

SHARDING

→ dividing data into smaller fragments + storing them on different distributed machine

TYPES**Range Based / Dynamic**

→ Each shard holds a range of values

→ Ex - date ranges, time ranges

→ when ordering + ranges are there in data

Hash-Based / Algorithmic / key based

→ use a hash function on data to get which shard

→ used when there is no natural order (ids)

Directory Based

→ you manually maintain a data → shard mapping

Geo-Based

→ shard based on geo locations or nearest to locations

Partitioning → same logic but logically within same system

PROS

→ Faster and efficient queries (given you pick good sharding strategy)

CONS

→ Outage possible as it is distributed

SOLN Master-slave design

handles writes

handles reads

→ Fixed number of shards → Need rebalancing / resharding

SOLN Consistent hashing

→ Skewed shards → Better sharding fn approach

→ joins across shards are costly

IMPROVEMENTS**SHARD + INDEXING**

Ex - shard on CITY + index on AGE

→ Here first shard on city + internally each shard is indexed based on age

Query - give me all persons in DELHI with age between 20-40

CONSISTENT HASHING

* Normal hashing - lets say you have N servers

$$\left(\frac{h(\text{data})}{N} \right) \rightarrow \text{slot-no}$$

hash fn no. of servers

you can't change this
if done needs rehashing everything

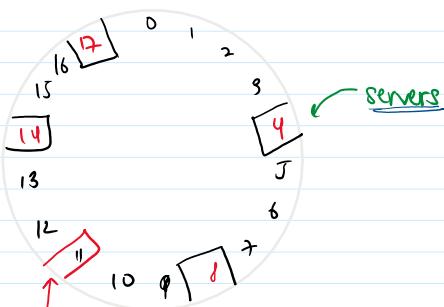
Consistent Hashing

→ Consider your slots are in circular form

→ $x \rightarrow \text{hash.fn}() \rightarrow (0-360^\circ)$ example

→ Each result of hash fn will lie on this circle

→ Some of these slots will be servers also

**HOW IT WORKS**

→ data $\rightarrow \text{hash}()$ \rightarrow slot-no = 3

→ Move clockwise and see next server to 3

→ 3 gets stored in 4 ie 4

New server gets added

→ 11 is now made a server

→ New indexing doesn't change (9,10 will go to 11 instead of 14)

Rebalancing is minimal

→ only 14 which was prev assigned to 14 needs to be rendered to 11 now

* Load factor = $1/N$ (ideal case) equal load to each server

→ But this isn't possible when you have sparse

SOLN Virtual servers

→ Add more server slots which brings load factor to $1/N$

→ Internally a single server handles (2+ slots)

Ex - R. 111 means 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1. 1.

→ Add more server slots which brings load factor to 'N'

→ Internally a single server handles (2⁷ slots)

Ex - ⑪ might be handled by a single physical server

ACID

Maintained in **Relational DBs** → Good at vert scale (Bad at vertical)
+ Fixed schema
+ Generally not distributed

Atomicity

→ A txn can either happen completely or not happen at all

Ex - Bank amt transfer

① Deduct frm ur acc → ② Add to the des' acc

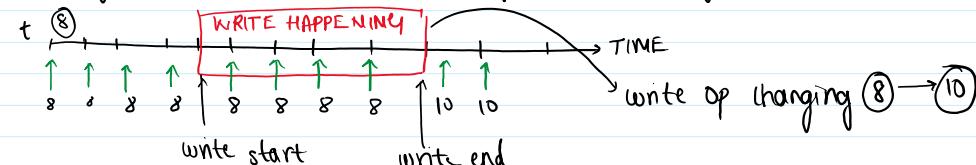
Now this shud happen completely. Ex lets say it got deducted from accnt but then failed

→ In such case revert back everything or nothing approach

Consistency → multiple reads at same time should give same results
why easy? since its not distributed no issue of
out-of-sync replicas

