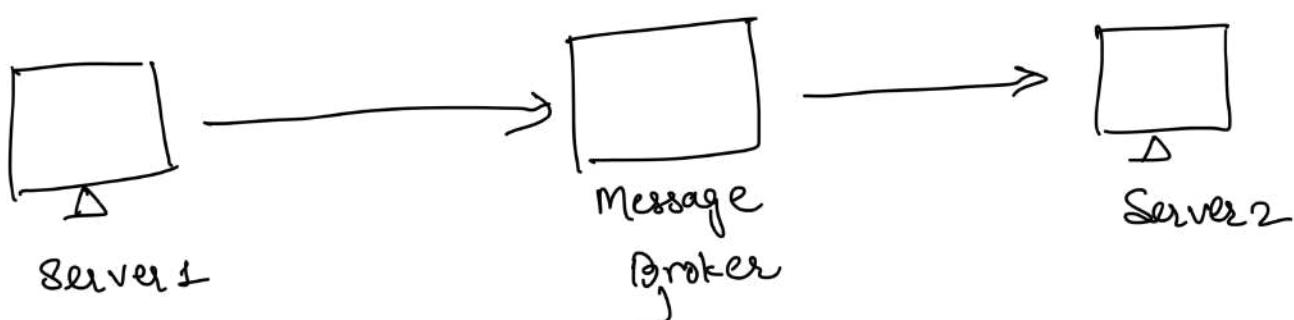
Data evolution → servers will be updated first and then clients

- ↳ backward compatibility for requests.
- ↳ forward compatibility for response.



Message Passing Dataflow

↳ asyn processing



- ↳ temporary storage, can store msg for sometime if receiver is not available → makes system reliable
- ↳ Receiver Server can consume msgs on its rate.

- ↳ one msg can be sent to multiple receivers.
- ↳ Sender doesn't wait for msg to be delivered to receiver

Example Kafka, RabbitMQ, Amazon SNS - SQS

Actor framework

- ↳ programming model for concurrency in single process.
- ↳ Each actor represents one client.
- ↳ Actors communicate with each other by sending msgs async.

Example

Akka, Erlang OTP

Subscribe for more such content.

YT channel - Ms Deep Singh

Happy Learning 😊

Designing Data Intensive Applications

CHAPTER 5

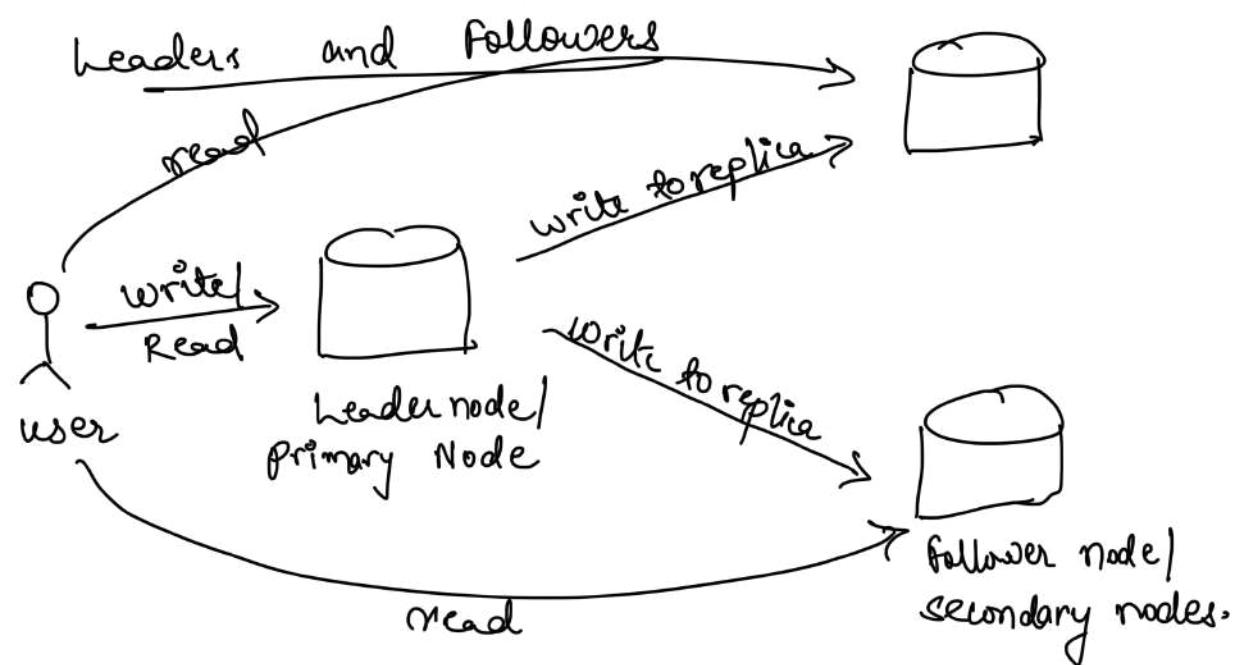
Replication

Keeping copy of data at multiple places.

why?

- ↳ keep data in some region as well \rightarrow reduce latency.
- ↳ Availability \rightarrow system is operational if certain machine goes down
- ↳ can serve more read TPS.

How to Replicate?



Synchronous vs Asynchronous Replication

- | | |
|--|--|
| <ul style="list-style-type: none">① User is sent ack once write is completed on leader & followers② Data is always consistent | <ul style="list-style-type: none">① User is sent ack once write is complete on leader node - follower nodes are updated in background② eventually consistent. |
|--|--|

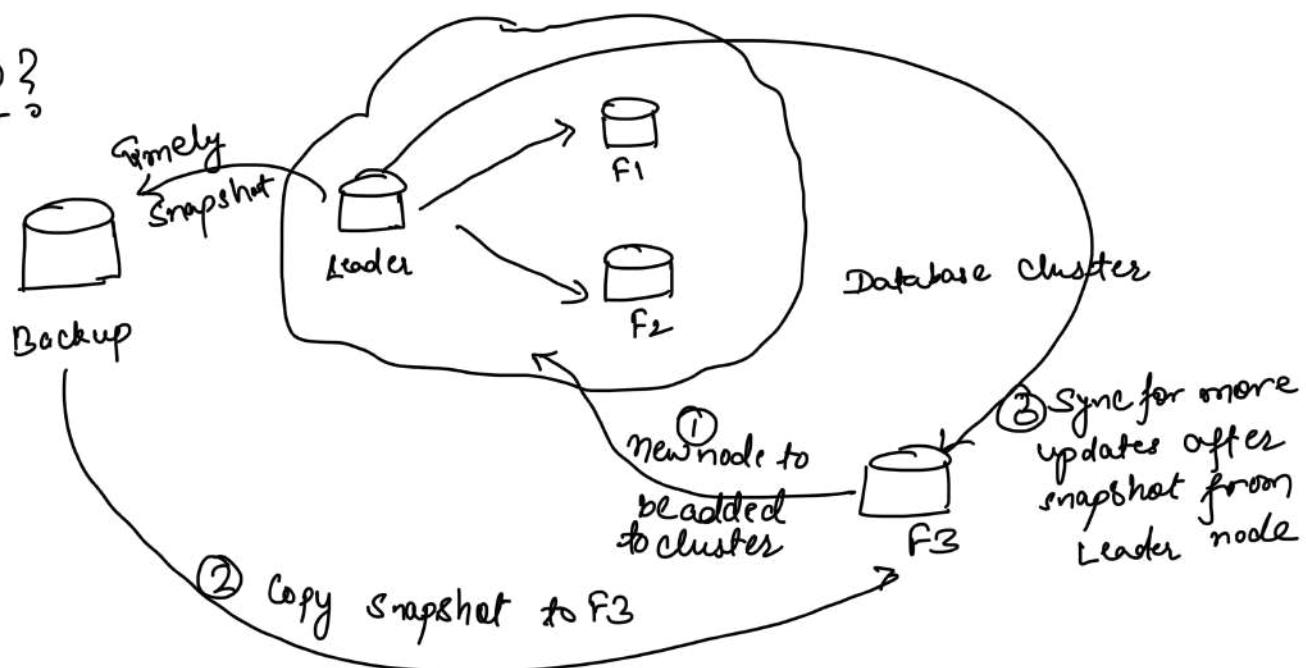
③ write can't be processed if any nodes goes down.

④ For write → leader node should be available.

Adding new follower

- ↳ in case more replicas are required to handle load.
- ↳ replace failed nodes in cluster.

How?



Handling Node Failures

① Follower Failure

- ↳ keep log of data changes on local disk. Once restarted, recover using log.
- ↳ sync with leader node for updates during downtime.

