

SHARDING

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 location

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

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

$\left(\frac{h(\text{data})}{N} \right)$ → you can't change this
hash fn no. of servers if done needs rehashing everything

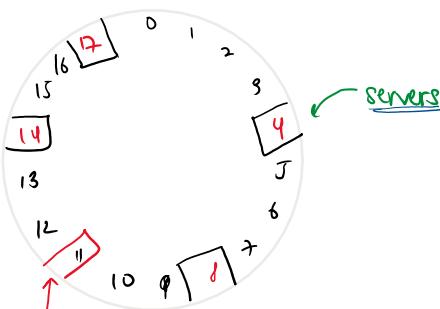
Consistent Hashing

→ 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 9,10 which was prev sharded to 14 need 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.11 which has 10 slots

- Add more server slots which brings load factor to 'N'
- Internally a single server handles (2^N slots)
- Ex:- ⑧, 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' accnt

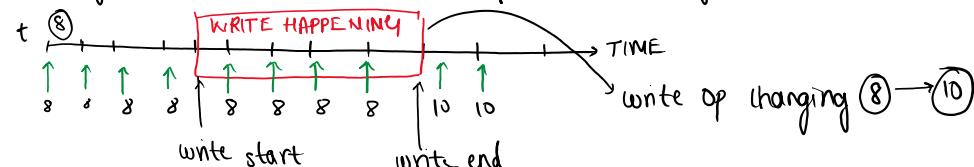
Now this shud happen completely. Ex lets say it got deducted from accnt but then failed

- 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 transactn 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 as 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



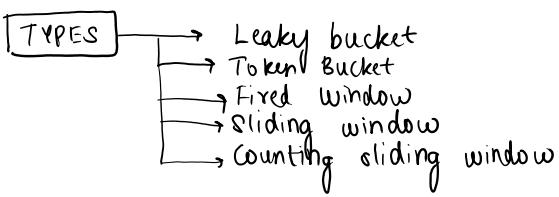
Eventually Consistent

→ Not immediately consistent → ACID ones are immediate consistent

→ 2 simultaneous reads may give diffent 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

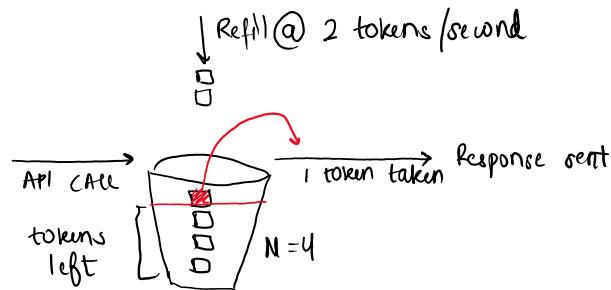
→ NO → send 429 response code

→ YES → Take one token from the bucket

⇒ 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

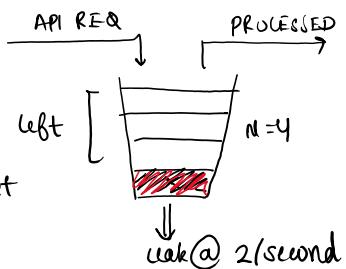
→ A request comes in, if bucket

→ 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

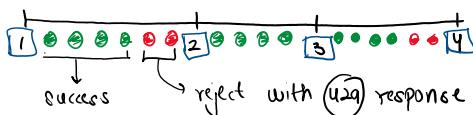
* IDEAL CASE → Rate of API hits \leq LEAK RATE



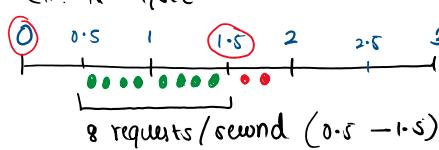
FIXED WINDOW

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

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



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



valid since our windows are 0-1-2-3
But truly not limiting at Boundaries

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

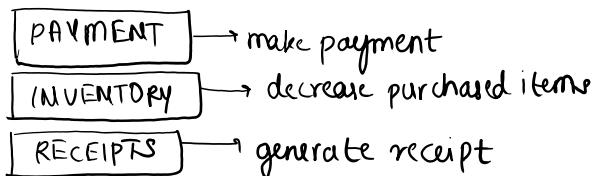
→ Here instead of 1000 logs you will only have 10 logs