Isolation

→ Every transaction happens independently without caring abt other parallel transactions



→ Here write op was changing val frm ⑧ → ⑩ but during the process any read doesnt know or care abt this and still reads at ⑧

→ Only after write operation is completed and state is updated, Read returns ⑩

Durability Every txn is logged and data is being persisted

BASE

→ shown in **NoSQL databases** → No fixed schema

+ vertically scalable

+ Distributed (data stored across machines)

Basically Available

→ Since data is distributed across servers it can handle outages

→ If one server crashes another replica/machine can give u the data

Soft states

→ States/values can auto change without txn or user queries

How? SYNC across replicas

→ Since data is stored across machines/replicas it may happen that data copies keep getting updated



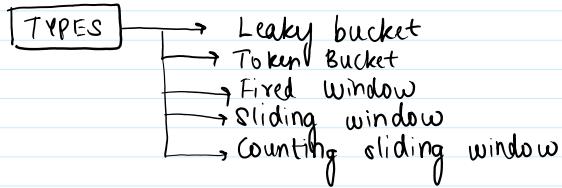
Eventually Consistent

→ Not immediately consistent → ACID ones are immediate consistent

→ 2 simultaneous reads may give different results sometimes

Why? SYNC might have been happening b/w 2 reads

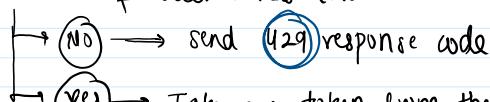
∴ TXNs are not always consistent due to distributed nature
BUT will become consistent after some time (once sync(updates are) completed)



TOKEN BUCKET

→ You have N tokens in a Bucket (Consider this as a currency)

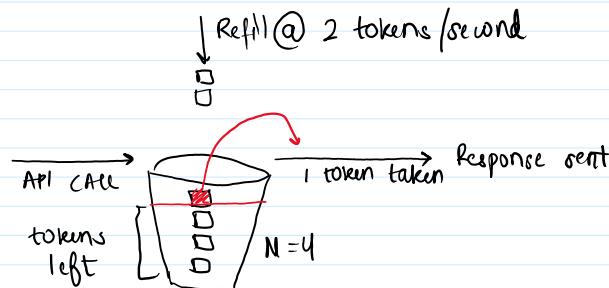
→ whenever a request comes
① Check if bucket has token



→ Refiller: → Refill the tokens in the bucket at certain rate

$$\text{refill_rate} \equiv \text{TPS} / \text{speed of your api}$$

- * Configurations →
 - 1 N = Bucket size
 - 2 Refill rate



LEAKY BUCKET

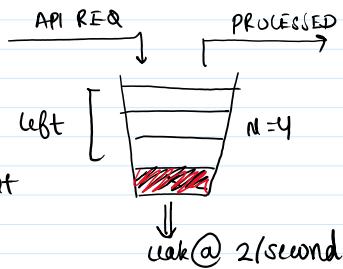
→ You have a bucket of size N

→ A request comes in, if bucket

- Full → 429 code
- else → Add 1 token to the bucket

→ Leak rate → You keep on removing tokens at this given rate
similar to refill rate in token bucket

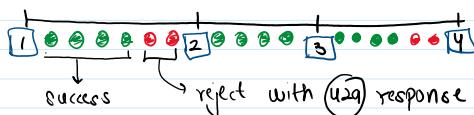
* IDEAL CASE → Rate of API hits \leq LEAK RATE



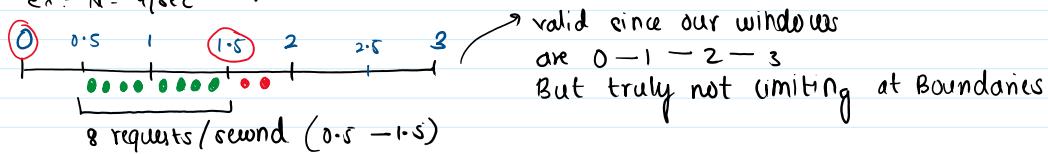
FIXED WINDOW

→ You maintain certain windows (INTERVALS) and each window can only contain N items