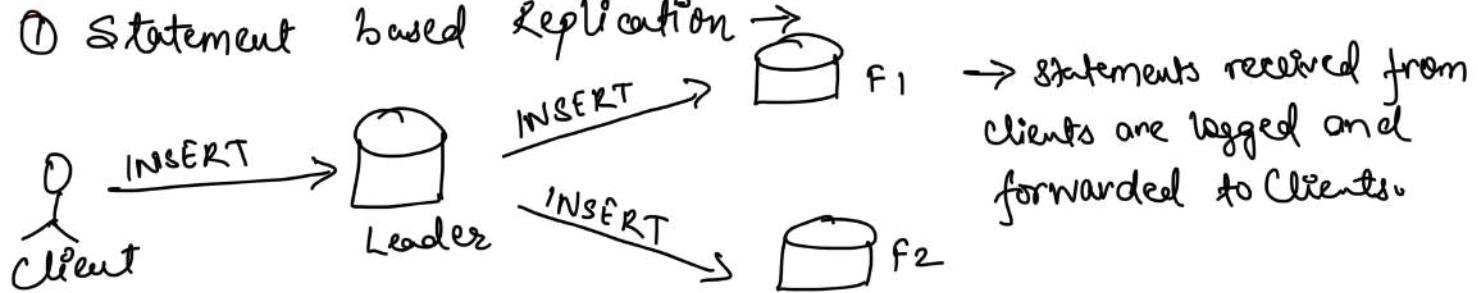
② Leader Failure

- ↳ A follower is promoted as new leader.
- ④ Identify leader has failed → health checks
- ⑤ Choose new leader → best candidate would be node which has most recent replicated data.

- ⑥ System use new leader → new leader accepts writes from clients.
 - ↳ If crashed node comes back, it is promoted to follower.
- {more details on above in coming discussions}

Replication Logs

① Statement based Replication →



→ statements received from clients are logged and forwarded to clients.

Problems

- ↳ functions like NOW() or RAND() can store different value on nodes.
- ↳ if query depend on existing data, they must be executed in same order on each replica. → else can cause different effect.

② WAL - Write ahead log

Append only file to maintain the operations on DB → discussed in CH-3

③ Logical (Row based) Replication

Use different log format for replication & storage engine.

↳ issue in #2 where WAL is used by replication & internal BTree.

① Row Insert → new values for all columns

② Row delete → primary key of row

③ Row Update → new values for columns that changed.

④ Trigger Based Replication

write custom code that is triggered for any operation in DB.

Replication Log

↳ application can see outdated data in case of async replication.
 ↳ eventual consistency.

① Read After Write consistency → "read your own writes"

The user who made an update, will see recent data on read.

This might not be the case for other users.

↳ a Decide when to read from leader and when from Followers.

⑥ for 1 minute after write, read only from leader.

what if same user reads info on multiple devices?

↳ can user requests routed to same datacenter to handle this.

② Monotonic Reads

Users gets different data on subsequent reads.

↳ can same user always read from same replica.

③ Consistent Prefix Reads

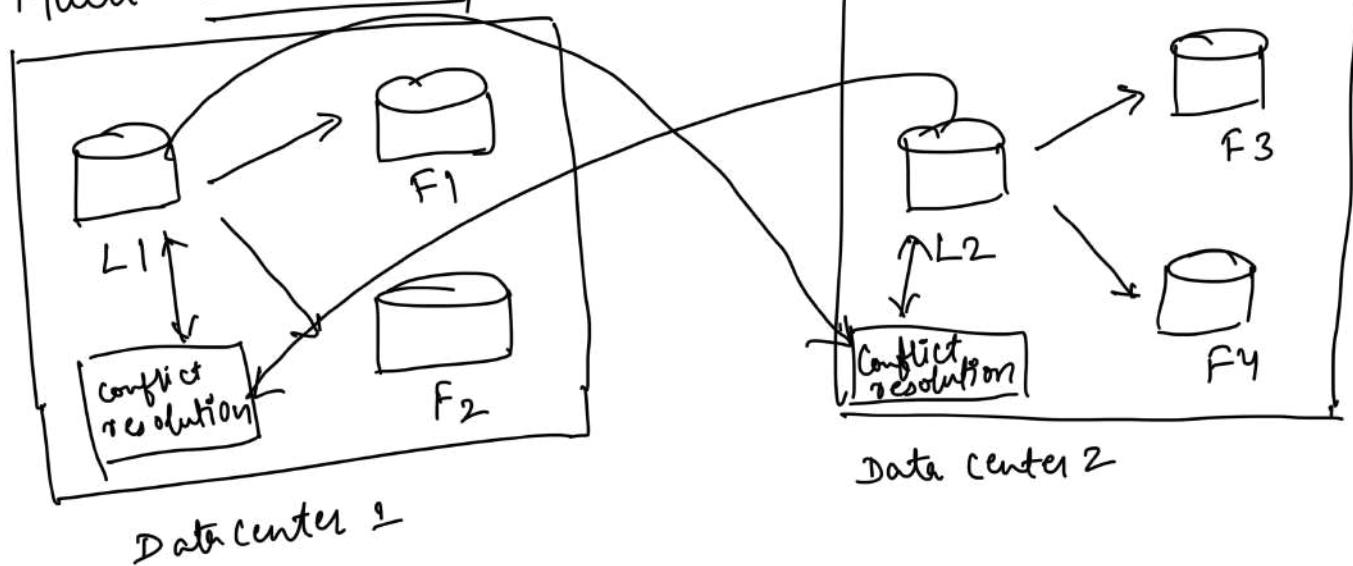
if writes are performed in certain order, reads should happen in same order. → violation of causality.
→ generally happens in partitioned databases.

Multi Leader Replication

↳ write can happen on any of leader nodes.

↳ Each leader simultaneously acts as follower to other leaders.

① Multi data center operation



Benefits

- ① latency improvement - write can happen in specific data center.
- ② fault tolerant → data center outage.
↳ network issues

Downsides

- ① same write can happen in both data centers concurrently.
↳ requires conflict resolution.

Examples

- ↳ Clients with offline operation - calendar application on multi-devices.
- ↳ collaborative editing - Google docs - conflict resolution for concurrent users.

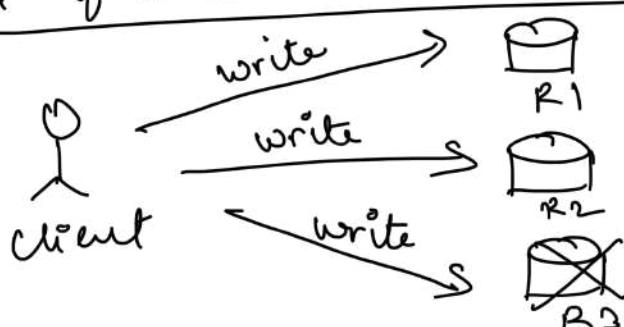
Handling of write conflicts

- ① Conflict Avoidance → avoid conflicts to occur.
 - ↳ requests from particular user are always routed to same leader
 - ↳ single reader in user's view
- ② Converging to Consistent State
 - ↳ assign each write a unique ID. example Timestamp.
 - conflicts can be resolved basis "Last Write Wins".
 - ↳ assign each replica unique ID.
 - writes from higher number replica takes precedence over other.
- ③ Custom Application Code
 - ↳ write your custom conflict handler.

Leaderless Replication

- ↳ any replica can accept writes from clients.
ex - Cassandra, Voldemort

What if a node is down on write?



- ↳ only R1 & R2 are available
- ↳ write is successful basis quorum
- ↳ on Read as well, request is sent to all nodes, and then recent value can be picked up basis version no.

How will be R3 in sync once it comes online?

- ① Read Repair → on read, we identified which Replica has stale data, so write back to R3.
- ② Anti-entropy process → background process to check differences between replicas.

Quorums

n - replicas

w - no of nodes to confirm on write for it to be success

r - no of nodes to be queried for read.

$$w+r > n$$

↑ we'll get up-to date reads.

General recommendation

$n \rightarrow$ odd number

$$w,r \rightarrow (n+1)/2$$

① If you set lower values of r and w ($r \neq w \leq n$) \rightarrow possibility of stale reads.

② even with $(w+r > n)$, there is possibility of stale reads.

③ Sloppy quorum \rightarrow

↳ in case of network disruptions

↳ will you return errors?

↳ or return response (maybe stale data) without reaching quorum.

↳ Sloppy quorum

reads and writes still require r and w successful responses but the response might not be from designated node.

↳ return temporary response and once node is back, restore it to correct state. ↑
handed handoff.