→ space = 1000 \rightarrow 10

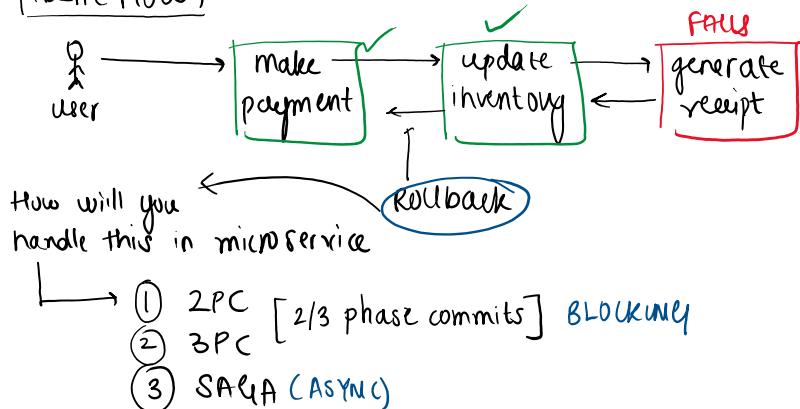
→ No of ts steps

DISTRIBUTED TRANSACTIONS

PREMISE → You have 3 services



IDEAL FLOW

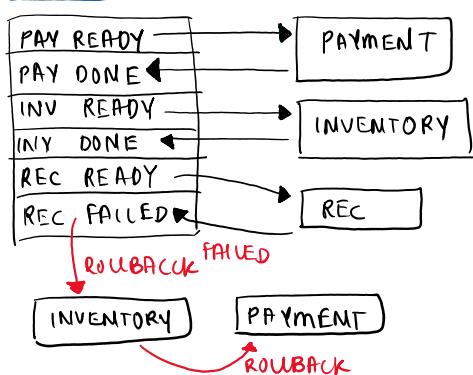


SAGA

* Asynchronous and non-blocking but complex to implement

→ Queues are used for Inter Process Communication

STEPS



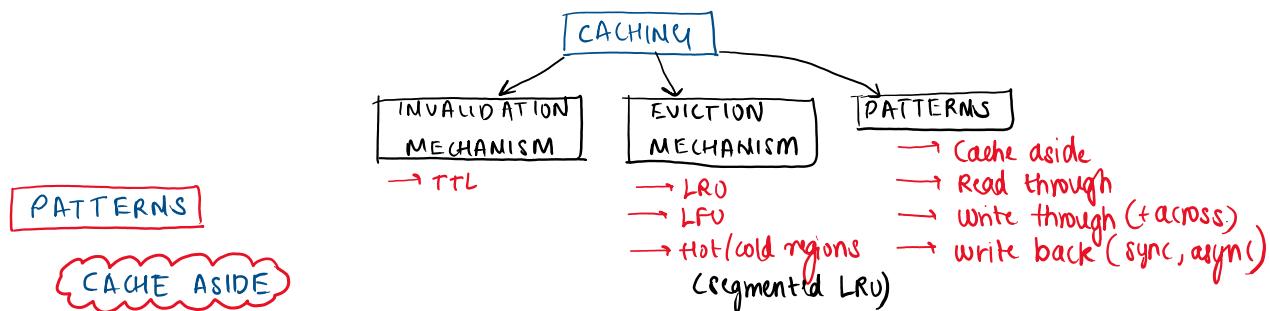
Implementation

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

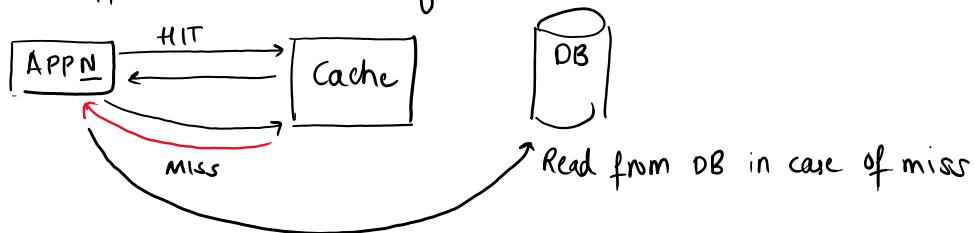
→ 2 Queues
 (1) Handles all successes (all mservice listen to this)
 (2) Handles all failures (OR listen to a specified partition)