→ Ex - Window of 1 min each $\quad N=4$



CON → Doesn't always guarantee rate limit
Ex: $N = 4/\text{sec}$



SLIDING LOG WINDOW

→ Instead of fixed window you use timestamps to dynamically calculate it

→ Rate = 4 req/sec

(Ex)

Request comes at 1.5, you count no of req b/w 0.5 \leftrightarrow 1.5
if num < 4

→ Insert this log (timestamp + reqid)

* Ensures proper dynamic rate limiting

(CONS)

→ you store each request as one entry

∴ Space needed = N

→ lets say your rate is 1000 rps \Rightarrow you will store 1000 log entries

SOLUTION

SLIDING COUNTER WINDOW

→ Instead of pushing each request as log item maintain a counter

→ Rate = 1000 rps

ts	counter
0.1	1
0.2	100
0.3	200
:	:
0.9	400

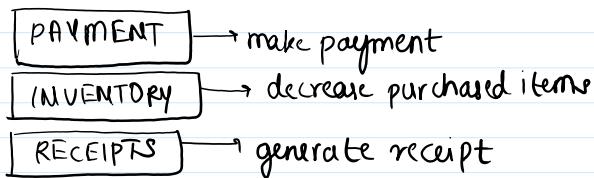
→ Here instead of 1000 logs you will only have 10 logs

→ space = 1000 \rightarrow 10

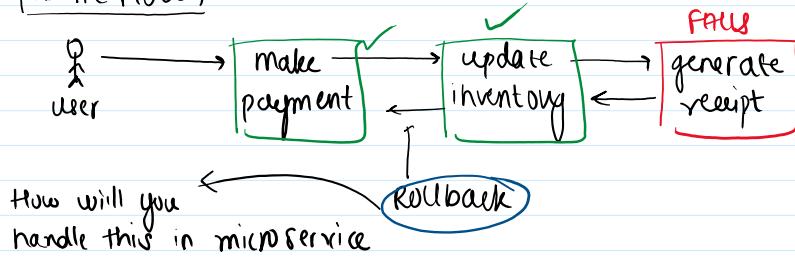
→ No of ts steps

DISTRIBUTED TRANSACTIONS

PREMISE → You have 3 services



IDEAL FLOW



How will you handle this in microservice

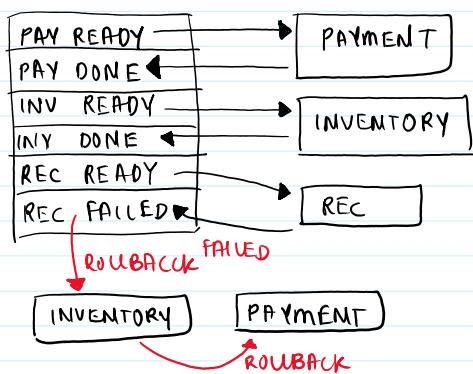
- ① 2PC [2/3 phase commits] BLOCKING
- ② 3PC
- ③ Saga (ASYNC)

