

Designing Data Intensive Applications

YouTube - <https://www.youtube.com/@MsDeepSingh>



CHAPTER 1

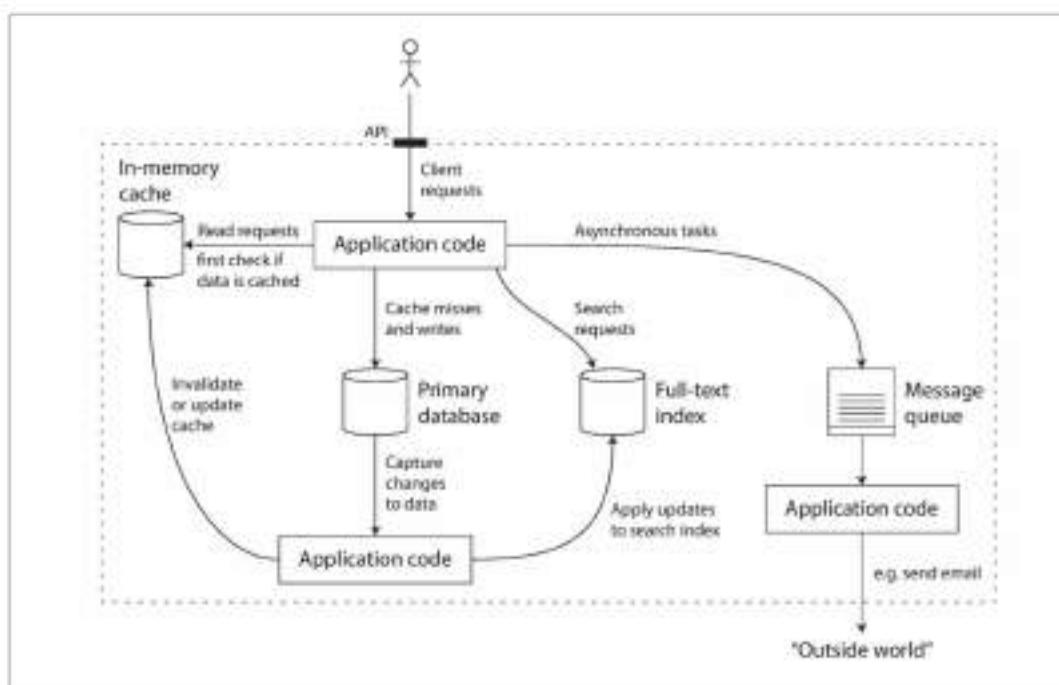
Reliable, Scalable and Maintainable Applications

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 our 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 big 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.

<http://www.linkedin.com/in/willmengole>

Bill Gates
Greater Seattle Area | Philanthropy

Summary
Co-chair of the Bill & Melinda Gates Foundation. Chairman, Microsoft Corporation. Voracious reader. And traveler. Active blogger.

Experience
Co-chair - Bill & Melinda Gates Foundation 1990 - Present
Co-Founder, Chairman - Microsoft 1975 - Present

Education
Harvard University 1973 - 1975
Lakeside School, Seattle

Contact Info
Blog: thegeeknotes.com
Twitter: @BillGates

user_table			
user_id	first_name	last_name	summary
251	Bill	Gates	Co-chair of... blogger
			region_id industry_id photo_id
us91	131		52812532

regions_table	
id	region_name
us7	Greater Boston Area
us91	Greater Seattle Area

industries_table	
id	industry_name
43	Financial Services
48	Construction
131	Philanthropy

position_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
867	251	Harvard University	1973	1975
868	251	Lakeside School, Seattle	NONE	NONE

contact_info_table			
id	user_id	type	url
133	251	Blog	http://thegeeknotes.com
138	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. → peer 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

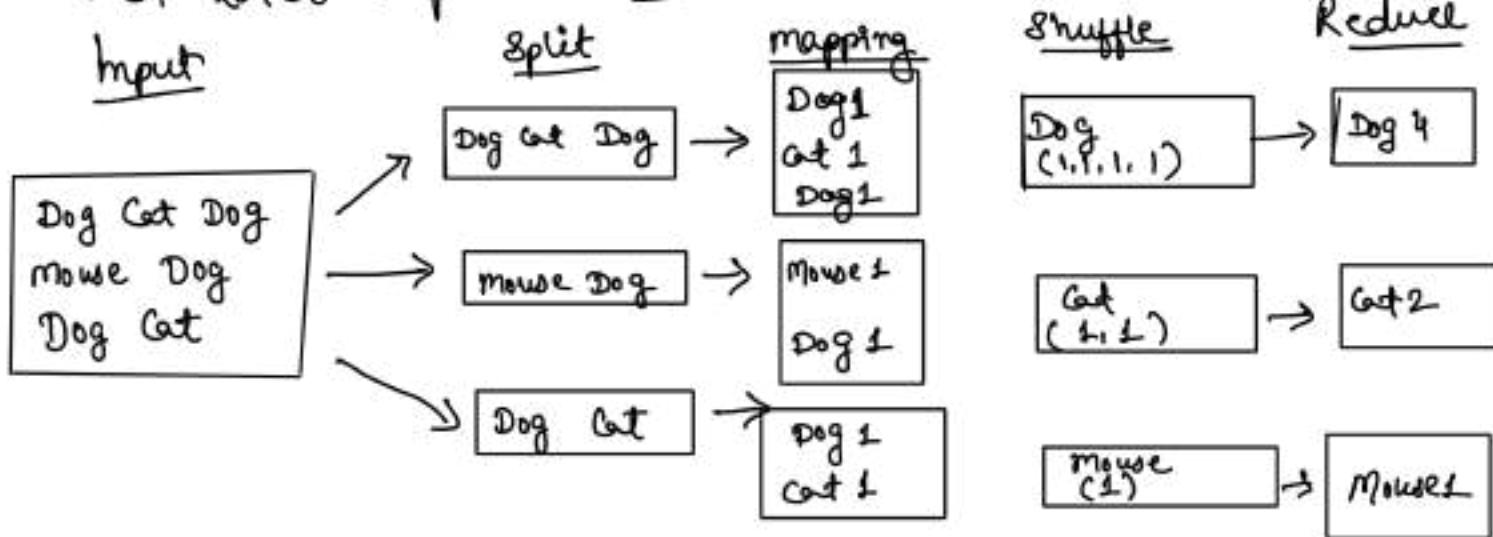
SQL - declarative query Language

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

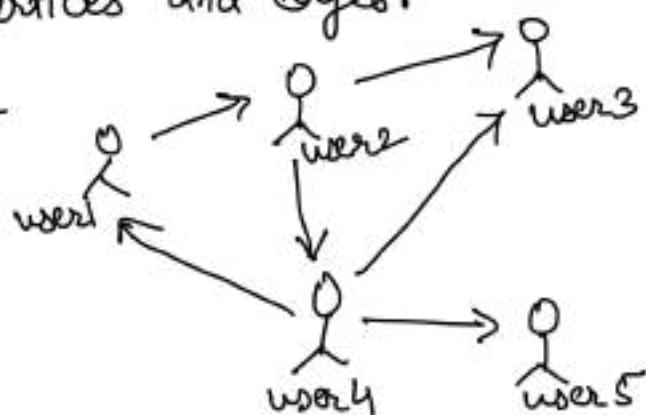


Graph Data Model

- ↳ 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



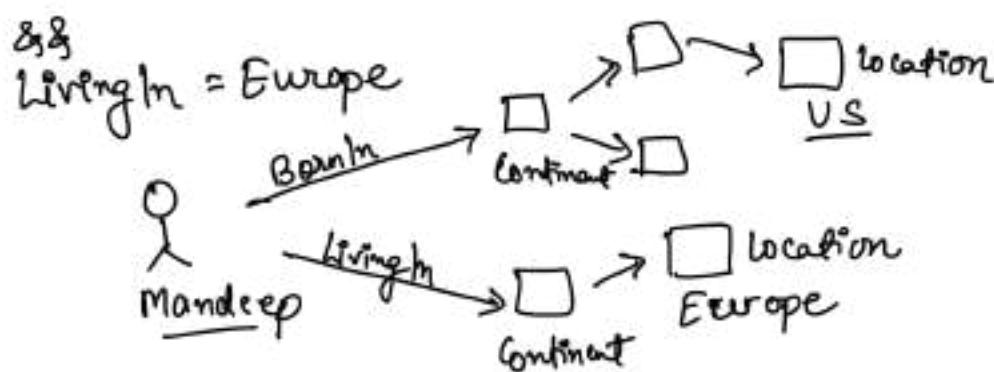
↳ Graph provides flexibility in data modelling as compared to RDBMS.

Example Neo4j

query language - cypher (declarative)

example - find all people who immigrated from US to Europe

↳ BornIn = US



Graph Databases Compared to Network Model

Network	Graph
① schema specifies record type could be nested within each record type.	① No such restriction. A vertex can have edge to any other vertex.
② To reach particular record, traverse the access path.	② You can directly go to any vertex by its unique identifier.
③ children of record are ordered set.	③ vertices & edges are not ordered.
④ Imperative queries	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

Now will data be stored in database?

How data will be retrieved from database?

Data 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)
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.

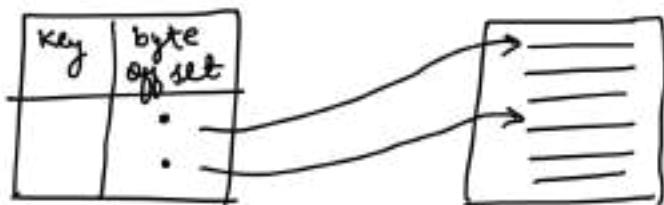
Hash Indexes → How will you index data on disk?

Assumption - Data storage provides only append operation on file

Indexing Strategy →

↳ Create in-memory hash map.

↳ Every key is mapped to byte offset in file.



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

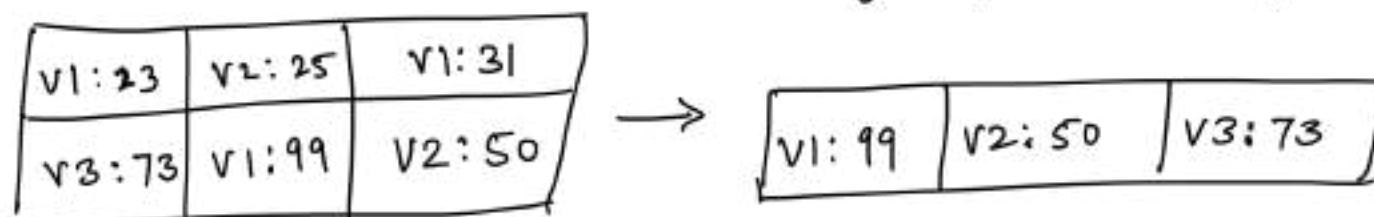
↳ In-memory Hashmap helps if there are limited no of keys.

↳ 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	<table border="1"> <tr> <td>V1: 5</td><td>V2: 3</td><td>V3: 3</td><td>V4: 7</td></tr> </table>	V1: 5	V2: 3	V3: 3	V4: 7			
V1: 5	V2: 3	V3: 3	V4: 7					
S2	<table border="1"> <tr> <td>V5: 9</td><td>V6: 11</td><td>V2: 5</td><td>V3: 9</td></tr> </table>	V5: 9	V6: 11	V2: 5	V3: 9			
V5: 9	V6: 11	V2: 5	V3: 9					
S3	<table border="1"> <tr> <td>V1: 11</td><td>V2: 9</td><td>V7: 23</td><td>V3: 21</td></tr> </table>	V1: 11	V2: 9	V7: 23	V3: 21			
V1: 11	V2: 9	V7: 23	V3: 21					
↓								
<table border="1"> <tr> <td>V1: 11</td><td>V2: 9</td><td>V3: 21</td><td>V4: 7</td><td>V5: 9</td><td>V6: 11</td><td>V7: 23</td></tr> </table>		V1: 11	V2: 9	V3: 21	V4: 7	V5: 9	V6: 11	V7: 23
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.