④ Two concurrent writes \rightarrow

multiple nodes will have confusion around which is most recent write.

⑤ Last write wins \rightarrow each replica stores most recent value and overwrites older values.

- ↳ as we say writes are concurrent, so write order is undefined. Attach timestamps to each write and choose most recent one.
- ↳ This might lead to date lose.
- ↳ Approach for handling this could be to assign version numbers to each write operation and subsequently merge basis that on re-writes.
- ↳ The old values are deleted directly, kept in database with help of delete markers - called tombstones.

Hope you enjoyed the chapter Summary.

Subscribe for more such content.

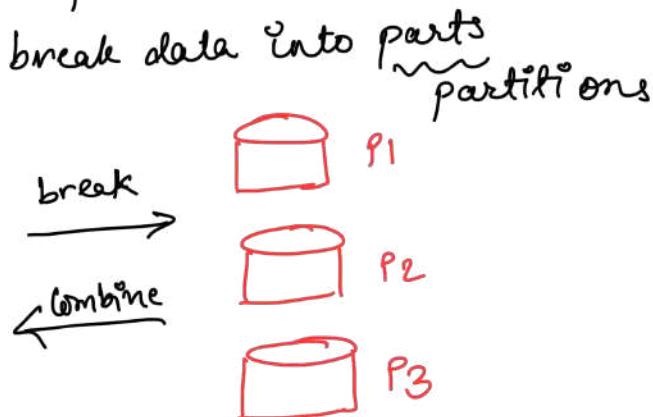
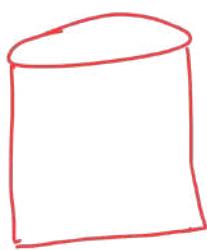
YT channel - Ms Deep Singh

Happy Learning 😊

Designing Data Intensive Applications

CHAPTER - 6

Partitioning



why? to scale the system. After certain point, data can't fit into single machine.

↪ the partitioned data should also be replicated on multiple nodes for fault tolerance. → high availability

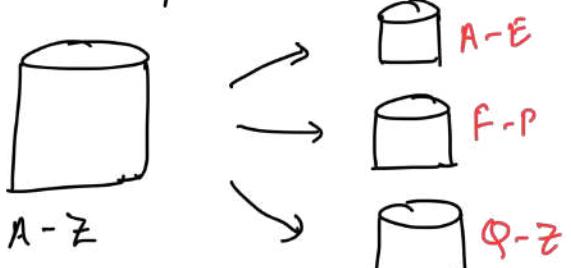
Partitioning of Key-Value Data

↪ Idea is to make sure data is evenly distributed among all nodes.

↪ If data is unevenly stored → called Skewed

① Partitioning By Key Range

Each partition can be assigned a continuous range of keys.



① within each partition, keys are sorted.
(sstable & LSM tree)

② range queries are easy.

③ Some access patterns can cause hotspots. example Timestamp.

Attach specific prefix
to timestamp to solve issue.

② Partitioning By Hash of Key

↳ to avoid skew and hotspots.

Eg. Cassandra & MongoDB use MDS.

↳ assign each partition range of hashes.

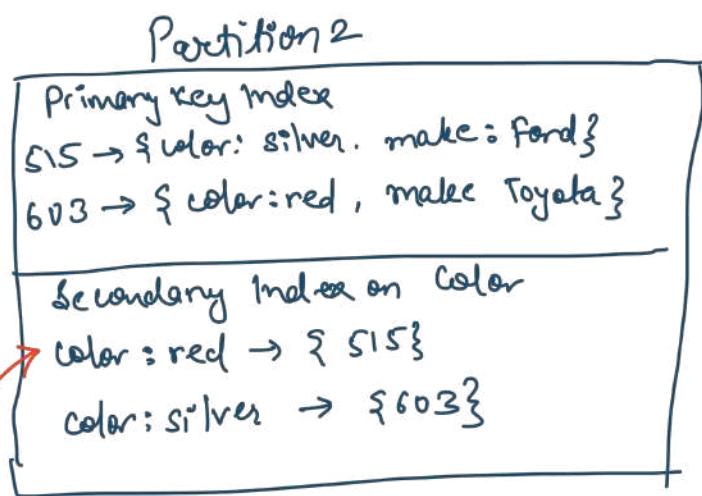
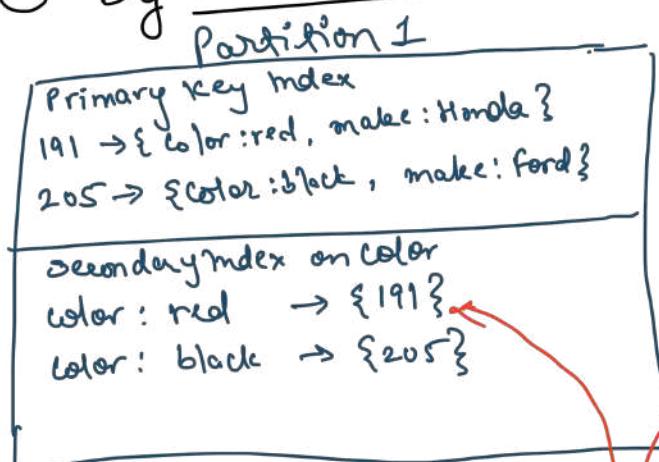
↳ key will be stored in partition whose hash falls in specific range.

Disadvantage

↳ capability of efficient range queries.

Partitioning Secondary Indexes

① By Document



for Red car → gather data from both partitions.

- Each partition is completely segregated, create local indexes.
- Read queries on secondary indexes could be expensive → gathering data from all partitions.
- used in MongoDB, Cassandra, Elastic Search

② By Term

↳ construct global index which covers data in all partitions.

↳ global index is also partitioned.

↳ you can use similar to range partitioning.

Benefits

↳ reads are more efficient. → hit only specific partition to gather data.

Downsides

↳ writes can be heavy. → write to single document can affect multiple partitions
↳ terms in doc located on different partitions

↳ example DynamoDB global Secondary Indexes.

Rebalancing Partitions

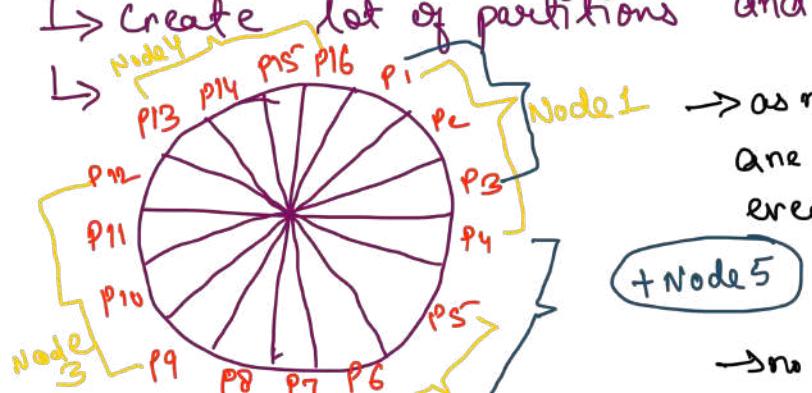
- ↳ add more CPU to handle load.
 - ↳ add more disks.
 - ↳ replace machine.
- ↳ requires movement of data from one node to another.

Strategies

① Hash Mod N → not effective approach as number of nodes (N) will change, hence changing node number where data will be stored.

② Fixed Number of partitions → used in Riak, Elasticsearch, couchbase, Voldemort.

↳ Create lot of partitions and assign these partitions to nodes.



→ as new nodes are added, few partitions are assigned to node to make distribution even.

→ no of partitions are fixed and only

Node

assignment changes.

→ partition size should be not too small OR not too large.
It's important that you come up with number after proper analysis. management overhead rebalancing on node failures is expensive.

③ Dynamic Partitioning

↳ Key range partitioned databases (HBase & Rethink DB) create partitions dynamically.

- as partition grows and exceeds certain limit, it is split into two parts.
- on data deletion, if required, partitions can be merged

④ Partitioning proportionally to Nodes

In #3 → no of partitions proportional to size of dataset.

In #2 → size of partition is proportional to size of dataset.

In Cassandra → no of partitions proportional to no of nodes.

↳ fix no of partitions per node.