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

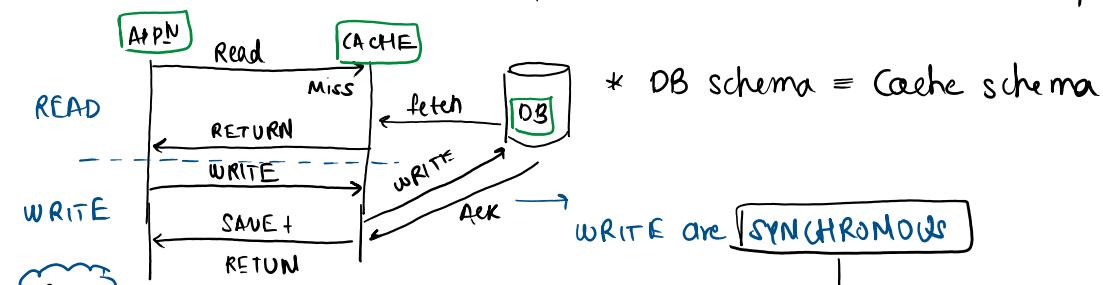
- * Good for HEAVY-READ apps like dashboards
- * If cache fails → you can directly hit the down DB each time
- * Since data stored in cache and DB is different, you can have different schemas

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 incur all miss

PROS

- Great for heavy reads · why? (on write cache writes to DB, waits for ACK and 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 \neq 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

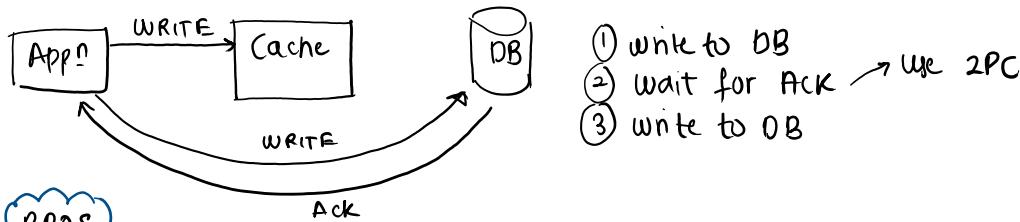
* when you do a get hit time: CACHE miss and cache will fetch from DB

WRITE BACK → SYNCHRONOUS — same as write through (BUT cache updates DB)

SYNCHRONOUS → ASYNCHRONOUS — NOT 100% consistent

→ Not app itself

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



PROS

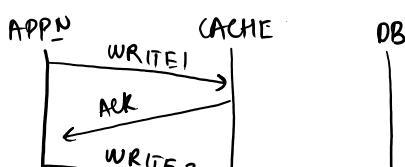
- Very high hits and low misses
- You always write to (DB + cache) each time · No lag

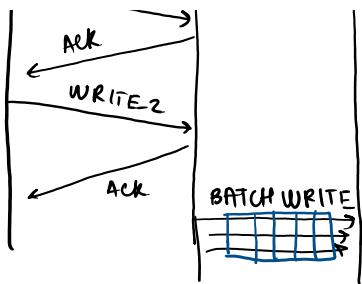
CONS

- 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
- sync issues if you immediately read again before DB write is done