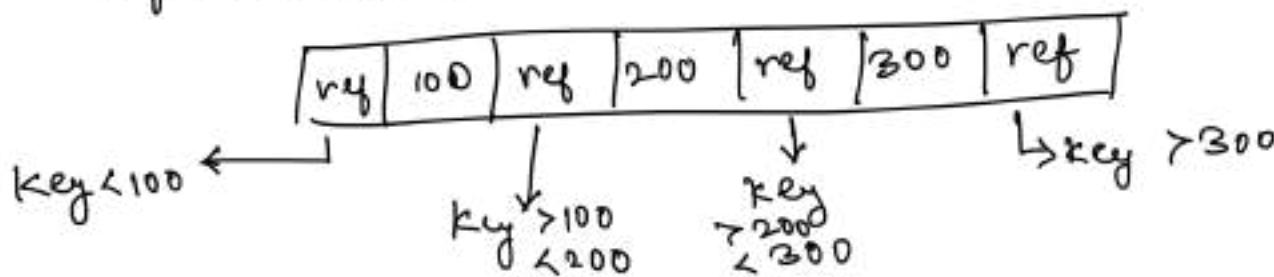
Size-Tiered	Leveled Compaction
HBase: Cassandra uses both.	LevelDB, RocksDB
smaller SSTables are successively merged into larger SSTables.	key range is split up into small SSTable and older data moved to separate levels.

BTree → It is not "Binary" Tree.

↳ key value pairs sorted by key. → on disk

↳ BTree breaks database into fixed size blocks.

↳ Each Block can be identified by address, using which one page refers 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.

LSM 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 to 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.

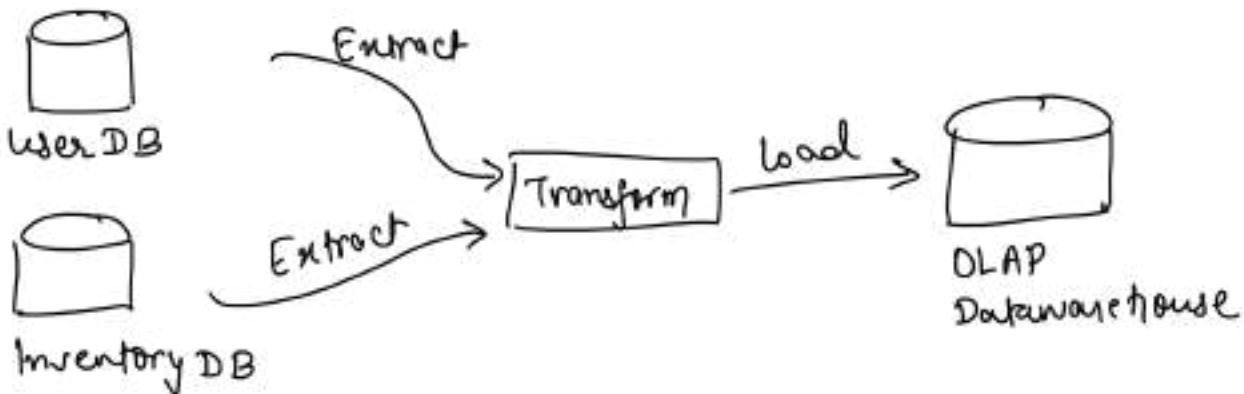
Your use-case → Transaction Processing or Analytics?