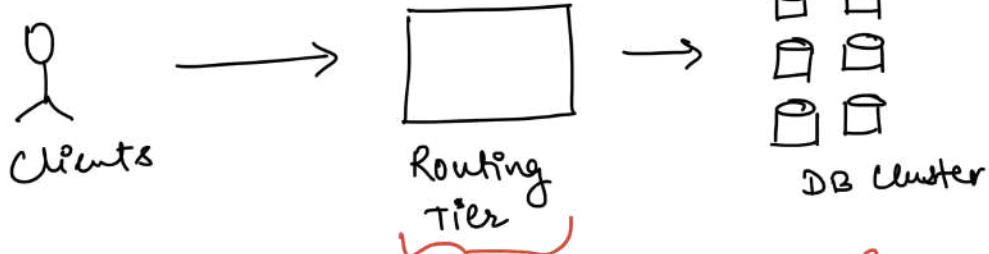
→ Add new node → partitions becomes smaller as data is moved.

- randomly choose fix no of partitions to split
- After split, take ownership of half of them and leave other half in place.

Request Routing → How will you serve read query as rebalancing happens and data moves from one node to another?

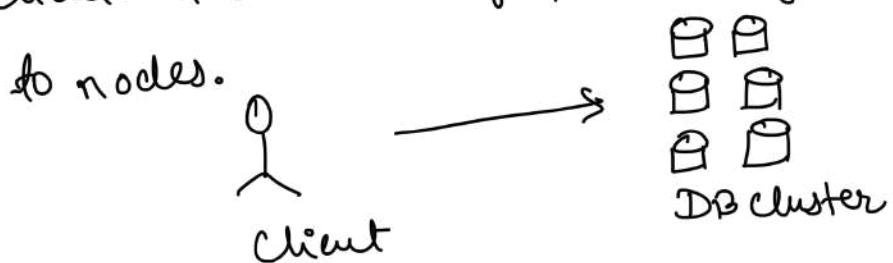
- route request to any node via round robin. If data not found, move to another node.

②



*determine which node [partition aware load
should handle request] balancer.*

③ Clients are aware of partitioning and assignment of partitions to nodes.



How will you know assignment of partitions to nodes?

↳ Use coordination Service example Zookeeper.

↳ ex. HBase, Kafka

Keeps mapping of partitions to nodes.

↳ Use gossip protocol between nodes. eg. Cassandra.

↳ own routing tier. eg. CouchBase we routing tier called moxi.

Subscribe for more such content.

YouTube Channel - msDeep Singh

Instagram | LinkedIn | GitHub - msdeep14

Happy Learning 😊

Designing Data Intensive Applications

CHAPTER - 7

Transactions

Youtube channel - MS Deep Singh

Subscribe for more...

Transaction → group multiple reads and writes into single logical unit
↳ executed as single operation
↳ either all succeeds or all fails.

ACID → Safety guarantees provided by transactions.

① Atomicity → atomic operation

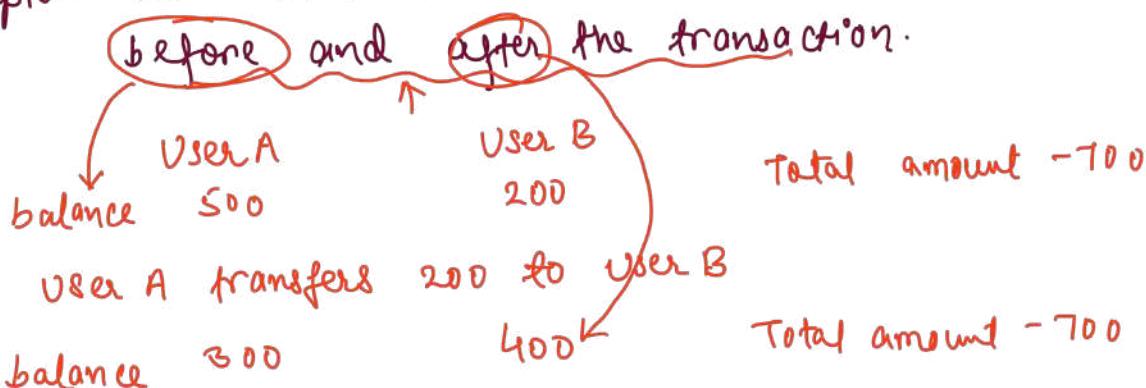
↳ can't be broken into smaller parts.

↳ successful if all parts of transaction are complete.

↳ for any fault, it will fail.

② Consistency → correctness of database.

Example - Transaction amount between two accounts should be balanced



↳ consistency is more of application property. Atomicity, Isolation and durability are database properties.

application's responsibility to achieve consistency utilizing A.I.D from ACID.

③ Isolation → concurrent transactions occur independently in DB without compromising consistency.
each transaction pretend that it is only transaction being performed
the result will be same if these transactions are carried out sequentially.

④ Durability → Once transaction is committed, it will be intact even in case of system failures.
ex-usage of WAL in case of B Trees.

Replication and Durability → no technique is 100% risk free.
↳ we can just take measures to reduce it

single object writes

↳ you're writing a 20KB JSON to DB, what if network issue after 10KB is written?

↳ storage engines in general provide atomicity & isolation at single object level.
log for crash recovery - B-tree example
lock on object.

multi-object transactions

use-cases
① In RDBMS, row in one table has foreign key reference to row in another table.

- ② multi field updates in document.
③ secondary indexes

Retrying aborted transaction

- ① operation completed on server but client didn't get response due to network error.

- ① If error is due to overload, retry might not help.
→ use exponential backoff & jitter.
- ② Retry for only temporary errors, not permanent.

Weak Isolation Levels

- ↳ A reads data and B modifies data at some time.
- ↳ A and B try to write data to same record concurrently.

Transaction Isolation → System works in manner as everything is working in sequence.

↑ Serializable isolation.

- ↳ As strong isolation support has performance cost, some databases tend to use weak isolation levels.

① Read Committed

- ↳ On Read, you'll see data that is committed.
→ no dirty reads
- ↳ On write, you'll only overwrite data that is committed.
→ no dirty writes

use locks } To avoid this issue, transaction B can wait till A is committed or aborted.

- ① use locks → can lead to more wait time.
- ② Hold both old and new value for which update is going on.

↑ return old value till new one is committed.

Disadvantage → Read Committed can cause anomaly on read

→ Read skew
Timing anomaly

Example

Account A

\$500

transfer \$100

Account B

\$500

Total amount mandeep has in two accounts - \$1000

For time being → \$400
if read from accounts

\$500

↑ New value not yet committed

Total amount - \$900 → This issue could occur in - time on Read.

Solution to above

② snapshot Isolation → Each transaction reads from

consistent snapshot of database.

extremely beneficial for long running read only queries
↑ backups / analytics

How to implement?

↳ write locks to prevent dirty writes.

↳ no locks required on reads.

↳ maintains several versions of objects (unlike 2 in read committed)
↳ Multi Version Concurrency Control (MVCC)

indexes with snapshot isolation

↳ index point to all versions → index query can filter out object versions not visible to current transaction.

↳ use append only / copy on write B-Tree pages.

↳ don't overwrite pages, create new copy of modified page.

↳ write transaction creates new B-Tree root
↳ old root is consistent snapshot.

Lost Update Problem



- ① read from db
- ② modify value
- ③ write back



two transactions perform concurrently
↑
one modification could be lost

Solutions

① Atomic write operations → atomic update operations.
↳ by acquiring lock

⑤ Explicit locking → handle locking mechanism in application.



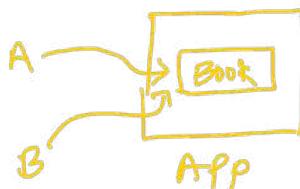
⑥ Detect lost updates → allow to execute in parallel and fail if lost update is detected.
can be retried

⑦ Compare and set → On update, check the current value is same as what it was read earlier.

Write Skew

Example - Book meeting room.

- ↳ meeting room should not be booked twice.
- ↳ Two users tried to book at same time.
- ↳ check availability → snapshot isolation
- ↳ room is available, so booked for both?
 - use serializable isolation
 - lock multiple rows that are accessed for transaction.



Phantom → write in one transaction change result of search query in another transaction.

Serializability

Serializable Isolation → transactions may execute in parallel but they have same effect as executed sequentially.
↑ no concurrency

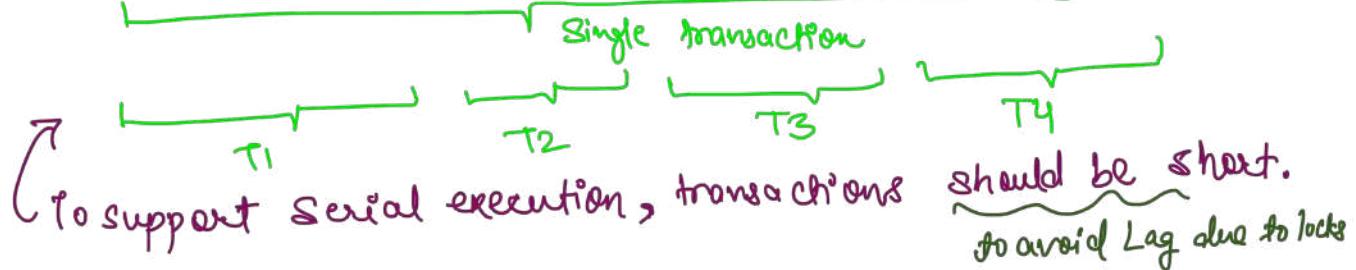
① Actual Serial Execution → more concurrency.