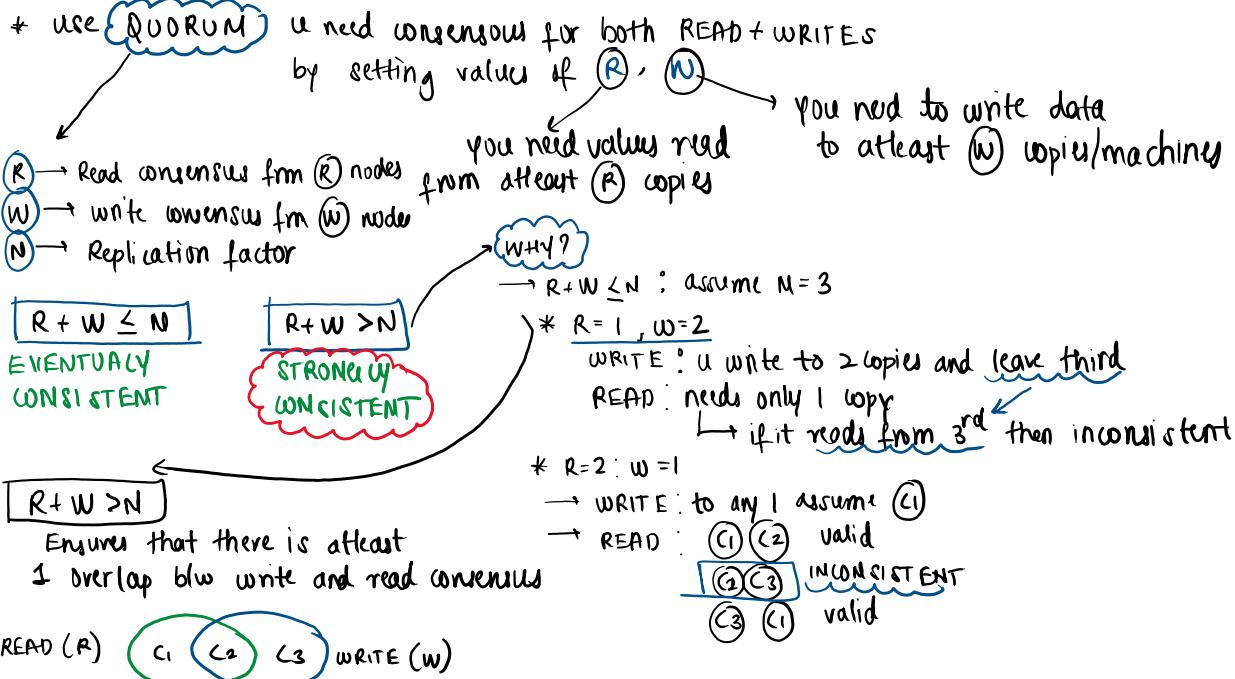
Bolt with READ THRU WITH CACHE

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
 - a) Default uses consistent hashing for allocation
 - b) Geographically based - route/store in nearest store
 - c) Identifier based - Ex: Based on country-id group and store in same server
or multiple countries into 1 continent server

REDUNDANCY



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 are linear/serialize in single transaction
- Consistent within transaction
- Better than causal but not as good as linear

CONS

Aggregates/joins are inconsistent

QUORUM

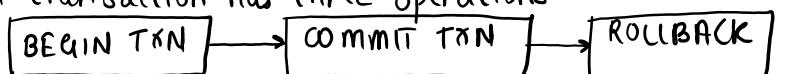
- Configurable consistency vs performance

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

Highly consistent
Eventually consistent

TRANSACTION ISOLATION LEVELS

A transaction has three operations



Start op's

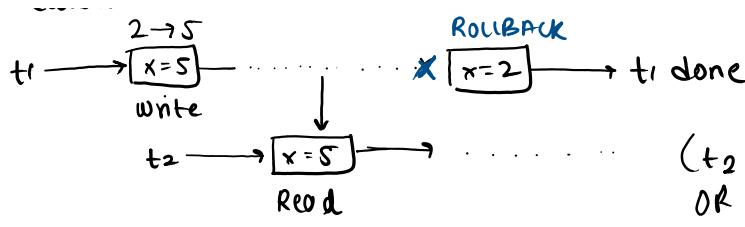
txn done and update res to DB

Rollback commit
in case of errors/failure

READ UNCOMMITTED

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





(t_2 never gets rolled back value of X)
OR some other op! might also have updated it