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- Ms Deep 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
}
```

↳ req/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

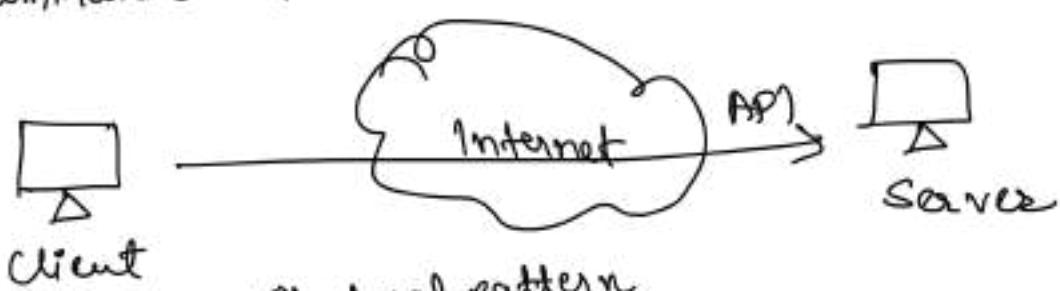
↳ You can specify a default value for newly introduced columns.

↳ 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

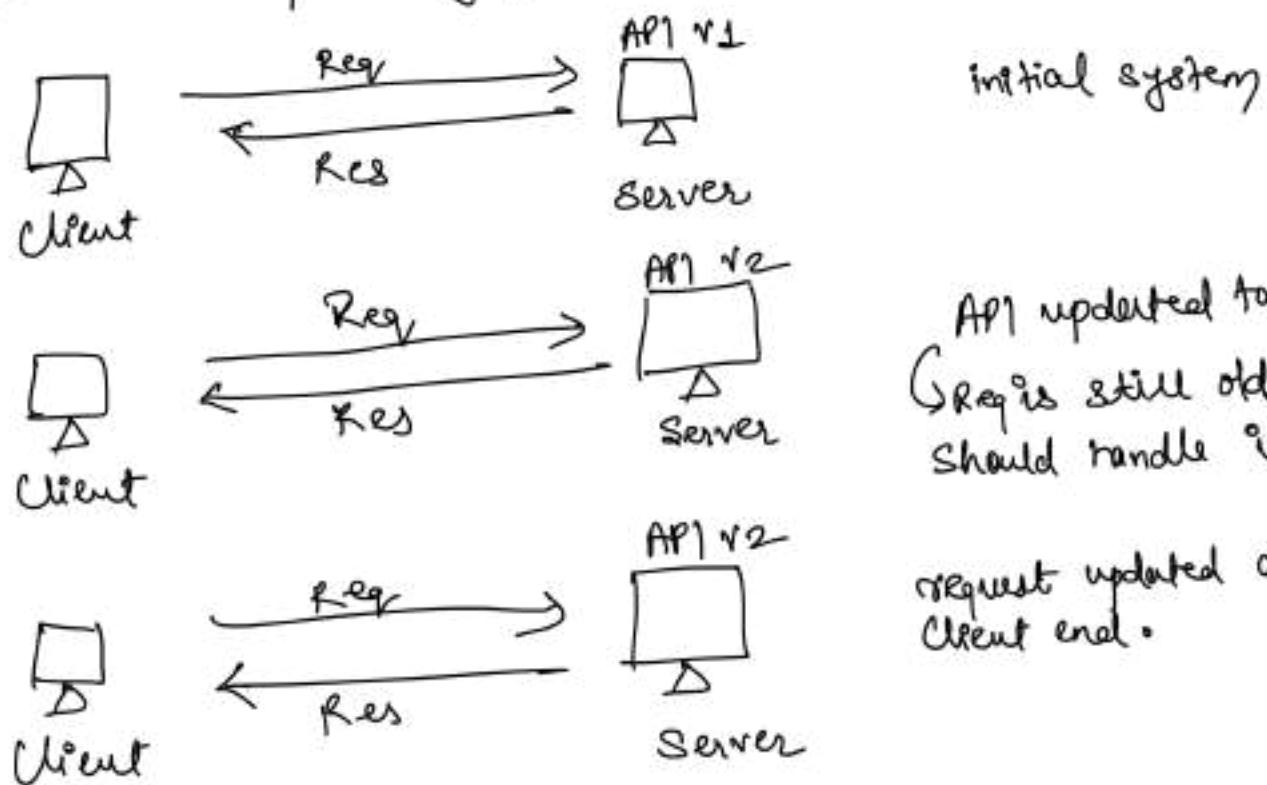
RPC → Remote Procedure Call

↳ make request to remote network service.

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

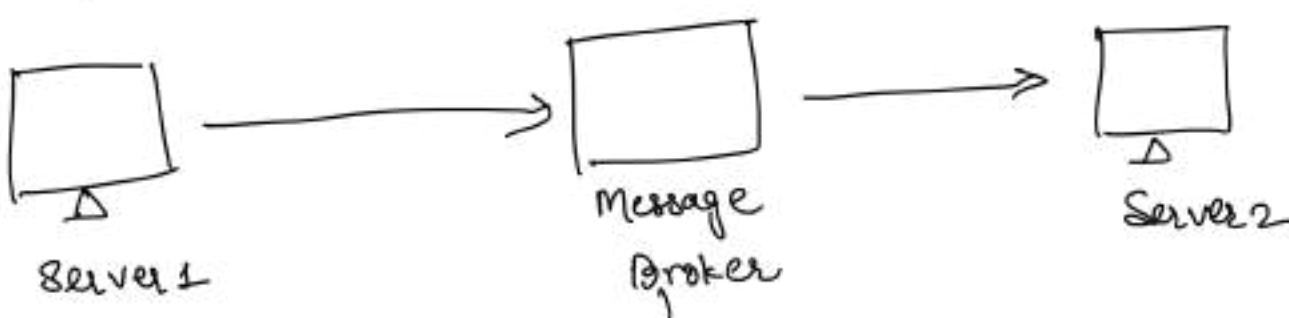
Data evolution → servers will be updated first and then clients

- ↳ 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

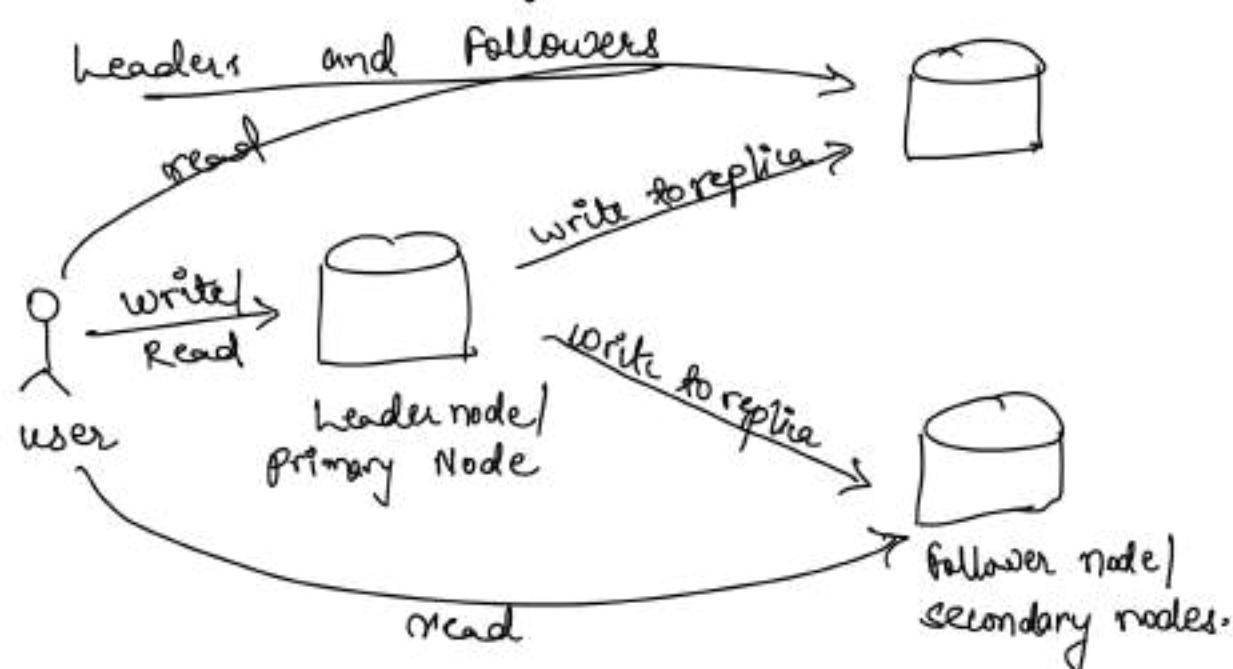
Replication

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

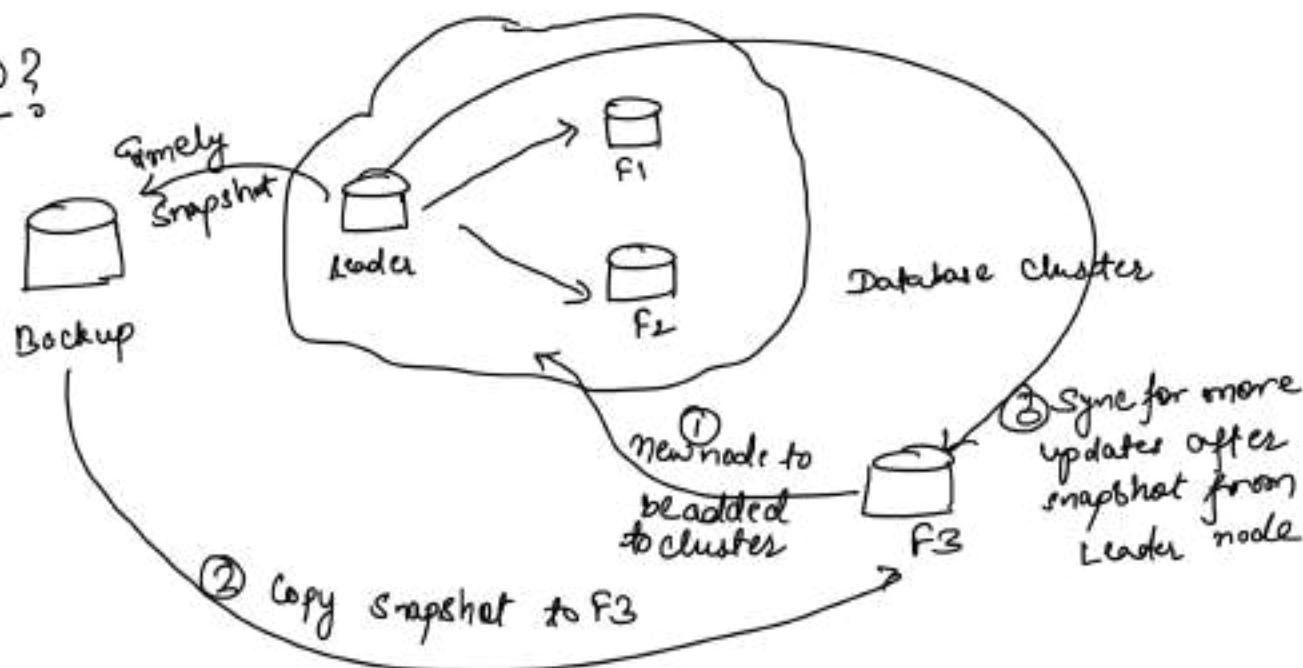
- | | |
|--|--|
| <p>① User is sent ack once write is completed on leader & followers</p> <p>② Data is always consistent</p> | <p>① User is sent ack once write is complete on leader node - follower nodes are updated in background</p> <p>② eventually consistent.</p> |
|--|--|

- 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.

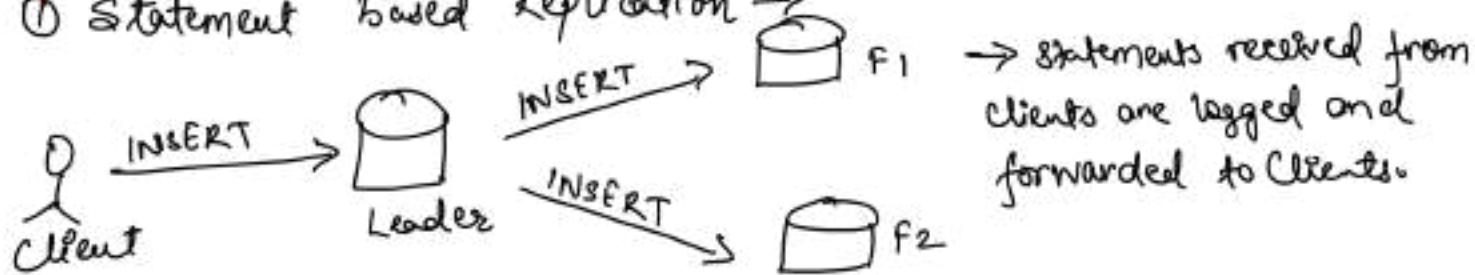
② Leader Failure

- ↳ A follower is promoted as new leader.
- ① Identify leader has failed → health checks
- ② Choose new leader → best candidate would be node which has most recent replicated data.

- ③ 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 (Rowbased) 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 sync 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.

↳ ② 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

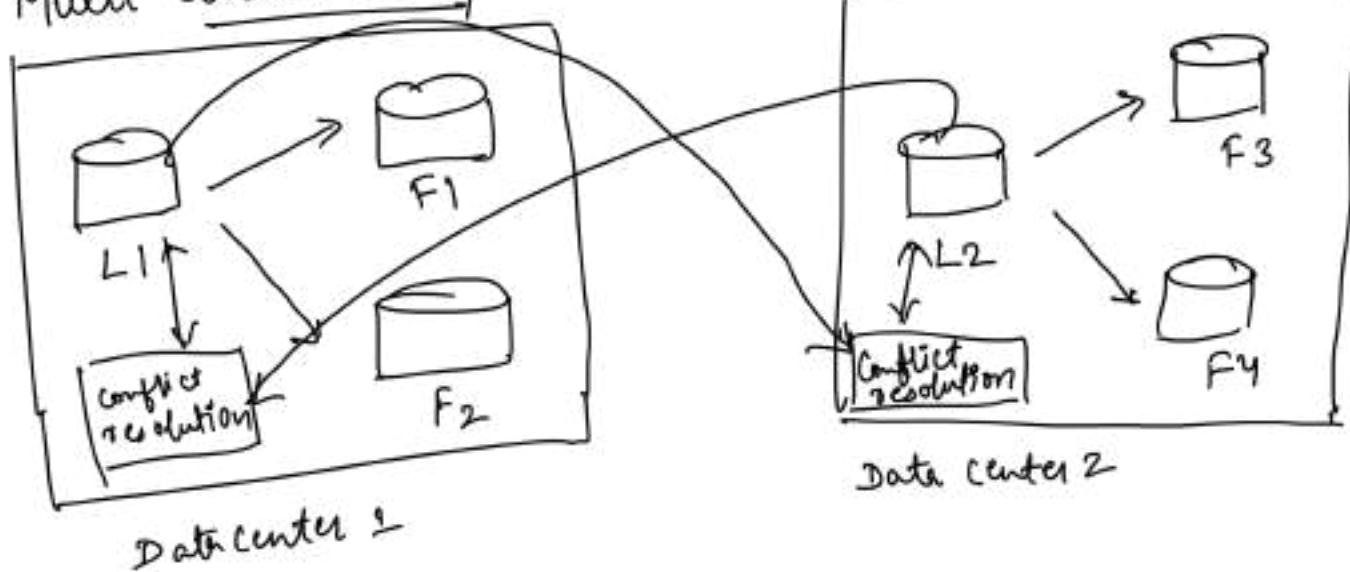
if writes are performed in certain order, reads should happen in same order. → violation of causality.
→ generally happens in partitioned databases.

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.

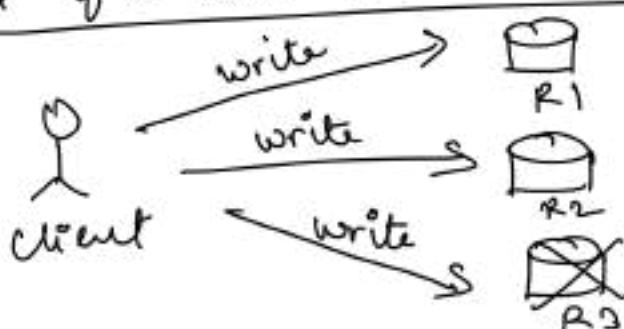
Handling of write conflicts

- 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 will be picked up basis version no.

How will R3 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 ($r \neq w \leq n$) → possibility of stale reads.

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

③ Sloppy quorum →

↳ 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 →

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

⑤ Last write wins → 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 data loss.
- ↳ 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.

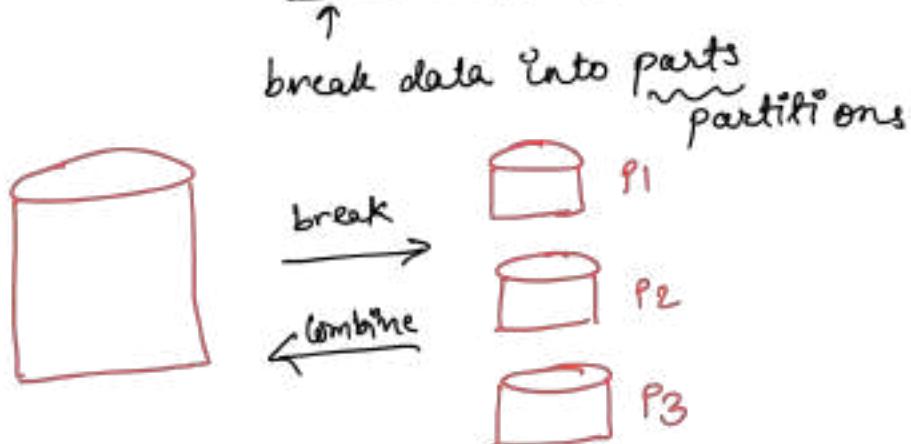
YT channel - Ms Deep Singh

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

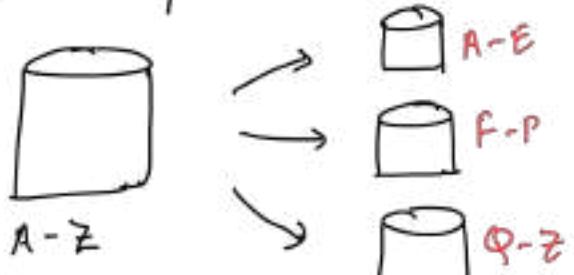
Partitioning of Key-Value Data

↪ 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 at continuous range of keys.



- ④ within each partition, keys are sorted. (sstable & LSM tree)
- ⑤ range queries are easy.
- ⑥ 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.

e.g. Cassandra & MongoDB use MD5.

↳ assign each partition range of hashes.

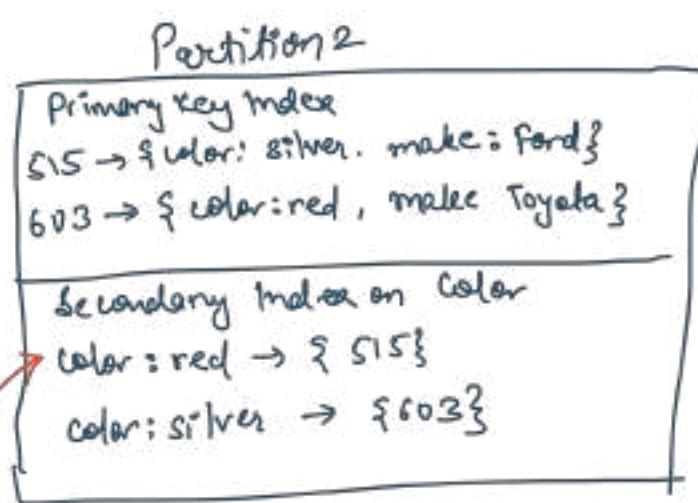
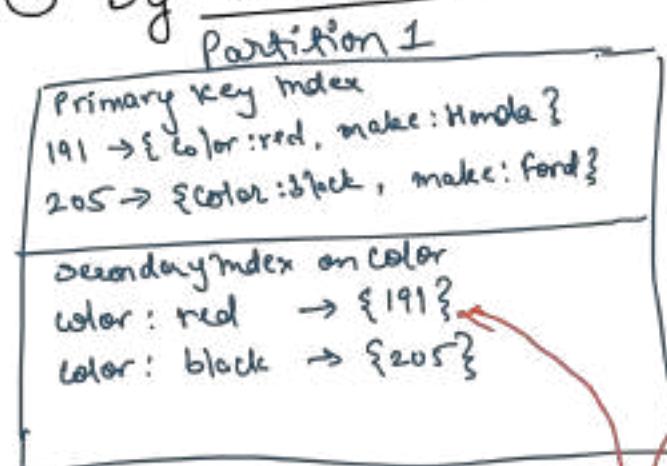
↳ key will be stored in partition whose hash falls in specific range.

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 to 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

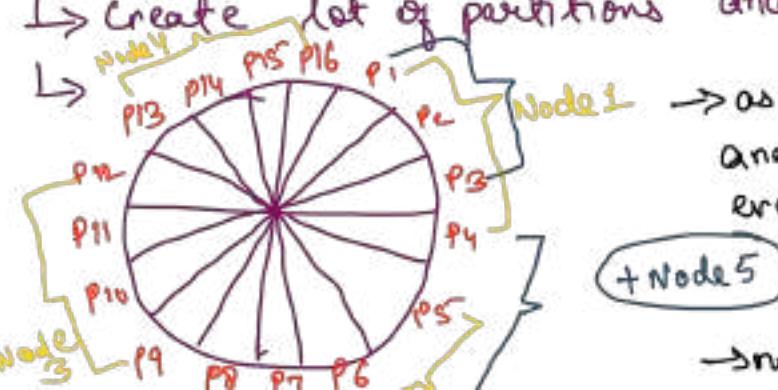
- ↳ add more CPU to handle load.
 - ↳ add more disks.
 - ↳ replace machine.
- ↳ requires movement of data from one node to another.

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.

+ Node 5

→ no of partitions are fixed and only

Nodes

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 & RethinkDB) 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.

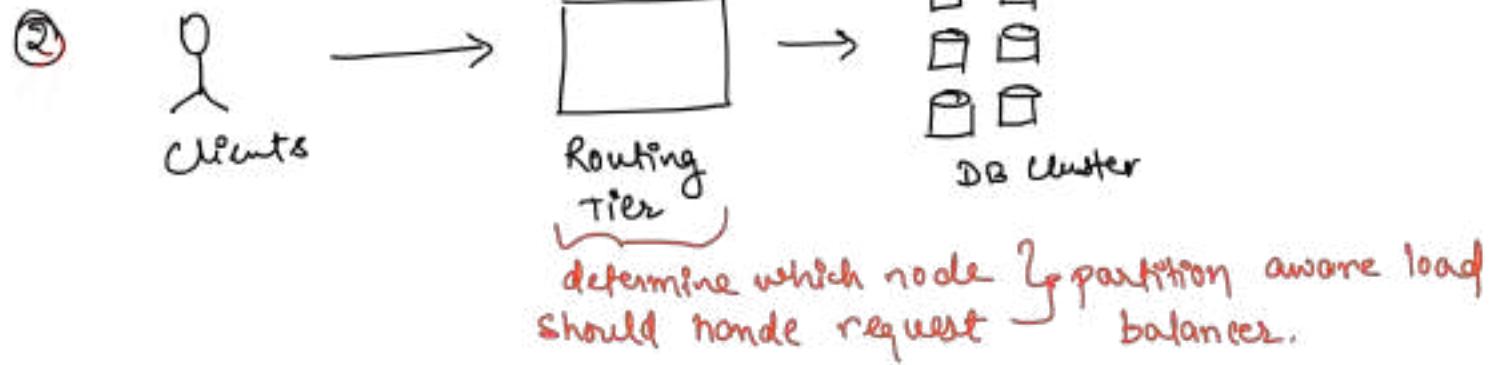
→ Add new node → partitions becomes smaller as data is moved.

ⓐ randomly choose fix no of partitions to split

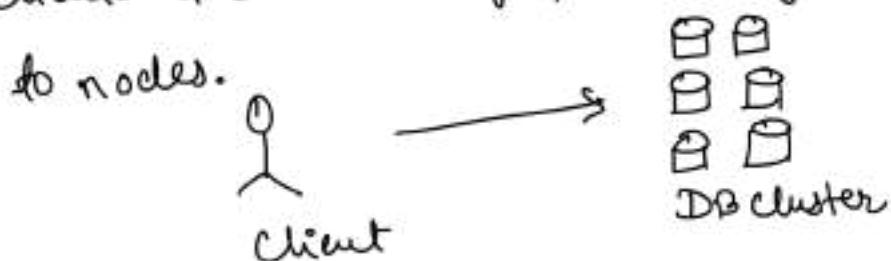
ⓑ After split, take ownership of half of them and leave other half in place.

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

① Route request to any node via round robin. If data not found, move to another node.



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 its 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 - MsDeep 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.

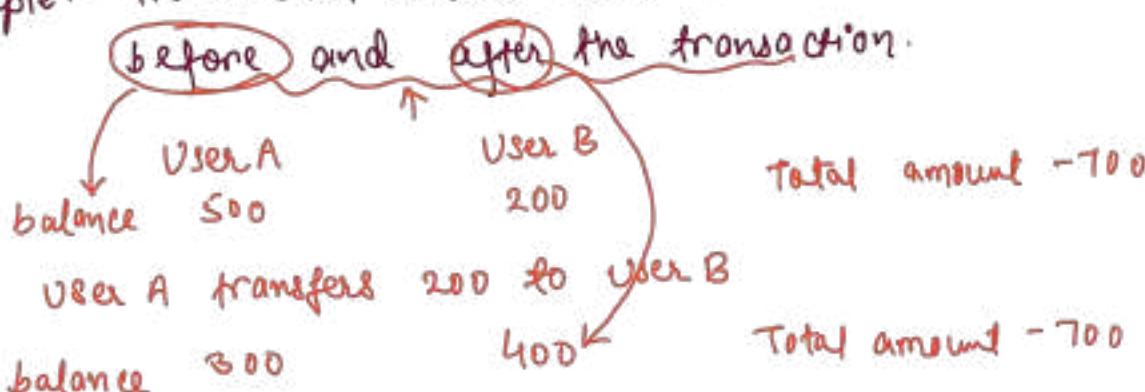
① Atomicity → atomic operation

↳ 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 Tree.

- 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

multi-object transactions

- use-cases
- ① In RDBMS, now 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.
- To avoid this issue, transaction B can wait till A is committed or aborted.
- use locks
- 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

Example	Account A	Account B	
	\$500 ↓ transfer \$100	\$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.
- maintaining 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.

Last Update Problem



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

} two transactions perform concurrently
 ↑ one modification could be lost

Solutions

- ② Atomic write operations → atomic update operations.
 ↑ by acquiring lock

③ Explicit locking → handle locking mechanism in application.



Both users trying to perform action at same piece.

not provided by DB
atomic operations not sufficient
not possible to handle via DB query
ie multiplayer game.

④ Detect lost updates → allow to execute in parallel and fail if lost update is detected.

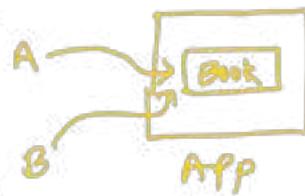
concurrent

⑤ 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 some effect as executed sequentially.
↑ no concurrency

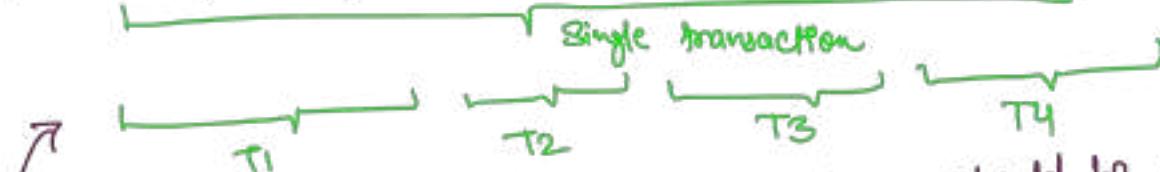
① Actual Serial Execution → more concurrency.

↳ no overhead of locking.

② Encapsulate transactions in stored procedures

↳ transactions should be short.

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



To support serial execution, transactions should be short.
journal Lag due to locks

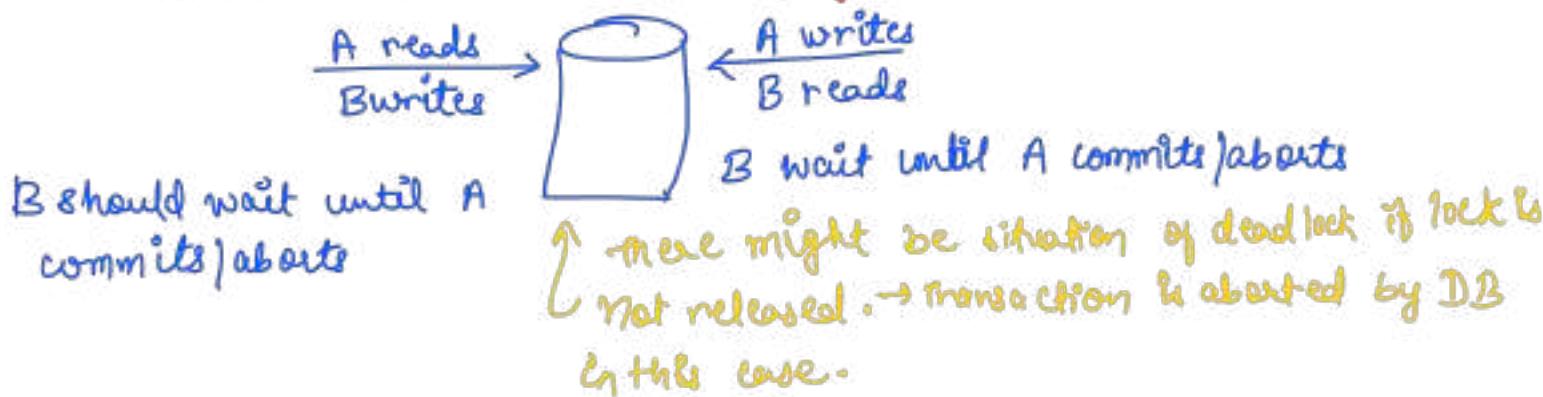
↳ application should submit entire transaction code to DB
ahead of time.

↳ stored procedure.

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.



B should wait until A commits/aborts

↑ 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.
- These have occurred due to any network faults.
- add timeout so as request don't keep on waiting.

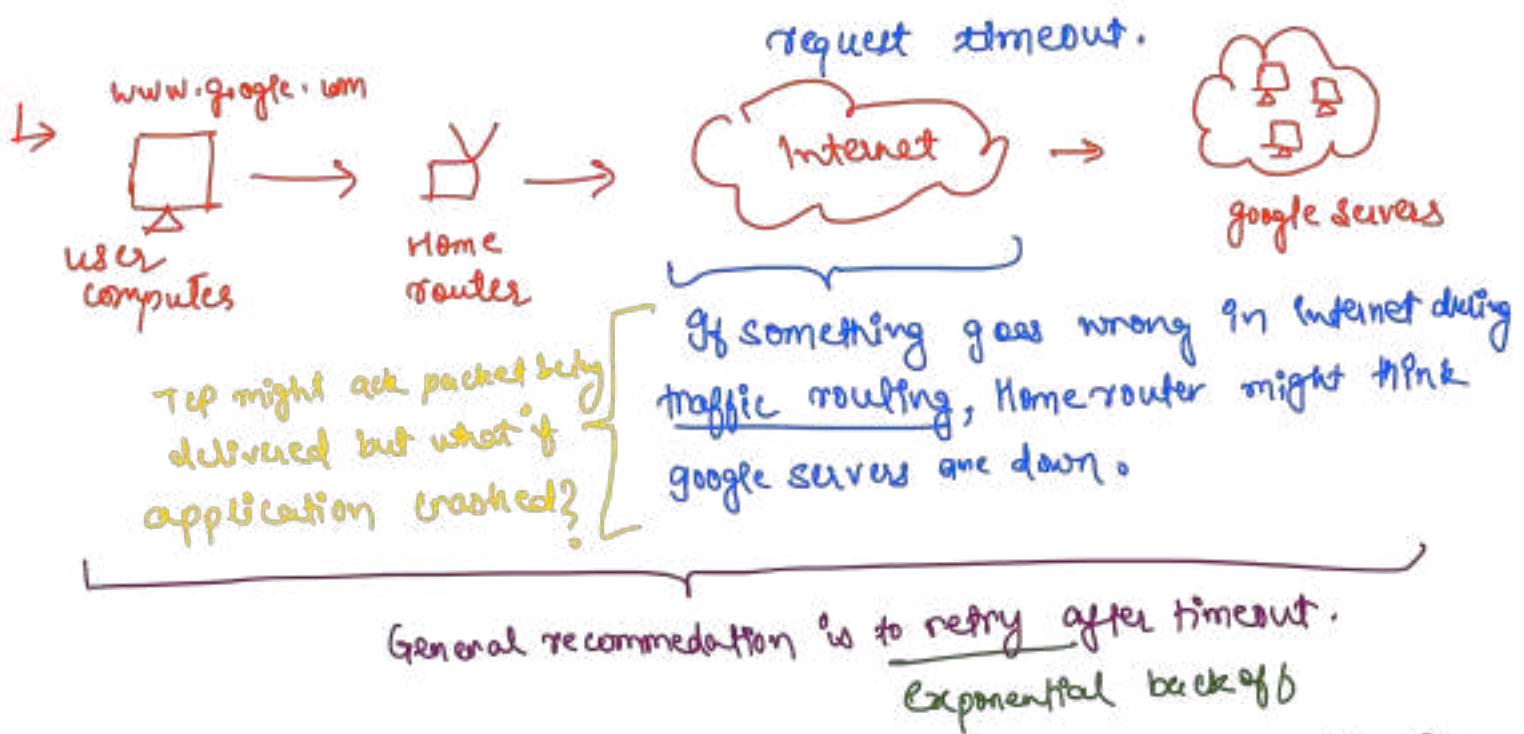
Detecting Network Faults

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?

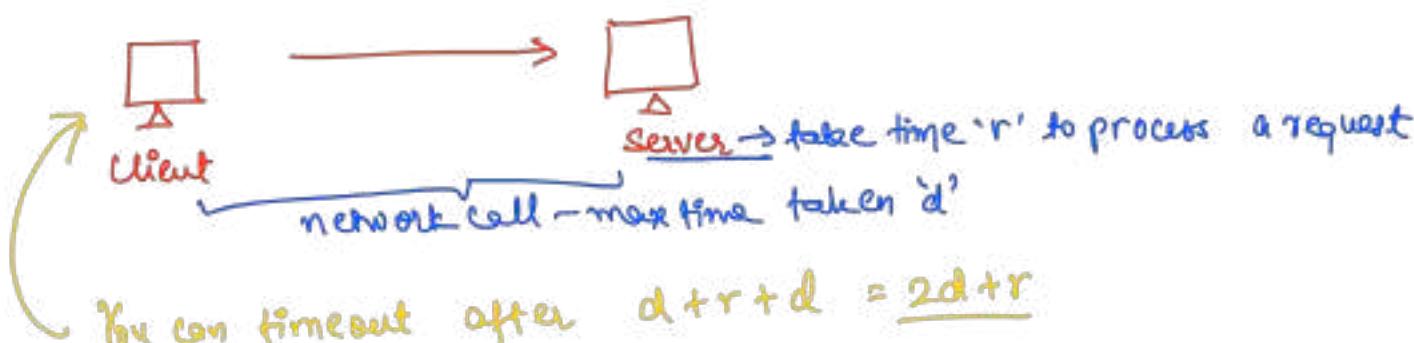
- ↳ e.g. node OS is running but node process crashed.
another node should takeover before



Timeouts and Unbounded Delays

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



Network Congestion and queuing

requests more than threshold has to wait to be served.