- ↳ no overhead of locking.

② Encapsulate transactions in stored procedures

↳ transactions should be short.

Book flight → look from flight → select → book seat → payment

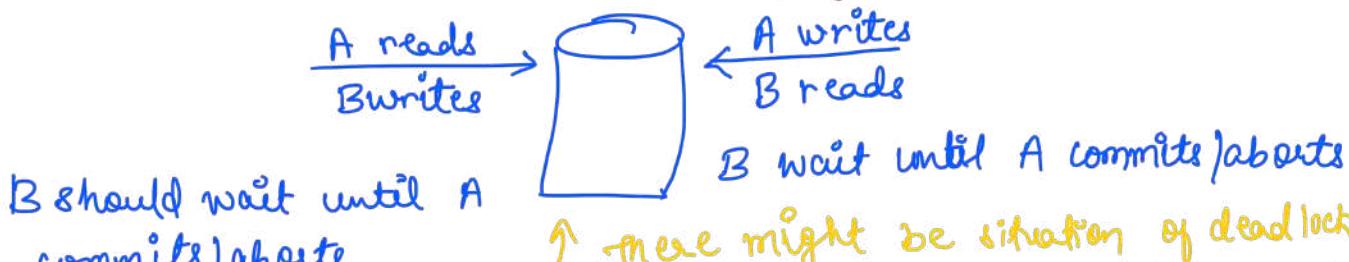


↳ application should submit entire transaction code to DB ahead of time.

no I/O wait
no concurrency overhead.

③ Two Phase Locking (2PL)

↳ Several transactions are allowed to write to an object unless no-one is modifying it.



↑ there might be situation of deadlock if lock is not released. → transaction is aborted by DB in this case.

See you again in next chapter..

Subscribe more more such content.

YouTube - MS Deep Singh

Happy Learning 😊

Designing Data Intensive Applications

CHAPTER 8

The Trouble with Distributed Systems

Subscribe for such content

YouTube - Ms Deep Singh

Faults and Partial Failures

- ↳ Single computer in-general works as per written software.
- ↳ In distributed system, there is possibility some parts are working, partial failure

Unreliable Networks

Assuming a shared nothing architecture, things that could go wrong

- ↳ Request is lost / waiting in queue.
- ↳ Remote machine failed / temporarily unresponsive.
- ↳ Response took more time than expected.

Each machine has its own memory/disk and not accessible to others.

} add timeout so as request don't keep on waiting.
It have occurred due to any network faults.

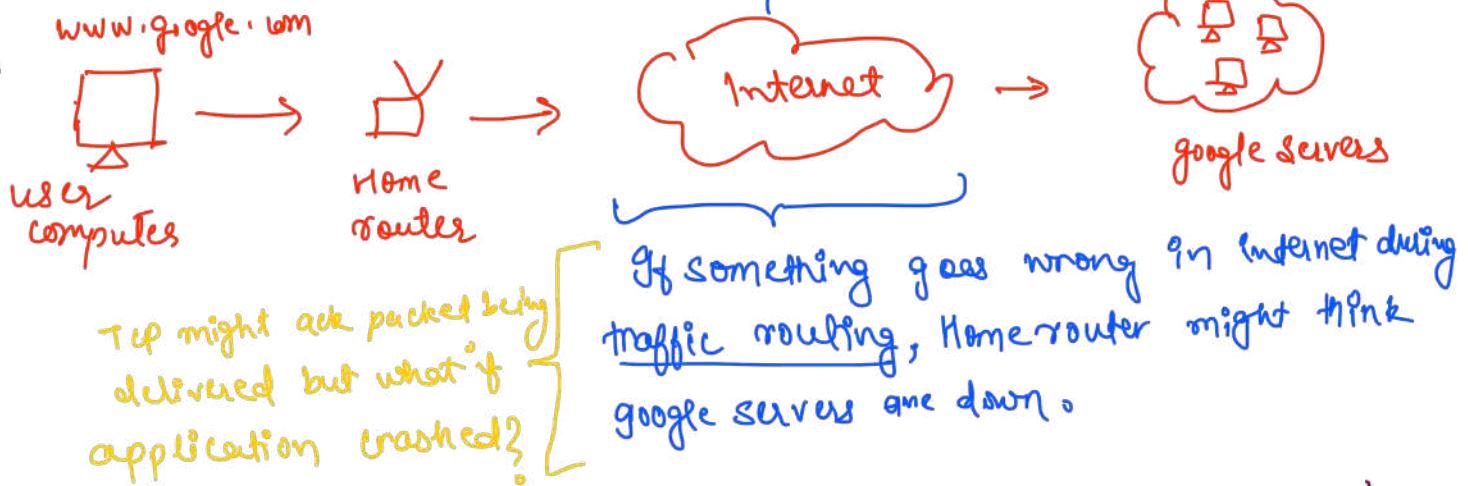
Detecting Network Faults

Automatic detection

- ↳ Load Balancer don't send request to faulty node
- ↳ If leader fails in DB, follower is promoted to leader.

What if we're not able to figure out node is dead or not?

- ↳ e.g. Node OS is running but node process crashed.
another node should takeover before



General recommendation is to retry after timeout.
exponential backoff

Timeouts

and Unbounded Delays

for asynchronous networks, there is no upper limit for packet to be delivered

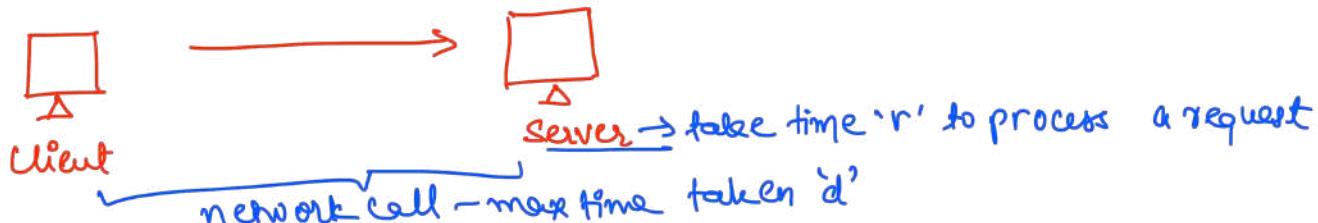
What should be timeout value?

not too big, not too small.

users have to wait long

node may be actually operational but slow to respond.

e.g. due to overload



You can timeout after $d+r+d = 2d+r$

Network Congestion

and queuing

until the requests are served, they are kept in queues

queuing can also be at client end

TCP flow control
backpressure

TCP vs UDP

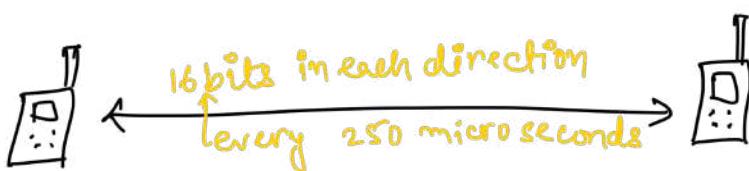
→ perform flow control

→ retry - retransmit lost packets

→ don't perform #1 and #2
→ we don't want delayed data
e.g. live cricket match.

Synchronous vs Asynchronous networks

fixed line Telephone network



- fixed guaranteed bandwidth is allocated
- This is synchronous network.
- ↳ no queuing - 16 bits reserved
- ↳ bounded delay → latency is ~fixed.

↳ This concept is not used in TCP.

Should we use packet switching over internet?

Should not be used for bursty traffic.

no → we don't have exact bandwidth requirements.
you'll end up wasting → depends on use-case.
bandwidth or allocating 10x bandwidth.

Unreliable Clocks

↳ Clock is used at multiple places in system. example - calculate latency of operation.

↳ As there are multiple machines in network, time calculation on each machine might vary.

for synchronization of clocks, mostly Network Time Protocol (NTP) is used

① Time of Day Clocks

↳ current date and time according to some calendar.

Example JAVA - System.currentTimeMillis()

wall-clock time
January 1, 1970
midnight UTC
return no of second/ms since epoch

↳ In scenarios of local clock is too far ahead of NTP, it will jump back to previous point in time.