SAGA

* Asynchronous and non-blocking but complex to implement

→ Queues are used for Inter Process Communication

STEPS



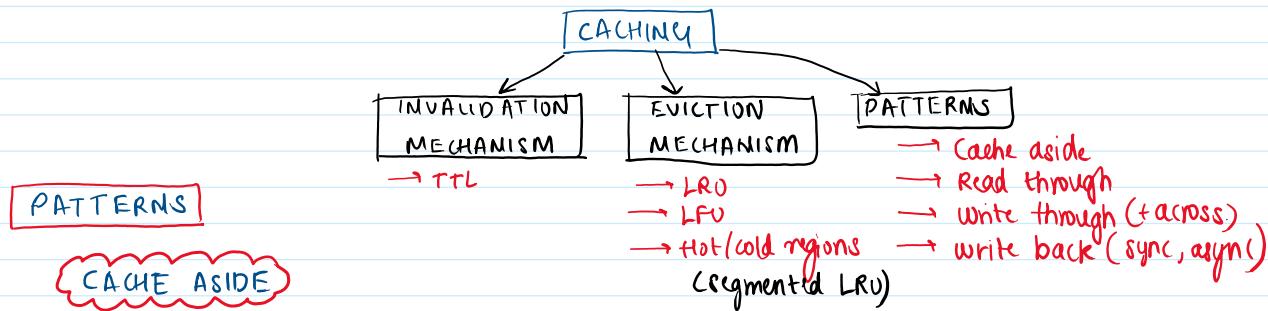
Implementation

→ 3 microservices
→ (3×2) Queue
 1 success 1 failure
 one bw each service
→ 2 Queues

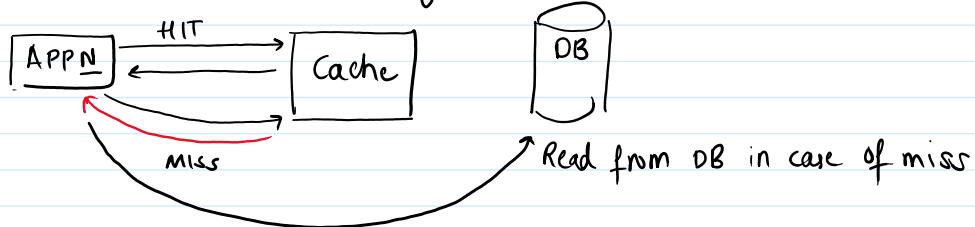
- ① Handles all successes (all w/service listen to this)
- ② Handles all failures OR listen to a specified partition

Caching

07 June 2024 23:04



→ Cache can't communicate directly with DB
Application does that for it



→ Initially data is loaded into the cache by app.
On every cache miss, the app directly fetches data from DB
→ On new data, it will always be Miss
until you update cache

PROS

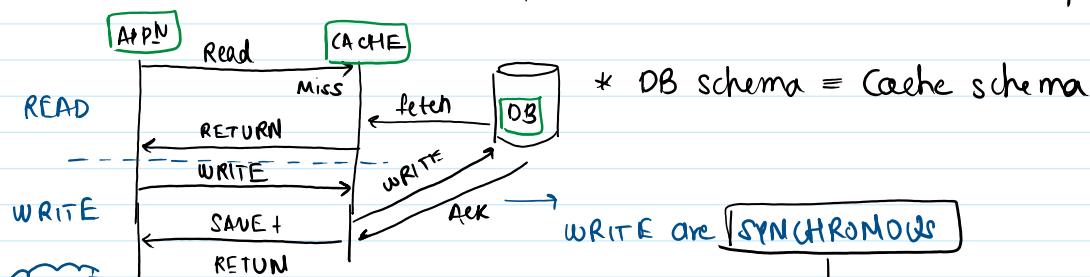
- * Good for HEAVY-READ apps like dashboards
- * If cache fails → you can directly hit the down DB each time
- * Since data stored in cache and DB is different, you can have different schemas

CONS

- * New data will be Miss
- * Inconsistency in data if it gets updated in DB

READ THROUGH

→ Cache can communicate directly with DB. App only accesses Cache
→ READ MISS: Cache itself updates value from DB
 → New val, update and then returns
→ Writing new data: you will write to cache + cache will write to DB
 → if this was update invalidate cache which will be updated on next miss



PROS

- Great for heavy reads why? (on write cache writes to DB, waits for ACK saves it in cache and then sends it back)
- You don't have all miss

- PROS**
- Great for heavy reads why? (on write cache writes to DB, waits for ACK saves it in cache and then sends it back)
 - You don't worry abt miss and fetch from DB
- CONS**
- Cache can't fail
 - Cache miss for new data (PRE-HEATING can solve)
 - Cache schema \equiv DB schema
 - Data inconsistency