until the requests are served, they are kept in queues
queuing can also be at client end
TCP flow control
backpressure

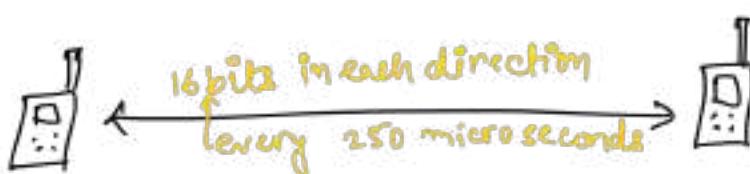
TCP vs UDP

- ↳ perform flow control
- ↳ retry - retransmit lost packets

→ Client performs #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 1 Gbps 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()

between no of second/ms since epoch

wall-clock time

January 1, 1970

midnight UTC

↳ In scenarios if 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.

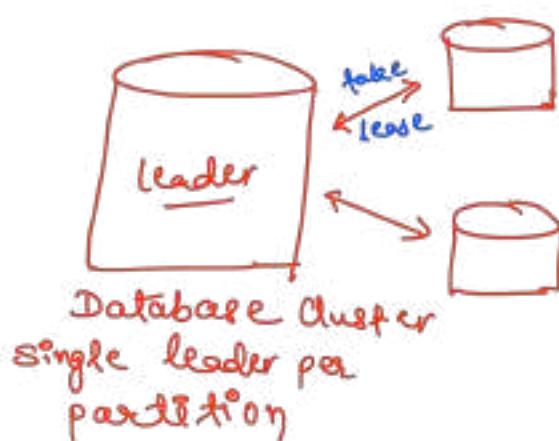
will depend on business your servers are spawned for.

Example - Last write wins algo on DB cluster.

for scenarios where ordering of events is important, logical Clock 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 → because leases

② To solve #1 - use monotonic clocks

what if there is unexpected pause in system

(e.g. example due to Garbage Collector (GC))

(e.g. OS access disk for accessing memory pages (virtual memory))

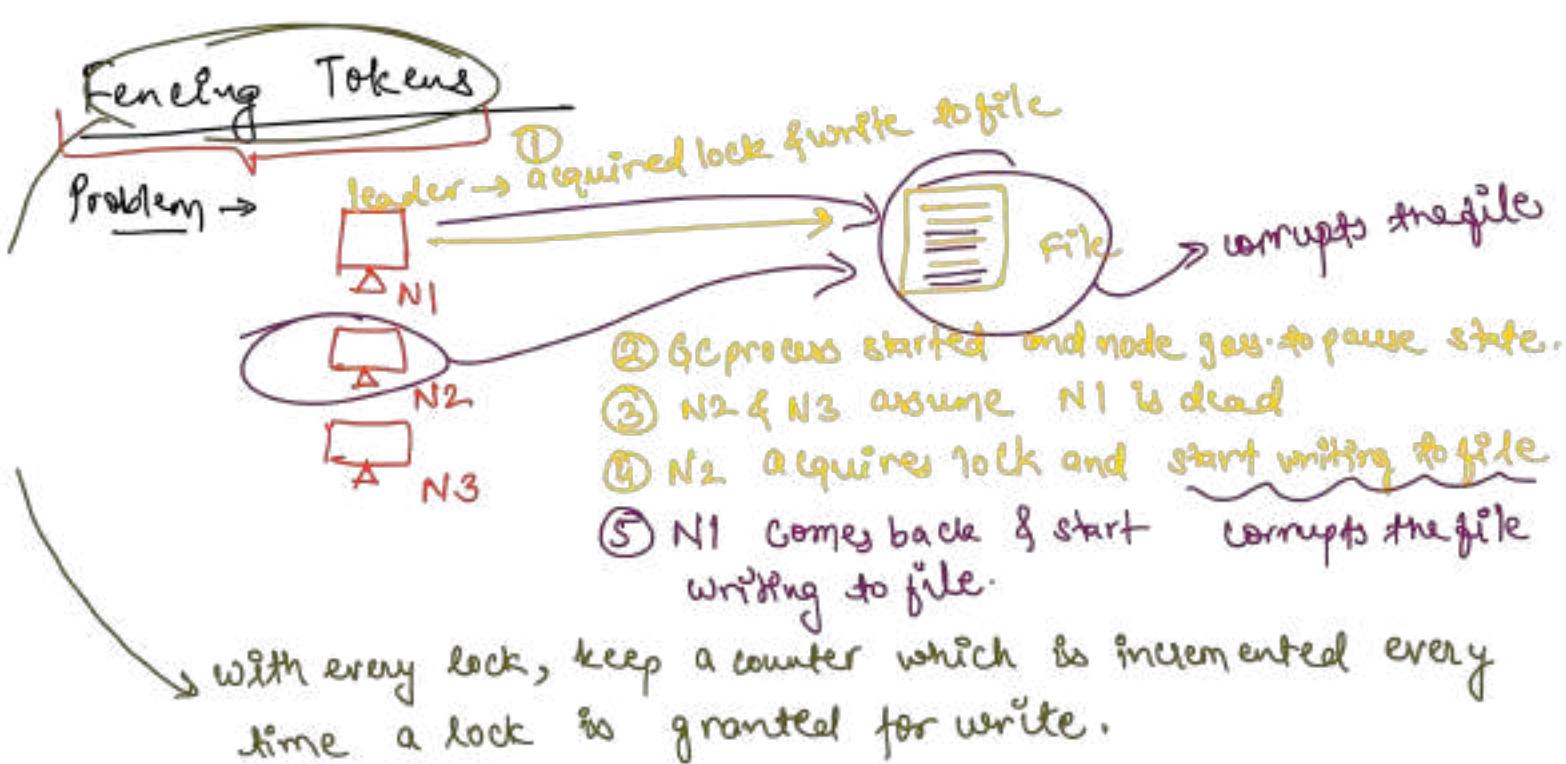
→ The thread is paused as I/O takes place

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 applications



In above scenario

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

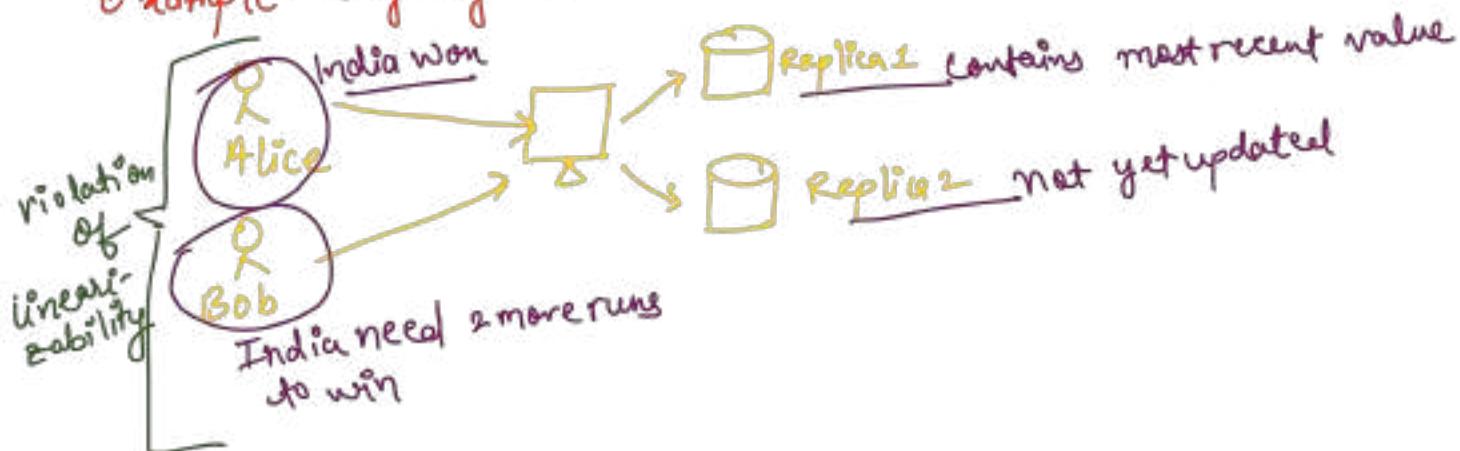
Consistency and consensus

Subscribe for more such content ...

YouTube - Ms Deep 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 all clients should be sent updated score.

Even if write is in progress.

Comparing Linearizability with Serializability

Fairness guarantee on reads

Consistency property of transactions

and writes of an object.

Transaction behaves as same if executed sequentially.

why linearizability?

① Locking and leader election

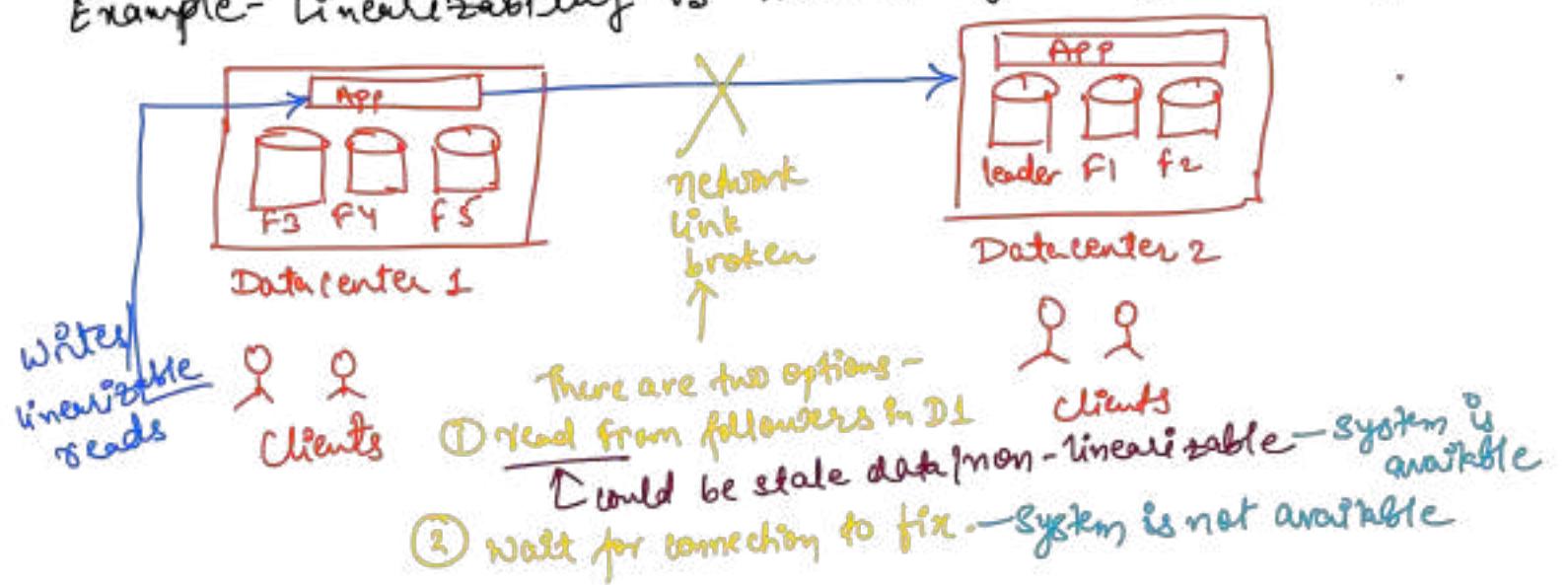
↳ for single leader replication - should only be one leader.
a node can acquire lock on becoming leader. not more → split brain
linearizable. Example Zookeeper, etcd.

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. Impractical, CA system is not possible.

Linearizability and network delays

- ↳ It is generally avoided (if 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

↳ cause

Causal dependency
b/w que & ans

→ question should be asked before so as it can be answered. → effect

Causality impose 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 timestamper can be used.

↳ from logical clock.

assign seq no to each operation → we have 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:

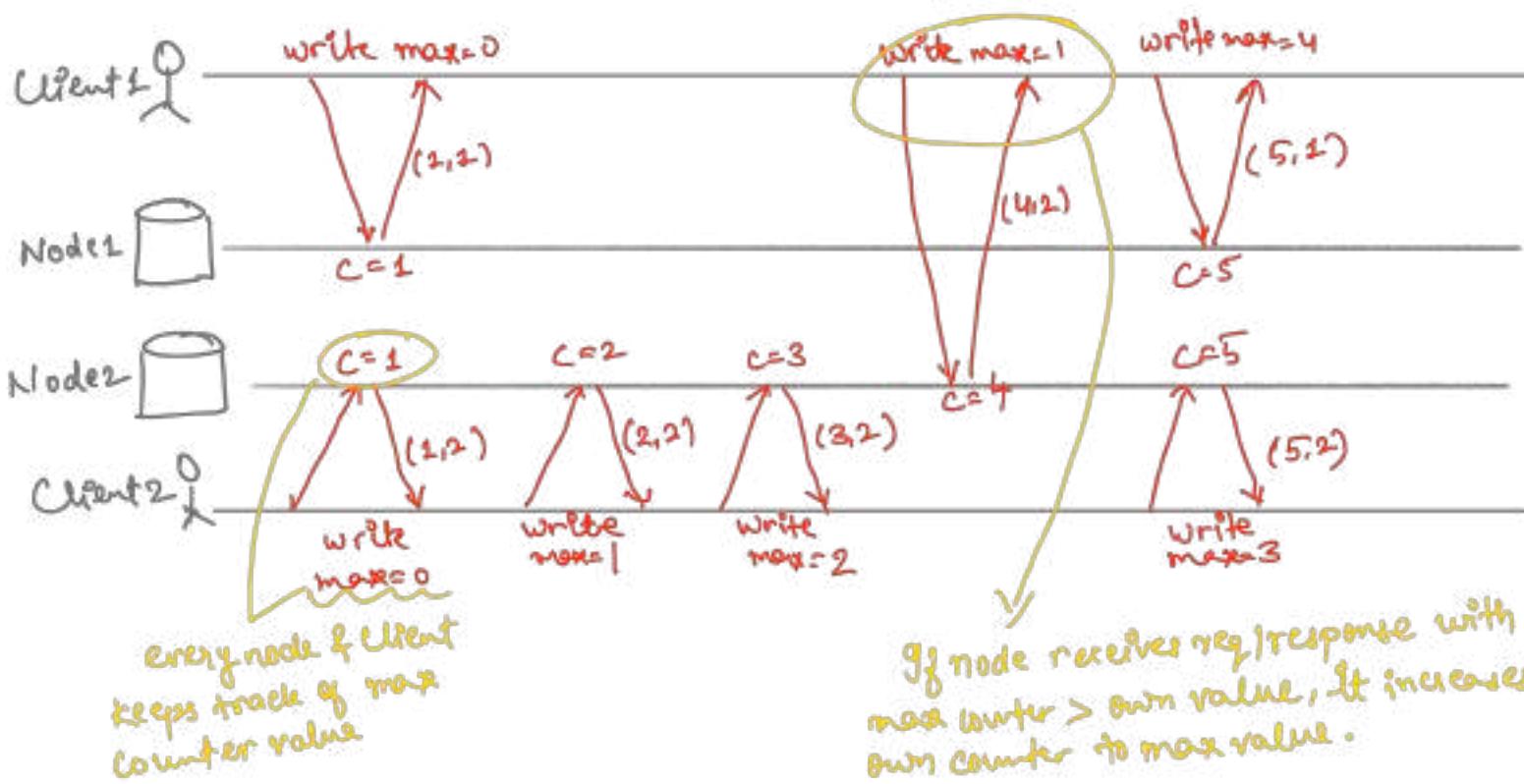
they are not capturing order of events across the nodes.

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?

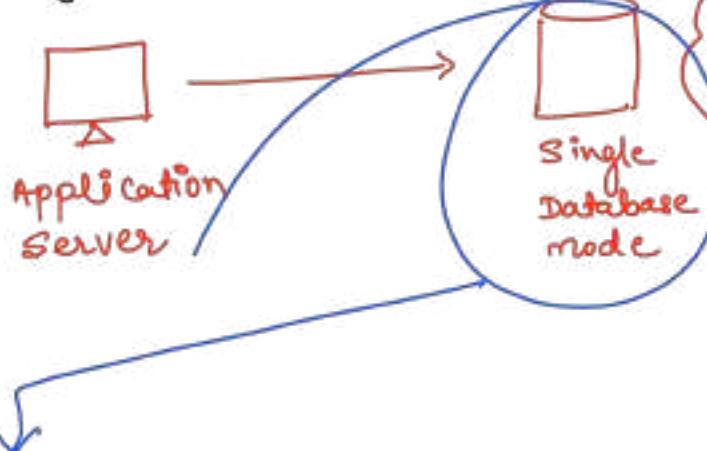
Total order Broadcast → also used in consensus algorithms.

- ↳ 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?

~~not part of 2PC~~

① Either database nodes interact with each to decide what to do with specific transaction.

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

② Coordinator dumps all actions to disk
So it can recover once back online. } what if disk is also corrupt?
} manual intervention by admin.

Distributed Transaction Examples

- ① Database internal → internal transactions among the nodes in cluster.
- ② Heterogeneous → different technologies distributed transactions
ex: 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
- 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

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 single leader replication

and it requires a leader

We need all nodes agree for leader.

↳ consensus

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 e.g. zookeeper, etc.

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

consensus

- 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 . . ~

Chapter 10 Batch Processing

Book Summary - Designing Data-Intensive Applications

Types of Systems

1. Online Systems — [services] — process request as soon as it's received and send a response
2. Offline Systems — [batch processing systems] — takes large amount of data → process → generate output
3. Near-real time systems — [stream processing systems] between online and offline — there might just be a ack response and output sent via events — Chapter 11

Processing data with Unix tools

Example — top 5 popular pages on your website

`cat /var/log/nginx/access.log` = read the log file

Awk '{print \$7}' = split each line into whitespace and output 7th field

Sort = sort alphabetically

Uniq -c = removes duplicate , -c adds a counter

Sort -r -n = sort by number of times page accessed

Head -n 5 = output first 5

Sorting vs in-memory aggregation

1. If dataset has lot of duplicates, it can fit into in-memory hash table
2. Else we can use something similar to SSTables and LSM Trees - combination of in-memory and disk.
3. Unix utility 'sort' automatically spills data to disk and parallelise sorting across multiple cpu cores

Unix Philosophy

1. Make each program do one thing well. To do a new job, build afresh rather than complicate old programs by adding new “features”.
2. Expect the output of every program to become the input to another, as yet unknown, program. Don’t clutter output with extraneous information.
Avoid stringently columnar or binary input formats. Don’t insist on interactive input.
3. Design and build software, even operating systems, to be tried early, ideally within weeks. Don’t hesitate to throw away the clumsy parts and rebuild them.
4. Use tools in preference to unskilled help to lighten a programming task, even if you have to detour to build the tools and expect to throw some of them out after you’ve finished using them.

Limitation of unix tools — they run on single machine.

MapReduce and Distributed FileSystems

MapReduce job —

- * takes one or more inputs and produces one or more outputs
- * Read and write files on distributed file system — HDFS

HDFS —

- * based on shared-nothing principle (No special hardware is needed, only computers connected by datacenter network)
- * A daemon runs on each node, exposing network service for other node to access files on that node.
- * NameNode — keep track of which file blocks are stored on which node
- * For failures — replication to multiple nodes

MapReduce Job Execution —

Mapper

1. Take the input record
2. Extract key and value
3. Write output to HDFS

MapReduce framework

1. takes key-value pair from mapper
2. Collect values belonging to same key
3. Write output to HDFS

Reducer

1. Take collection of values as input
2. Produce output . e.g. count of characters

General Details —

1. Number of mappers = number of input files.
2. Number of reducers is configured by person executing the job
3. Framework use hash of the key to Ensure all same keys end up on same reducer