↳ not a good solution for calculation of elapsed time

② Monotonic Clocks

→ guaranteed to always move forward

JAVA - System.nanoTime()

↳ Not dependent on synchronization between different nodes.

↳ NTP will not cause to jump back or forward.
↳ no synchronization is required.

Example on How NTP can go wrong?

> NTP daemon is misconfigured on servers.

↳ firewall is blocking NTP traffic.

Clock offsets between servers should be monitored.

↳ In case clock on one node drift too much from others → node can be declared dead.

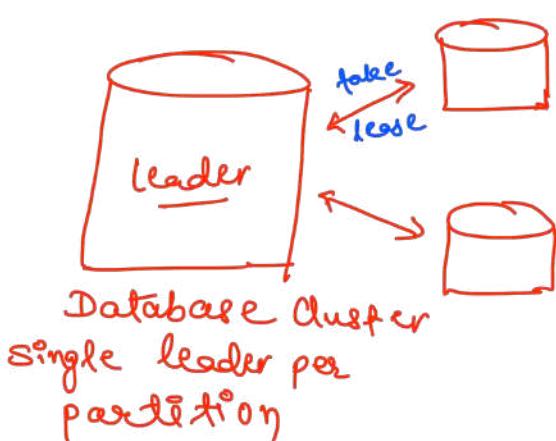
will depend on business your servers are spawned for.

Example - Last Write Wins algo on DB cluster.

For scenarios where ordering of events is important, Logical Clocks can be used.

why? → based on incrementing counters, it helps in identifying relative order of events instead of time elapsed.

Example → How System Pause can cause issues



who is leader → A node in coordination can take lease for certain timeout

needs to be updated once timeout expires
How? via synchronized clocks
if nodes don't update lease, it is assumed dead and another node should take over.

In this perfect setup, what could go wrong?

① lease time depends on time taken from another machine.
what if clocks are out of sync → can cause issues

② To solve #1 - use monotonic clocks

what if there is unexpected pause in system
① example due to Garbage Collector (GC)
② OS access disk for accessing memory pages (virtual memory)

→ The thread is paused as I/O takes place

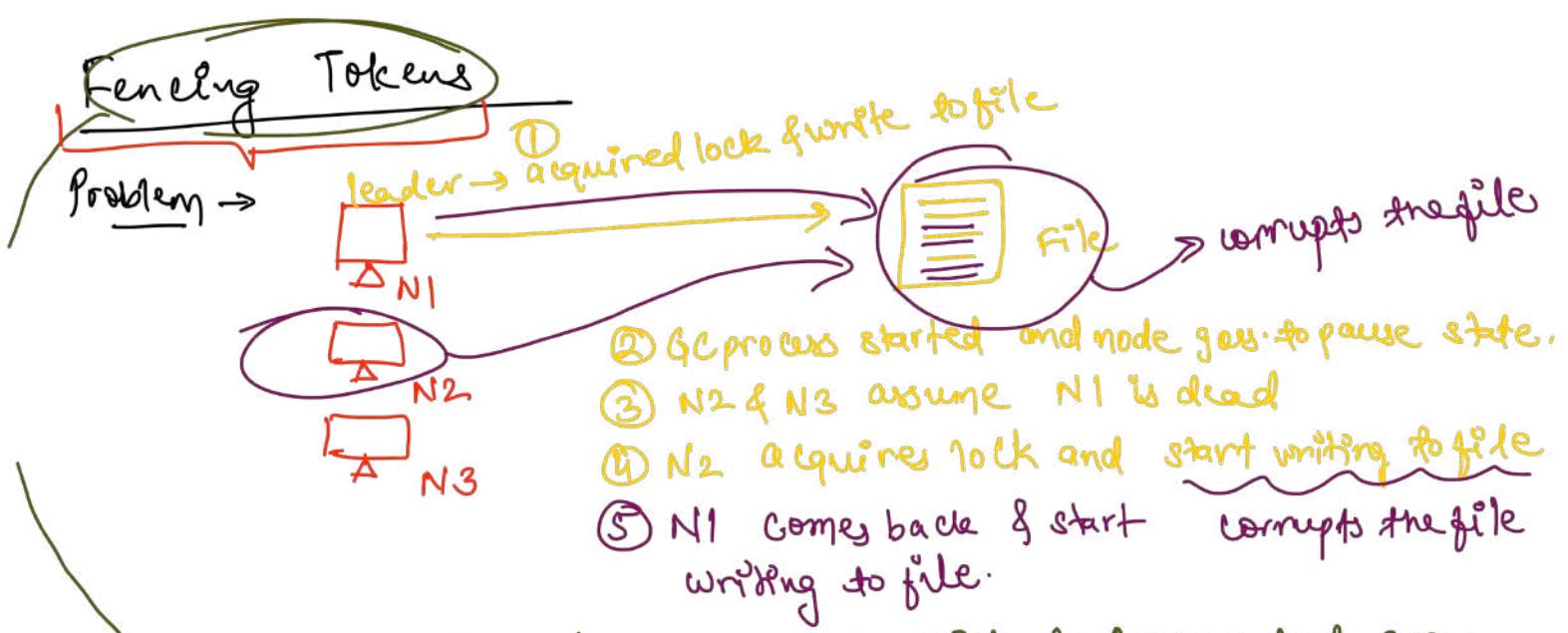
If lot of time is spent in doing I/O rather than actual processing → Thrashing

③ Can we treat GC as

planned outage } → no requests served by node in this time.

② Use GC only for short-lived objects and restart processes periodically.

→ will help in reducing impact on application.



with every lock, keep a counter which is incremented every time a lock is granted for write.

In above scenario

N1 is granted 10
N1 paused → as N1 tries to write back with token as 10
N2 granted 11 it is rejected.

what if N1 update the token at its end and send 11 instead of 10.
By Zantine Faults

safety and liveness → response is sent by server if certain conditions are met.
↳ example - quorum of nodes.

Even if entire system crash / network fail,
algorithm should always return
correct result.

See you in next chapter...~

Happy learning 😊

Designing Data Intensive Applications

CHAPTER - 9

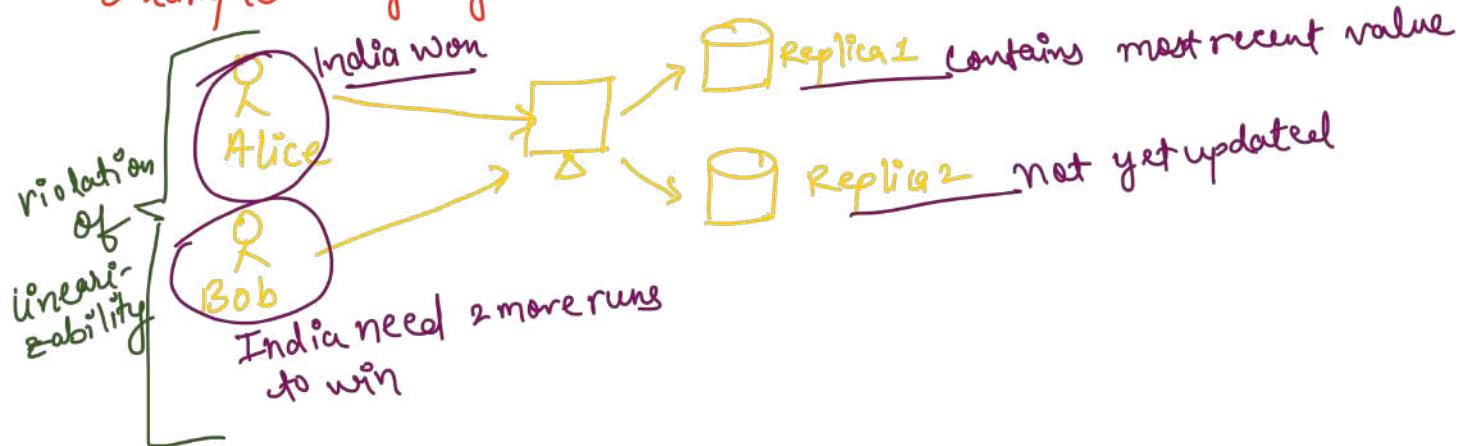
Consistency and Consensus

Subscribe for more such content . . .

YouTube - MsDeep Singh

Linearizability → There could be multiple replicas in database but to client it appears like a single node. provide strong/atomic consistency.

Example - ongoing cricket match



Making system linearizable

$x \rightarrow$ current score of India

$$x = 55$$

↓

$$\underline{x = 61}$$

If updated value is sent in response to one client, then Clients all clients should be sent updated score.