WRITE THROUGH

- used with READ THROUGH (100% CONSISTENT)
- Here for writes you directly write to DB and invalidate cache (if updated)
 - Next hit will be cache miss and cache will fetch from DB
(saves latency of app \rightarrow cache \rightarrow DB for write operations)

WRITE AROUND

→ You ONLY WRITE TO DB and forget

* when you do a get hit time: CACHE miss and cache will fetch from DB

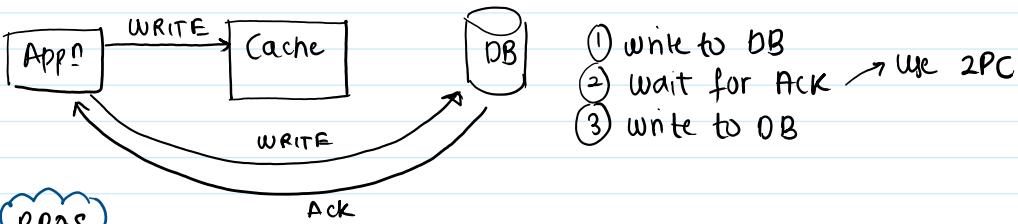
WRITE BACK

- SYNCHRONOUS** — same as write through (BUT cache updates DB)

- ASYNCHRONOUS** — NOT 100% consistent

→ Not app itself

- App! itself writes to both DB and cache whenever a write happens



PROS

- Very high hits and low misses
- You always write to (DB + cache) each time. No lag

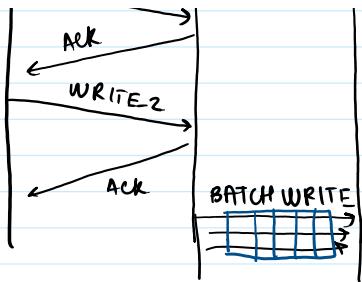
CONS

- High latency due to 2 writes and ACK wait
- Need 2PC to make sure DB update is ACID

ASYNCHRONOUS

- App! only writes to cache and then forgets
- Cache then asynchronously writes data into the database independently
 - done via QUEUES
 - you can publish 1 event for each write OR Batch n writes and at once write to DB





- lesser write latency and lower inconsistency
- better performance
- CONS** → sync issues if you immediately read again before DB write is done

Bolt with READ THRU WITH CACHE

why? → when write happens u just write to cache and forget
→ Now you can directly read that data from cache (even though its still not in DB)

GLOBALLY DISTRIBUTED CACHE

- ① There can be (N) different cache servers distributed across clusters
 - A cache gateway redirects requests to appropriate servers
- ② Default uses **consistent hashing** for allocation
- ③ Geographically based - route/store in nearest store
- ④ Identifier based - Ex: Based on country-id group and store in same server
or multiple countries into 1 continent server

REDUNDANCY

- + use **QUORUM** u need consensus for both READ + WRITES by setting values if R , N
- you need to write data to atleast (W) copies/machines
- R → Read consensus from (R) nodes
- W → write consensus from (W) nodes
- N → Replication factor

$R + W \leq N$

EVENTUALLY CONSISTENT

$R + W > N$

STRONGLY CONSISTENT

→ $R + W \leq N$: assume $N=3$

* $R=1, W=2$

WRITE : u write to 2 copies and leave third

READ : needs only 1 copy

→ if it reads from 3rd then inconsistent

* $R=2, W=1$

→ WRITE : to any 1 assume (1)

→ READ : (1) valid

(2)(3) INCONSISTENT
(3)(1) valid

$R + W > N$

Ensures that there is atleast 1 overlap b/w write and read consensus

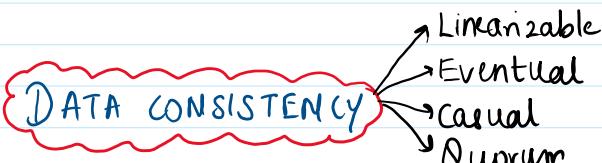
READ (R)