4. Sorting is done in stages (because large data size) similar to SSTables/LSM Trees
5. Shuffle —
 1. Maps writes output to file.

2. Scheduler informs reducer work is done.
3. Reducers connect to mappers to download sorted key-value pair for their partition.

MapReduce Workflows

Multiple mapReduce jobs chained together in workflows. Output of one job is used as input to another job.

Example workflow tools — Oozie, Airflow

Reduce Side-Joins and Grouping

MapReduce job is similar to full table scan — read entire dataset. Example — two datasets

User 101 clicked button
User 105 accessed page
User 101 load page
User 102 clicked on hyperlink

Log Entries - user activity

101	Msdeep14
102	Mandeep Singh
105	MsDeepSingh

User Info Table

For collecting data — infer who is user from log entries

1. Access log entry one-by-one and call database —> not efficient
2. Local Cache — Take a dump of database to HDFS where Map-Reduce is running —> Better Approach

Sort Merge Joins

MapReduce Execution

Mapper1 - extract user id (key), activity details (value) from user activity log entries

Mapper2 - extract user id (key), user details (value) from user database

Output —

- * the MapReduce framework partitions mapper output by key and then sort key-value pair —> all records (user activity & user database) with same user id become adjacent
- * This is used as reducer input
- * Reducer can process all records for a user id in one go — keeps only one user record in memory at a time & no network calls.
- * Called sort-merge join [Mapper output sorted by key; reducers merge the sorted lists of records from both sides of join]
- * The join work is done by Reducers; also referred to as **Reduce-Side Joins**.

The related data is at the same place — helpful in achieving GROUP BY kind of operations.
e.g. counting page views

What about celebrity user?

Bringing all data at single place for single user will break (creates skew/hot spots) — these records called lynchpin objects —

- * one reducer doing lot of work as compared to others
- * Can take lot of time to MapReduce task completion
- * So instead, mappers send the hot key to random reducers based on some hash — parallelisation

Map Side Joins

Mapper reads file from HDFS and writes back to HDFS — no reducers or sorting.

1. Broadcast hash joins
 1. Large dataset joining with smaller dataset (can fit into memory of each mapper)
 2. Essentially, smaller dataset is broadcasted to all mappers
 3. Pig - replicated join; Hive - MapJoin
2. Partitioned Hash Join
 1. Inputs to map-side join are partitioned in same way
 2. e.g. userId%10 — so 10 partitions
 3. Hive - bucketed map joins
3. Map-side merge joins
 1. Partitioned in same way and sorted based on same key

Output of Batch Workflows

Result of data processing?

1. Building Search Indexes
2. Key-value stores as batch process output
 1. Spam filters, anomaly detection, recommendation systems
 2. The output is saved such as it can be queried by some key. e.g. suggested connections for a user on LinkedIn

Comparing Hadoop to Distributed Databases

1. Diversity of storage
 1. Databases are more structured, Hadoop work on files- can be written using any data model or encoding
 2. Dump anything to files and later figure out how to process it - Don't have to worry about schema design
2. Diversity of processing models
 1. MapReduce allows to run your own code on large datasets and don't depend on database software.
 2. Build tools on top of MapReduce as needed. e.g. Hive
3. Designing for frequent faults
 1. If the task is aborted, it can be retried
 2. This enables better resource utilisation — allows tasks with more priority to run prior.

Beyond MapReduce

Writing Map-Reduce jobs from scratch is hard and time-consuming.

Materialisation of Intermediate State

[input] [MapReduce] [output]
 [input] [MapReduce] [output]

Output of one job acts as input to another job. Writing out this intermediate state of files is called materialisation.

Issues —

- * MapReduce job can only start when preceding job is completed — more execution time
- * Intermediate files are replicated on HDFS — overkill for temporary files

To fix these; multiple **Dataflow engines** such as Spark, Tez, Flink — Handle an entire workflow as one job instead of breaking into multiple sub-jobs

- * sorting is done only when needed
- * No unnecessary map tasks
- * Locality optimisations — task consumer and producer on same machine so data transfer over network
- * Can write to local disk or memory instead of HDFS — no replication needed

Fault Tolerance

- * MapReduce writes intermediate states to HDFS, task can be restarted on another machine on failure
- * Spark don't write intermediate state to HDFS; so if machine is lost; data is recomputed

- * For this it keeps track of how data was computed — use Resilient Distributed Dataset (RDD)
- * The recomputed data should be same; if not; then entire job should be re-run
- * If computation is CPU intensive; materialisation of intermediate states is chapter

Materialisation — as compared to MapReduce, write only needed intermediate states and not all.

High Level APIs and Languages

Writing MapReduce jobs is laborious task - so lot of new platforms came up — hive, Pig, etc.

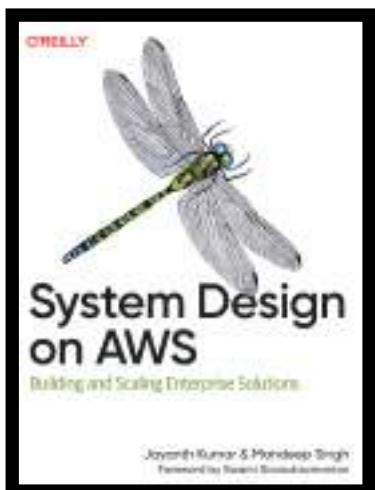
Writing code - requires more expertise and work

Provide declarative interfaces makes it easy — and internal execution engine can also decide the way to execute for faster task completion.

- * Spark, Hive and Flink have cost-based query optimisers — that can change the order of joins and reduce the number of intermediate states.
- * Instead of reading entire record from disk, read only required fields

Subscribe on YouTube — @MsDeepSingh

Read my O'Reilly book "System Design on AWS"



Chapter 11 Stream Processing

Book Summary - Designing Data-Intensive Applications

In batch processing systems, the input is bounded.

But the data could gradually arrive over time — unbounded data.

Stream — data is made available incrementally over time. e.g. stdin and stdout of Unix

Transmitting Event Streams

Input - sequence of records or events

- * Event is a immutable object containing details of something happened at some of time.
- * e.g. order is placed on e-commerce website
- * Event is generated by one producer (or publisher) and can be processed by multiple consumers (or subscribers)
- * Related events are grouped together into topic or stream

Messaging Systems

- * Operates on publish/subscribe model
- * Multiple producer nodes can send messages to a topic and allows multiple consumers nodes to receives messages in a topic.

What happens if producers send messages faster than consumers can send them?

- * Drop messages
- * Buffer messages
- * Apply back pressure — blocking sender from sending more messages