Even if write is in progress.



Comparing Linearizability with Serializability

recovery guarantee on reads

consistency property of transaction

and writes of an object.

Transaction behaves as same if executed sequentially.

why linearizability?

① Locking and leader election

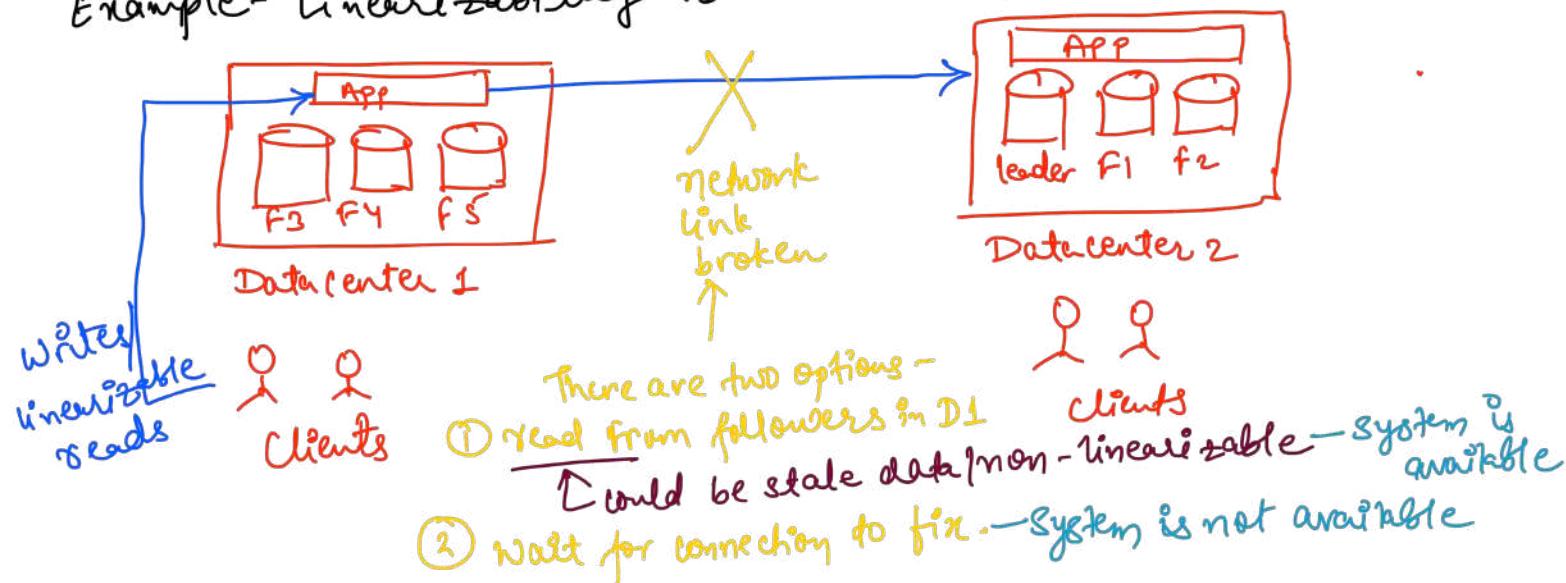
↳ for single leader replication - should only be one leader not more → split brain
a node can acquire lock on becoming leader.
linearizable. Example Zookeeper, etcd.

② Uniqueness Guarantees → Example username should be unique for user. should be linearizable operation.

Implementation

- ↳ Keep only single copy → operation will be atomic by default.
System is not fault tolerant.
- ↳ Consensus Algorithms → can be used to implement linearizable system

Example - Linearizability vs Availability. → Single leader replication



CAP Theorem → Choose either 2 of Consistency, Availability and Partition Tolerance. In practical, CA system is not possible.

Linearizability and network delays

- ↳ It is generally avoided (if our use-case permits) to improve performance.

It is slow even if there are no network faults.

Ordering Guarantees → operations are executed in particular order

① Ordering and Causality

↳ ordering helps in preserving causality.

→ Example

② Consistent Prefix Reads

→ question should be asked before so as it can be answered. → effect
casual dependency
↓ w que & ans

Causality imposes ordering in events

↑ cause comes first, then effect.

→ system obeys this ordering
↑ referred as causally consistent.

Linearizability vs Causality

Total order of operations

allows two elements to be compared. e.g. $13 > 5$

we can always say which occurred first.

Database abstracts out that there are no concurrent events.

Partial order

The elements can be incomparable or one could be greater than other.
e.g. $\{a, b\}$ and $\{c, d\}$

→ two operations are ordered if they're causally related.

↑ one happened before other.

↑ incomparable if concurrently occurred.
example Git version control system.

↳ Any linearizable system will preserve causality correctly.

Stronger than causal system
hence performance can be slow.

Sequence Number Ordering

↳ maintaining all causal dependencies would be large overhead.

→ for event ordering, sequence no or timestamps can be used.
↑ from logical clock.

assign seq no to each operation → use counters for each operation.
operation & these no can be compared.

Sequence number for non-causal example multi leader replication

① Each node generate its own numbers.

node should have some identifier so as numbers are unique across cluster
eg one node generate odd and other even.

② Attach Timestamp → used in last write wins algorithm.

③ Assign range of numbers to each node.

Generated seq numbers are not causally consistent:

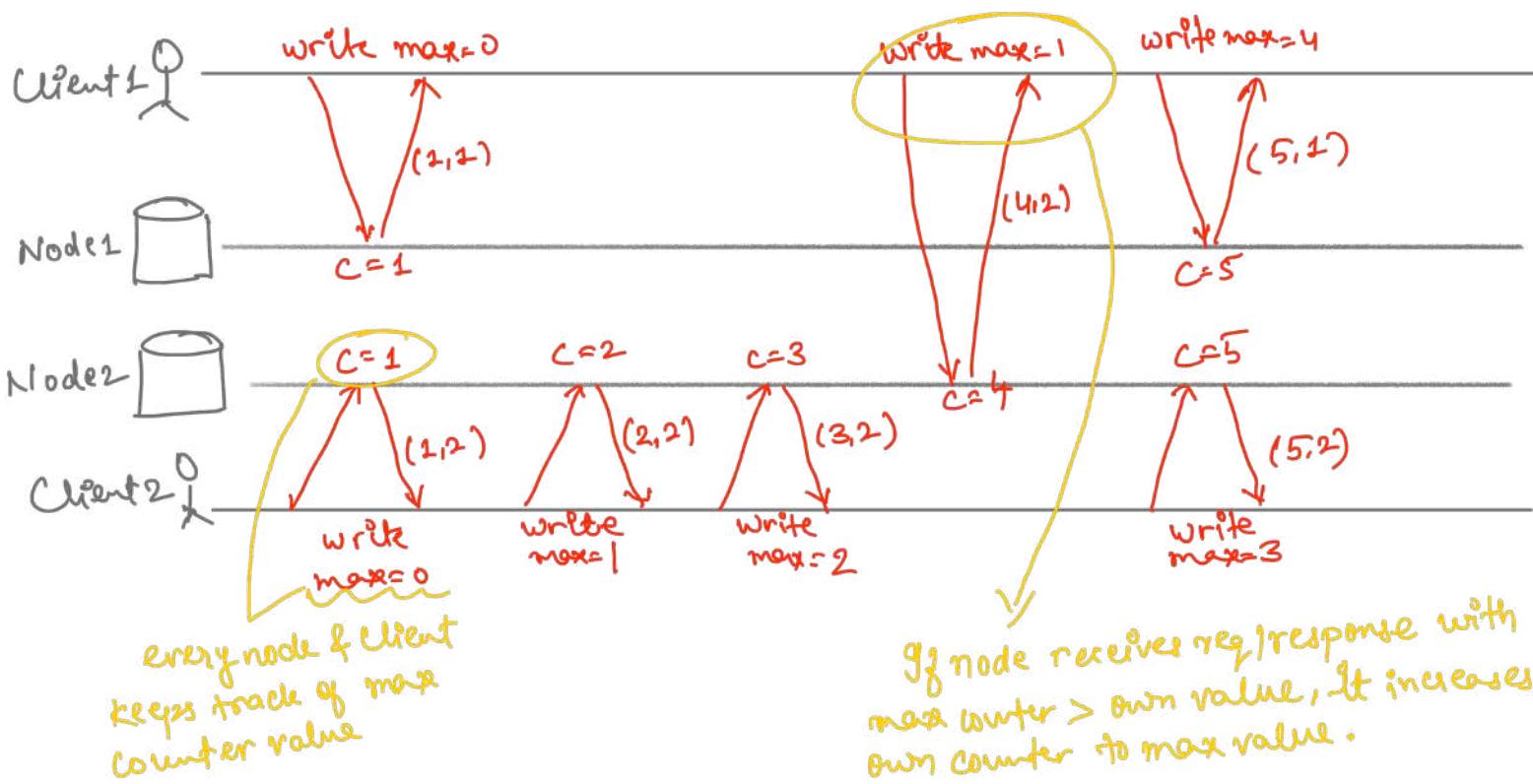
they are not capturing order of events across the nodes.