WRITE (W)

Isolation Levels

08 June 2024 14:54



LINEARIZABLE

- single thread sequential execution of Queries
- No chance of any mismatches
- **PERFECT** consistency

EVENTUAL

- DB might not be always consistent BUT eventually consistent
- multiple Queries/read/write can happen of multiple threads/servers
 - ↳ parallel proxy so is faster
- * But sync takes some time after transactions

CASUAL

- Allows parallel transactions BOT operations in single transaction are linear/serialize
- Consistent within transaction
- Better than causal but not as good as linear

CONS

Aggregates/joins are inconsistent

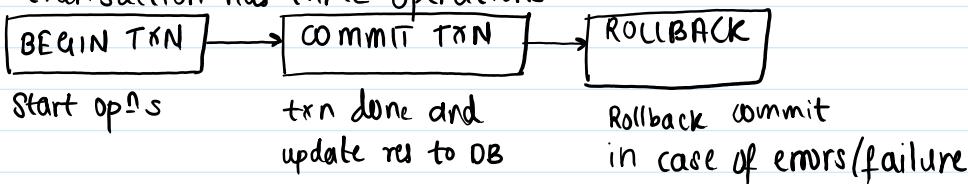
QUORUM

→ Configurable consistency vs performance

$R + W > N$	Highly consistent
$R + W \leq N$	Eventually consistent

TRANSACTION ISOLATION LEVELS

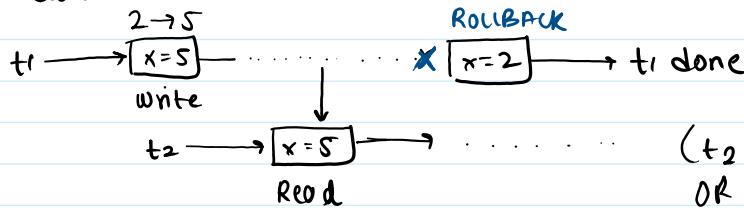
A transaction has three operations



READ UNCOMMITTED

- Each txn directly writes to DB ≡ No concept of commits (No data replication)
- Causes **DIRTY READS**





READ UNCOMMITTED

- Txn can only read values which are committed by other txns
- Slower since you have commit step → it won't roll back later and don't directly update DB

- * In both read uncommitted/uncommitted you can't ensure REPEATABLE READS
- if you have 2 reads of same var in the txn you can get 2 different values (someone else must have updated it)

REPEATABLE READS

SNAPSHOT ISOLATION

- when any row is read by a transaction
 - its LOCKED until that transaction ends
- other txns which update this locked value can locally commit and keep it as new version and update in DB once lock is released

SERIALIZABLE

- All operations in txn are executed sequentially in order
- if 2 txns don't have any common ops then they can run in parallel
 - Both writes and reads
- use serialized locks

*least isolation
most efficient*

*Least efficiency
most isolation*



CONFLICT-FREE REPLICATED DATA TYPE

KINDS OF CRDTS

→ You are in a multi-writer DB duplication setup

→ Now at same time the same key may get updated

WHICH VALUE WILL YOU PICK?

① Last write wins (very bad)

② Store all writes and decide which to pick later

③ Let DB take that call

= CRDT

INFO ABOUT CRDTS

→ The conflict merge will happen eventually (NOT immediately)