Direct messaging from producers to consumers

- * Without any intermediary nodes
- * UDP multicast in financial industry. e.g. stream stock market feeds
- * Web hooks — consumers expose a endpoint and producers can push messages directly.

Issues —

- * assumption producers and consumers are online; or application logic should be aware to handle faults
- * Consume if offline? It will miss messages

Message Brokers

- * message broker or message queue
- * Runs a server; producers and consumers connect as clients
- * Producers write message to broker and consumers read them from broker
- * Durability is handled by broker

Comparison with Database

- * database keep data until explicitly deleted
- * Data search capabilities such as indexing

Multiple Consumers — multiple consumers read message from same topic

- * Load Balancing — each message delivered to one of the consumers, so consumers can share the work.
- * Fan-out — each message delivered to all consumers.
- * These two patterns can be combined as well

Ack & redelivery

- * consumer crash - messages were delivered by consumer but never processed by it.
- * Consumer should send a ack that message is consumed
- * If broker is unsure of message delivery, it retries on its end
- * If message ordering is important (messages can be reordered in case of retries/redelivery), recommendation is to keep separate consumer per queue.

Partitioned Logs

- * Message brokers are built around transient messaging mindset — delete after delivered to consumers
- * Databases/Files — delete on user request
- * If messages delivered, they can't be recovered — in AQMP/JMS-style messaging

Can we have — durable storage + messaging notification facility???

Using Logs for message storage

- Log — append only sequence of records on disk
- * this approach is used to implement message broker
 - * Producer sends message by appending to end of log
 - * Consumer receives it by reading log sequentially
 - * Similar to unix tail -f

For more scale —

- * the log can be partitioned across multiple machines
- * Topic can be defined as group of partitions to carry messages of same type

- * Each partition, broker assigns monotonically increasing number, (offset).
- * Messages within partition are totally ordered
- * No ordering guarantee across different partitions
- * So , all messages that need to be ordered consistently, route them to single partition

Example — Kafka, Amazon Kinesis Streams

Logs compared to traditional messaging

- * In general, the broker can assign entire partition to nodes in the consumer group
- * Consumer reads messages in partition sequentially, in single-threaded manner.

Downsides —

- * number of nodes sharing the work can be almost the number of log partitions in the topic
- * If single message is slow to process, it also blocks others

If the downsides outweigh over pros of log based message brokers, use JMS/AMQP

Consumer Offsets

- * All messages with offset less than consumer's current offset have been processed
- * Broker doesn't need ack for each msg
- * Just periodically record consumer offsets
- * This helps in increasing throughput

Disk Space Usage

- * Keep on appending to log — will run out of disk space
- * Log is divided into segments
- * Old segments are deleted (or moved to archive store)

Databases and Streams

- * The event can also be any operation on the databases — writing to a database. This event can be captured, stored and processed.
- * Replication log is stream of database write events, produced by leader

Keeping Systems in Sync

Same data can be stored at multiple places in different forms for different use cases. such as OLTP as source of truth, search index for search capabilities, cache for frequent queries, etc.

- * Full database dumps
- * Dual writes — application is responsible to write to multiple data sources
 - * Concurrency issues
 - * One write fails and other succeeds — probable solution, Atomic Commit & Two-Phase Commit

Change Data Capture

Process of observing all data changes in database — dedicated video in System Design playlist

Database —> Parse DB log —> Publish event
e.g. debezium by parsing MySQL binlog

Log Compaction

- * Let's say the broker only keeps data for last 14 days. So if a new data source is added built on top of primary database, we need to build new data source on top of some snapshot
- * Log Compaction — storage engine periodically looks for log records with same key; removes duplicates; and keep most recent update of key.
- * For CDC system, if every change as primary key; and every new update replaces previous value of key; then it's sufficient to keep most recent write.
- * So you can rebuild new data source; by reading log-compacted topic from offset zero without taking snapshot.

Event Sourcing

- * Stores all changes to application state as a log of change events (similar to CDC)
- * But the idea is applied at different level of abstraction
 - * CDC — database is used in mutable way; updating or deleting records at will
 - * Database don't have to be aware about CDC, it is independent operation
 - * Event sourcing, application logic is built on basis of immutable events
 - * Event store is append-only
 - * Updates/deletes are discouraged or prohibited

Deriving current state from event log

CDC — the entry in database represents current state

Event Sourcing — You need to look all the event history associated with specific entry and then compute current state

- * This can be done based on some application logic
- * Recomputation again and again can take time
- * You can take snapshot or cache is somehow faster response time.
- * e.g. series of debit and credit events for your bank account to compute available balance.

Issues —

- * if there are so many updates, computing current state is challenging task
- * Sometimes data deletion is a requirement so immutability don't serve the purpose. e.g. deletion of data for GDPR compliance

Commands and Events

- * Request arrives to the system, it is a command
- * System evaluates the request
- * If the system decides to process the request, it becomes an event — durable and immutable
- * E.g. booking a ticket - command
 - * Check if seat is available
 - * If yes, generate an event; which can further trigger other systems such as invoicing

Processing Streams

Process the stream once you receive it.

- * write to database, cache, search index, etc
- * Push events to users e.g. push notifications
- * Process one or more input streams to produce one or more output streams; called derived streams
 - * The program handling this transformation — operator or a job

Uses of Stream Processing

- * Fraud detection system
- * Trading systems — price changes in financial market

Complex Event Processing