Solution - Lamport Timestamps

{ counter, node Id }

one with greater timestamp is greater

of same time, greater node Id is greater



The Same username problem

Two users trying to create account with same username at same time.
↳ username creation operation is not collected yet, Lamport timestamp will not help.

↳ node has generated operation but don't know what they are?

↳ A node should check with all other nodes if any request with same username.

↑ What if network fault between node communication.

To uniquely identify usernames, total ordering of operations is not sufficient

↳ you should when total order is finalized.

Total order Broadcast.

How can all nodes in cluster agree on same total ordering of operations for single leader database? - one node as leader handles this.
what if leader fails?

Total Order Broadcast

→ also used in consensus algorithms.

↳ protocol for exchanging msgs between nodes.
→ Reliable delivery

→ Totally ordered delivery → delivered to every node in same order

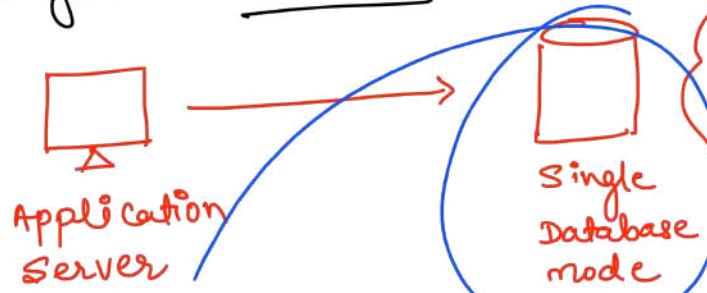
→ If msg is already delivered, node is not allowed to insert it back in earlier position.

Distributed Transactions and Consensus

→ Several nodes to agree on something
e.g. leader selection

Atomic Commit and Two Phase Commit (2PC)

① Single node Commit



- On transaction commit
① Update WAL on disk
② append commit record to log on disk

what if DB crashes before #1 or after #2?
↳ rollback changes

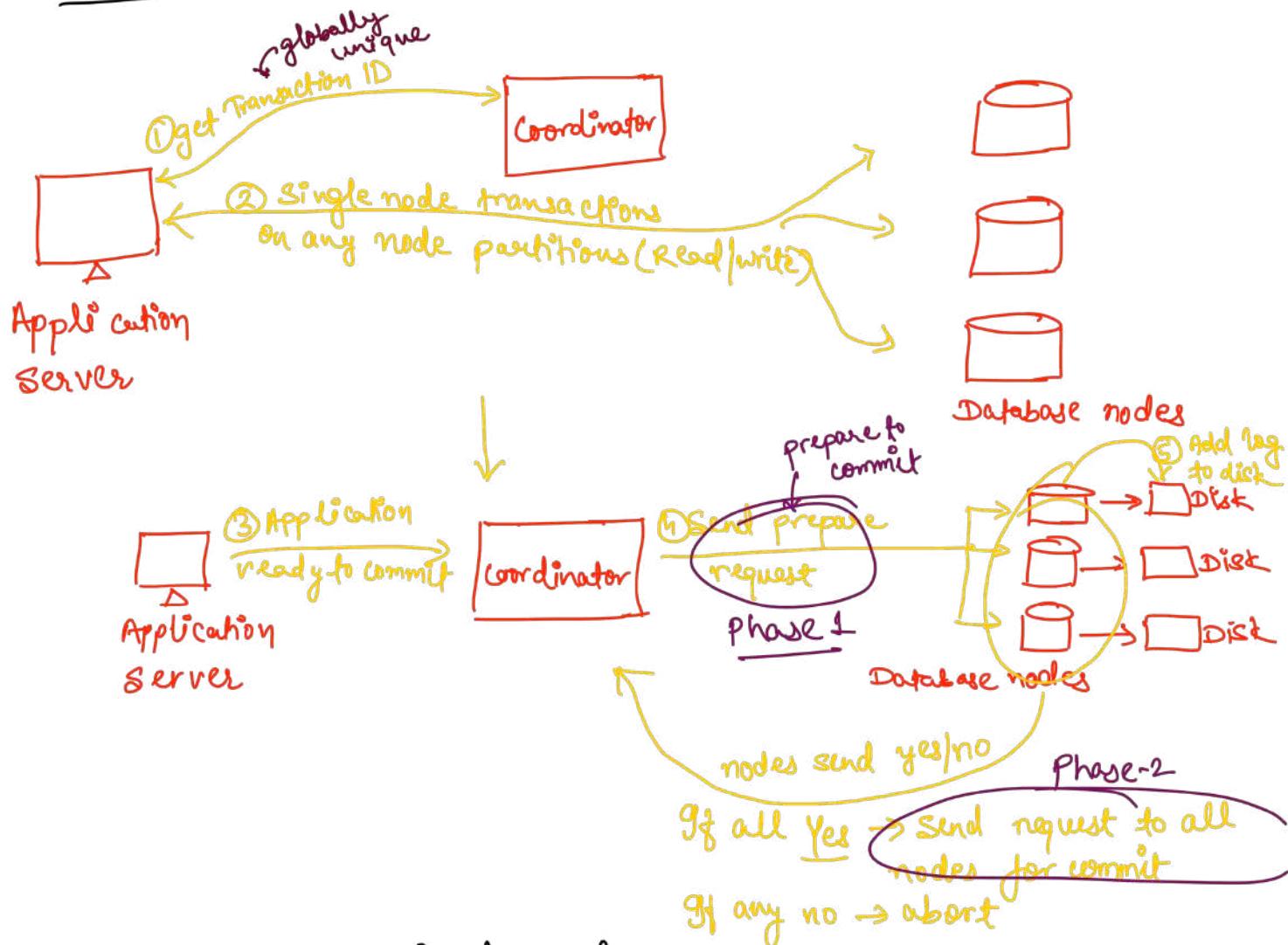
what if crashed after #2?
↳ recover from disk once node restarts.

what if multiple Nodes?

↳ we can't send commit msgs to all nodes. → what if it fails on few nodes? → network delay, node crash, etc.

Difficult to handle.

Two Phase Commit (2PC) → the commit/abort is split into two phases.



what if coordinator is down?

2 Coordinator is down after sending commit request to one node → other nodes will not know what to do with transaction but wait.

- ① Either database nodes interact with each to decide what to do with specific transaction.
- ② Coordinator dumps all actions to disk
 - what if disk is also corrupt? → manual intervention by admin.
 - so it can recover once back online.

Distributed Transaction Examples

- ① Database Internal → internal transactions among the nodes in cluster.
- ② Heterogeneous distributed transactions
 - ex: different technologies
 - two databases from different vendors, message broker & database

extended Architecture

XA Transactions → 2PC implementation across heterogeneous tech

SC API for interfacing with transaction coordinator.
in Java - Java Transaction API → network driver or client library
The driver can send required info to coordinator (prepare | commit | abort)

Fault Tolerant Consensus → safety properties

- ① Uniform agreement among nodes → no two nodes decide differently → consensus algo properties
- ② Integrity → node decides only once
- ③ Validity → nodes decides value v, then v was proposed by some node.
- ④ Termination → every node decides a value
↳ Liveness property
that doesn't crash → system should reach a decision even if node crashes

Single Leader Replication and consensus

How is leader chosen?
in case of failures → manually by operation team
or automatic selection

Can split brain problem occur?

Consensus algo are based on total order broadcast.
like a single leader replication
and it requires a leader

We need all nodes agree for leader.

↳ liveness

Spiral loop
Solution → Epoch numbering for quorum
protocols define epoch number → within each epoch, leader is unique.
If leader dead → election is given an incremented epoch number → increasing
monotonic helpful in resolving conflicts.

Coordination Services eg. zookeeper, etc

- ① Linearizable atomic operations → for concurrent operation on nodes, only one succeed using atomic compare and set operation.

② Total ordering of messages

↳ fencing tokens are used to avoid conflicts on lock lease.

↳ Zookeeper uses monotonically increasing transaction ID & version number

③ Failure Detection → via heartbeats.

④ Change notifications → new nodes added / removed

See you in next video —

Happy learning ☺

Don't forget to subscribe ~