→ Which DB uses CRDTS? → Riak

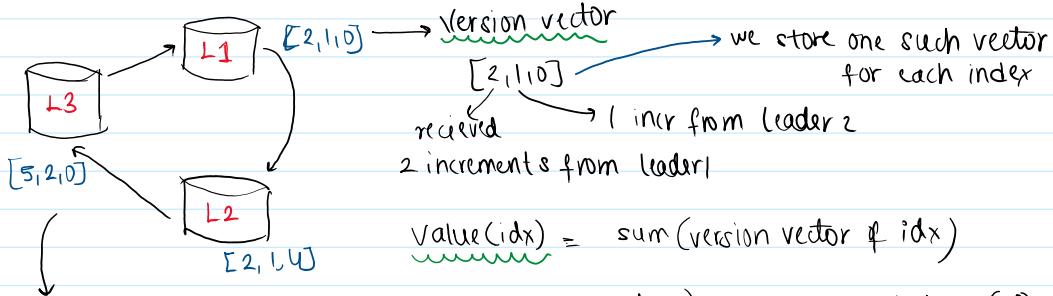
Redis (sets in Redis Enterprise)

OPERATIONAL CRDTS

→ You design a distributed counter on a multi-leader/writer setup → $\text{INCR}(idx)$ = increments value of idx by 1

but this can be simultaneously done on separate leaders at a time

VERSION VECTORS



→ Here instead of passing list of size(n -leaders) we just call $\text{incr}(s)$ on the other leaders

→ Ex: - Instead of $[5, 0, 3]$ being saved

you would have just $\text{incr}(s)$

③ $[5, 2, 0] \rightarrow \text{incr}(L1, +1) \times 5 \rightarrow \text{incr}(L2, +1) \times 2$ → This takes $O(1)$ instead of $O(n\text{-leaders})$.

OPERATIONAL CRDT

SOME CONS

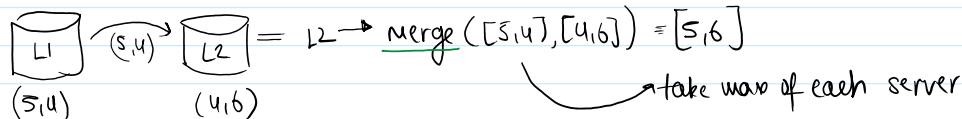
① You need to ensure no drops and no duplicates (when incr for example)

② Need to be causally consistent (same order followed)

(Ex: if you make a set using op! CRDT → $\text{remove}(idx)$ can be done only after $\text{add}(idx)$ has happened)

STATE BASED CRDTS

→ Instead of sending just the increments, SEND ENTIRE CRDT



The merge needs to be → ① COMMUTATIVE

$$m(a, b) = m(b, a)$$

The merge needs to be →

- (1) COMMUTATIVE → $m(a,b) = m(b,a)$
- (2) ASSOCIATIVE → $m(a, m(b,c)) = m(m(a,b), c)$
- (3) IDEMPOTENT → $m(a,b) = m(m(a,b), b)$

→ this handles in case of duplicate msgs received

HOW DO U COMMUNICATE AMONG WRITE-LEADERS

using state (CRDT)

- (1) Dont need order (ASOC + COMM)
- (2) duplicates refine (IDEMP)

GOSSSIP PROTOCOL

→ No extra/special infra needed

TYPES OF CRDTS

NAME	OPERATIONAL CRDT
Counter	incr(leader)
Decreasable Counter	incr(leader), decr()
Sets	add(x), remove(x)

↓ Hard to implement

STATE-BASED CRDT

inrcs: [0, 2, 1]

incs: [0, 2, 1] decs: [0, 1, 0] = [0, 1, 1]

final counts



Which Databases?

22 December 2024 00:11

DB INDEXES

B+ Trees

- Stored on disk
- slower writes
- faster reads
- NO size issues

LSM TREE + SS Tables

- LSM : memory
- SS-Table : disk (only when needed)
- slower reads (read across SS-table)
- faster writes (write to memory)
- delete = writing a tombstone

B+ Trees: [How do B-Tree Indexes work?](#) | Systems Design Interview: 0 to 1 with Google Software Engineer
LSM Trees + SS Tables: [LSM Tree + SSTable Database Indexes](#) | Systems Design Interview: 0 to 1 with Google Software Engineer