- * short — CEP
- * Describe pattern of events that should be detected based on some regex or filtering logic
- * If match found, system emits a complex event — with details of event patterns that was detected
- * e.g. SQLstream, Samza

Stream Analytics

Aggregations over large number of events

e.g. comparing current statistics to previous time intervals ; such as TPS or latency over 5 minutes interval

Apache Storm, Spark Streaming, Flink, Kafka streams, etc.

Maintaining Materialised views

- * Deriving the alternative view onto some dataset that you can query it efficiently
- * e.g. search index built on top of database

Reasoning About Time

- * Use time windows; such as average over last five minutes

Event Time vs Processing Time

Processing messages may delay —

- * Queueing
- * Network faults
- * Performance issue in message broker; etc

Confusing event time and processing time can lead to bad data. E.g. consumer went down for sometime, once it comes back, it processes the backlog of messages

- * Measuring rate based on processing time — there is spike in requests
- * Measuring rate based on event time - normal/steady

Which clock?

Assigning timestamps to events is difficult task if events are buffered at several points in the system

For capturing user events based on device clock, the user might have deliberately setup wrong time on device; so we can't trust it.

Solution? Use server time when event was received.

But what if the the server receives event after some time (minutes/hours/days) because the device was offline.

Not a full proof solution.

Solution — log three timestamps

- * #1 Event occurring time as per device clock
- * #2 Event sent time to server as per device clock
- * #3 Event received by server; as per server clock

Offset = #3 - #2. [assuming negligible network latency]

Event occurrence = #1 + Offset

Types of Windows

- * Tumbling window
 - * Fixed length, every event belong to exactly one window

- * Hopping window
 - * Fixed length
 - * Allows windows to overlap in order provide some smoothing.
- * Sliding window
 - * Contains all events that occur within some interval of each other.
- * Session window
 - * No fixed duration
 - * Defined by grouping together all events of same user that occurred closely together
 - * Window ends; if user is inactive for sometime.

Stream Joins

Join datasets by key

Stream-Stream join (window join)

1. Detect recent trends in searched-for URLs
2. System store query and result returned for every search request.
3. Further, if something for search result is clicked, another event is logged.
4. Join — combine click and search query event ; key = session id

Implementation —

1. Stream processor maintains state. e.g all events in last one hour; indexed by session id
2. For new search/click event, added to the index
3. Stream processor checks other index for events with same session id
4. For matching event; system emit an event saying search result was clicked.

Stream-table join (stream enrichment)

User activity events

Enriching the user activity events from user profile database

1. You can query db for every new user activity event
2. Or keep a copy of database locally to stream processor. — avoid network call

For local copy;

The database can keep on changing as stream is long running
So the local copy can subscribe to database changes — CDC

Table-table join (materialised view maintenance)

For twitter user's home timeline — expensive to iterate over all people user is following ; find their tweets and merge

So build timeline cache; per user-inbox

- * user u publishes new tweet; added to timeline of every user who is following u
- * For tweet deletion, it is removed from all users' inbox
- * User u1 starts following u2; recent tweets by u2 are added to u1's inbox
- * Reverse for unfollow action

Implementation —

1. Stream of events of tweets, follow relationships
2. Database with set of followers for each user

Basically joining two tables — tweets and follows

Fault Tolerance

How stream processors can tolerate fault?

Microbatching and checkpointing

- * Break stream into small blocks
- * Block like miniature batch process; say one second
- * Called microbatching
- * Used in Spark Streaming

Another approach

- * periodically generate rolling checkpoints of state; and write to durable storage
- * Used in Apache Flink
- * For any crashes, it can restart from recent checkpoint

Atomic commit revisited

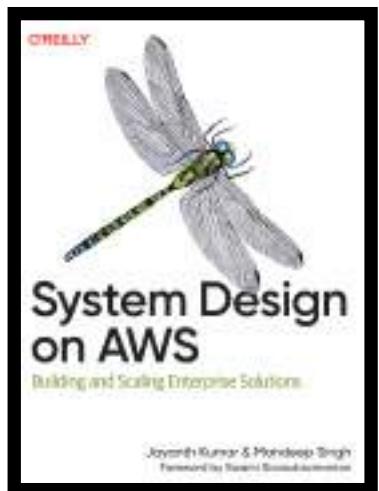
Outputs and side-effects of processing an event take effect if and only if the processing is successful

Idempotence

Perform operation multiple times; it has same effect as it is processed once.

Subscribe on YouTube — @MsDeepSingh

Read my O'Reilly book "System Design on AWS"



Chapter 12 The Future of Data Systems

Book Summary - Designing Data-Intensive Applications

Ideas and approaches that can improve the ways we design and build applications.

Data Integration

- * Every solution for a problem has pros, cons and trade-offs.
- * Appropriate choice of tool depends on circumstances.

Combining Specialised Tools by Derived Data

- * A single tool might not solve all the problems so might need combination of multiple. E.g. different datastore for different needs - relational, search indexes, caching, analytics....
- * What one person considers pointless features might be very much important feature for someone else.

Reasoning about Dataflows

- * same data is stored in different stores for different access patterns
- * You need to be clear about input/output; what's source and destination for the data.

E.g. Database -> CDC -> Search Index

If there is single source of truth for the derived data, it's easier to reconstruct in case of faults.

Derived Data vs Distributed Transactions

Distributed transactions - can be helpful to keep data consistent — generally in sync or atomic operation

Derived data - generally written in asynchronous way

Combining Specialised Tools by Derived Data

- * Integration of OLTP db + full-text search index for arbitrary keyword queries
- * For simple applications, PostgreSQL full-text indexing feature can be used.
- * Search indexes are generally not used as durable system of record.
- * One person might consider something as pointless feature but it could be central requirement for someone else.

The limits of total ordering

- If system is small, construction of totally ordered event log is relatively easy. Limitations with scale —
- * all events should pass via single node which decides ordering. You might need to partition if single machine can't handle the volume.
 - * If servers are multi-region; network delays can make synchronous coordination inefficient — can lead to undefined event ordering.
 - * Microservice architecture — two events generated in separate services

Consensus algorithms/total order broadcast work well with single leader node; it becomes difficult in case of multi-nodes

Ordering events to capture causality

- * If there is no causal link between events, lack of total order is not an issue.
- * Facebook; two users in relationship and broke up
 - * User A removed user B from connections and posted something online
 - * User B should not see this post.
- * If there are separate systems for messaging and friendship; the ordering of events can be lost.
- * e.g. message event is processed first and unfriend event later.
- * So the notifications can be join between messages + friend list.
 - * Logical timestamps can be helpful. The recipients should still have handling for out of order events.
 - * Log an event to record state of system and assign unique identifier. This identifier can be used to record causal dependency.
 - * Conflict resolution algorithms.

Batch and Stream Processing

- * The output of batch/stream processing are derived datasets

Lambda Architecture

1. Record incoming data — append immutable events to dataset
2. Generate read-optimised views from the events
3. This architecture paradigm purposes two systems in parallel
 1. Batch processing systems such as MapReduce (consumes the events and produce corrected version of derived view)
 2. Stream processing system such as Storm (consumes event and produce (approximate) update to the view)

Issues with Lambda Architecture:

- * Operational overhead of two systems
- * Output from separate systems is combined for user requests — can be hard for complex joins

Unifying the Two Systems — need below features

- * Ability to replay historical events
- * Exactly-once semantics for stream processors
- * Tools for windowing by event time

Unbundling Databases

Difference between Unix and databases

- * Unix is thin wrapper on hardware resources - files
- * High level abstraction, user don't have to worry about implementation details.

Creating an Index

- * Take a snapshot of existing data and store them in index in sorted fashion
- * Continue on this with any new data added

Meta-database for Everything

A single database can't all solve all the data needs.

- * Federated Databases : Unifying Reads — single interface to access all underlying data stores.
- * Unbundled Databases : Unifying Writes — synchronising writes across the underlying data stores.

Making Unbundling Work

- * Traditional approach to synchronise writes — distributed transactions
 - * Every system has different transaction semantics, so distributed transactions could be hard.
 - * This might not be a good approach for transactions spreading multiple data stores
- * Practical Approach — async event log with idempotent writes
 - * Loose coupling between components

Unbundled vs Integrated Systems

- * Complexity is high for managing multiple infrastructure (learning curve, operations)
- * Single integrated software can be easy to manage.
- * If feasible, we can try to use single software for multiple needs as long as it serves the use case. <Don't think of too high scale in beginning, you may never reach there>

Designing Applications Around Dataflow

Separation of Application code and State

- * Application code is deployed with tools such as Docker, Kubernetes, etc.
- * Database is used as durable storage
- * Most web applications are stateless — easy to add/remove servers

Dataflow : Interplay between state change and application code

Application code can respond to state changes — by consuming events

Stream Processors and Services

Example — customer is buying something priced in USD but paying in INR; you need currency conversion logic

- * Microservices — query exchange-rate service/or database to get currency rate
- * Dataflow approach — subscribe to stream of exchange rate updates and keep record of exchange rate in local database.

- * dataflow approach is faster
- * Dataflow is more robust, independent of failures of another service

Reads are Events too

Generally write requests are seen as events but you treat reads as well as events

- * read event goes to stream processor
- * Processor emits the results of read to output stream.

Multi-partition Data Processing

- * For queries touching single partition, using streams for reads could be overkill
- * Could be useful for complex queries needing multi-partition data

Aiming For Correctness

Databases

Just because database has safety properties such as serialisable transactions; the issue can still happen if there is bug in application (write incorrect data or delete some records)

Exactly-once execution of an operation

- * Processing twice is a form of data corruption.
 - * Idempotency is one way to overcome this.

Duplicate Suppression

Suppress duplicates

- * TCP use sequence numbers to figure out lost/duplicated packets.
- * Retransmit lost packets and remove duplicates.
- * e.g. database connection - transaction is tied to client connection
 - * Network issue at client after sending COMMIT and before getting response from db server
 - * Client can't know if transaction committed/aborted
 - * Client can retry but it's outside of scope of TCP duplicate suppression (because new connection)

- * Two-phase commit is solution to 1-1 mapping b/w TCP connection & transaction.
- * Is this enough ?
 - * No
 - * What if connection issue between user and application server
 - * No response received by user
 - * So user retry
 - * So duplication can still happen

Uniquely identifying requests

Consider end-to-end flow of request.

- * generate unique identifier such as UUID
- * This id is used by all the components involved in serving the request

Applying end-to-end thinking in data systems

You need end-to-end solution to fix the problem;

One component handling it really well don't solve the problem

Doing The Right Thing

Every system is built for a purpose — and every action has consequences (intended and unintended)

*** Predictive Analysis**

Data analysis to predict weather

*** bias and discrimination**

*** Feedback loops**

* ...

Privacy & Tracking

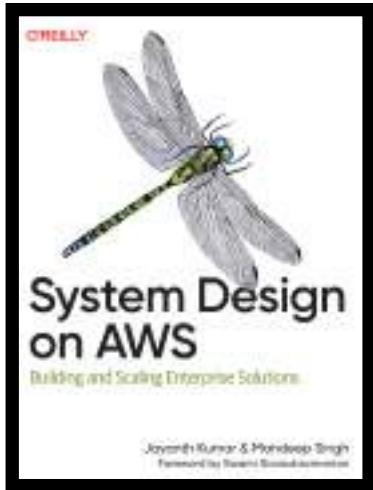
*** Consent**

*** Privacy and use of data**

*** Data as assets and power**

Subscribe on YouTube — @MsDeepSingh

Read my O'Reilly book "System Design on AWS"



📖 Passionate about System Design and want to learn how to build large-scale systems using AWS? "System Design on AWS" is the ultimate resource for you. Get your copy -

O'Reilly: <https://oreil.ly/ruQbc>

Amazon.in: <https://amzn.to/4jExkte>

Shroff Publishers: <https://www.shroffpublishers.com/books/9789355428035/>

Amazon.com: <https://amzn.to/3D3BIBJ>

DragonflyDB (Free chapters): <https://www.dragonflydb.io/content/aws-system-design>