REPLICATION

- ① Single leader → 1 Node handles writes / can read from many replicas
- ② Multi leader → You can write to multiple nodes + Read multiple
 - changes of write-merge conflicts
- ③ Leaderless → You write/read from any node/replica
 - uses multicast/gossip protocol among nodes



MySQL

- Relational/Normalized data (uses 2Phase-commits for consistency)
 - ACID + transactional guarantees
- very slow but sturdy
- Uses B+ Trees + single leader replication

Mongo DB

- Document based NoSQL DB (data written as large documents)
- BTrees + transaction support

NOTE Nothing very special but acts like SQILish NoSQL DB

- transactional support

APACHE CASSANDRA

- Wide column DB : has shard key + sort key
 - flexible schemas
 - easy to partition
- Multi leader / Leaderless replication (configurable using Quorum)
 - super fast writes but uses last-write-wins for write conflicts
- Index : LSM Tree + SS Tables
 - fast writes

USE CASE Apps with high write volume / speed · consistency is not so important

RIAK

- key-value based NoSQL DB
- multi-leader / leaderless BT supports CRDTs (manual code for merge conflict)
- SSTables + LSM trees

USES : high write throughput + controllable write conflicts

HBASE

- Wide column DB — shard key + sort key
- Single leader replication
 - Built on top of HDFS : durable and reliable
 - slower for frequent ops
- LSM Trees + SS tables
- column oriented storage
 - column compression + data locality in column

INRAPS : In-row data + Full column reads

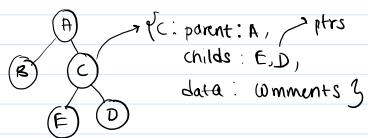
→ column oriented storage
→ column compression + data locality in column w

Uses : Large data + fast column reads

GRAPH DATABASES

NATIVE GRAPH DB (Neo4j)

→ you store each node and memory ptrs to its adjacent nodes within the representations



NON-NATIVE GRAPH DB

→ you store data in some tabular format and later perform joins / opns to get relations

id	par-id	comment
A	-	good post
B	A	agree, +1
C	A	ur joking
E	C	+1
D	C	ur right

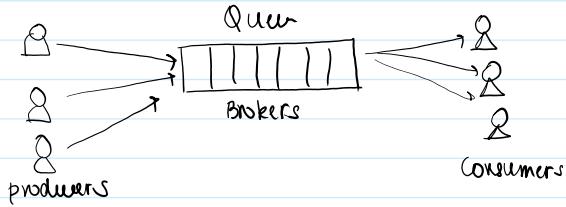
- + Faster lookups since addresses are stored
- Inserts/updates are slower
 - ① You need to traverse to insert location
 - ② Create node and find mrs/par/children
 - ③ Create links and store

- + Insertions are easy → just adding rows
- Queries/getting relations are tough / slower
 - since you don't have addresses you need self joins (in this case)

Streaming

22 December 2024 00:40

- Generally works on concepts of PRODUCERS + CONSUMER



MESS AUF BROKERS

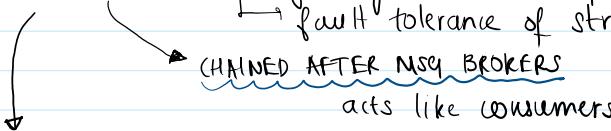
IN-MEMORY (RabbitMQ)

- stateless + in memory
- PUSH model
 - sends msgs to consumers + waits for ACK
- At-least once semantics
- very low latency

LOG-BASED (kafka)

- stateful + events written to a log
- PULL model
 - Consumers read at their own pace
 - uses offset-based model
- at least once, at most once + exactly once semantics
- high throughput and durable

STREAM PROCESSING



Spark streaming

- micro batch processing

Flink

- realtime on event at a time

APACHE FLINK

- guarantees that each msg affects state only once
- Flink checkpoints your state and dumps it to storage (SS)
 - o you can replay the events from your log-based queue
 - o Done by a job manager/zookeeper