READ UNCOMMITTED

- Txn can only read values which are committed by other txns
- slower since you have commit step → it wont roll back later and dont directly update DB

* In both read uncommitted/uncomitted you cant 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

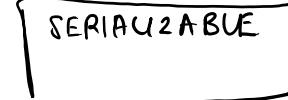
- when any row is read by a transaction
 - its LOCKED until that transaction ends
- other txns which update this locked value can locally commit and keep it as new version and update in DB once lock is released

SERIALIZABLE

- All operations in txn are executed sequentially in order
- if 2 txns dont have any common ops then they can run in parallel
 - Both writers and reads
- use serialized locks

least isolation
most efficient

Least efficiency
most isolation



← EFFICIENCY

ISOLATION →

CRDT

19 December 2024 19:34

CONFlict - FREE REPLICATED DATA TYPE

[KEY CRDTS]

- You are in a multi-writer DB duplication setup
- Now at same time the same → WHICH VALUE WILL YOU PICK?
- key may get updated

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

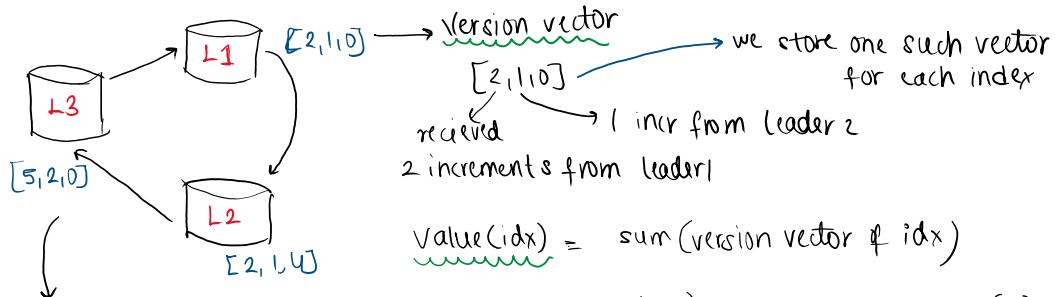
Redis (sets in Redis enterprise)

[OPERATIONAL CRDTS]

→ You design a distributed counter on → $\text{INCR}(\text{idx})$ = 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 $[5, 2, 0] \rightarrow \text{incr}(L1, +1) \times 5 \rightarrow \text{incr}(L2, +1) \times 2 \rightarrow$ This takes $O(1)$ instead of $O(n\text{-leaders})$.

[OPERATIONAL CRDT]

SOME CONS

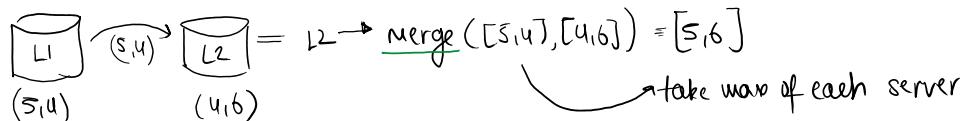
① You need to ensure no drops and no duplicates (when incr for example)

② Need to be causally consistent (same order followed)

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

[STATE BASED CRDTS]

→ Instead of sending just the increments, SEND ENTIRE CRDT



The merge needs to be → ① COMMUTATIVE $\rightarrow m(a, b) = m(b, a)$

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 (RBT)

- (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

final counts

inrc: [0, 2, 1]

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

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

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

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](#)



MySQL

- Relational/Normalized data (uses 2Phase-commits for consistency)
 → ACID + transactional guarantees
- very slow but sturdy

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 Quorom)
 - 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

INNS : In-row data + full column reads

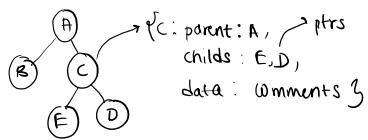
- column oriented storage
 - column compression + data locality in column

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



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

NON-NATIVE GRAPH DB

- You store data in some tabular format and later perform joins/ops to get relations

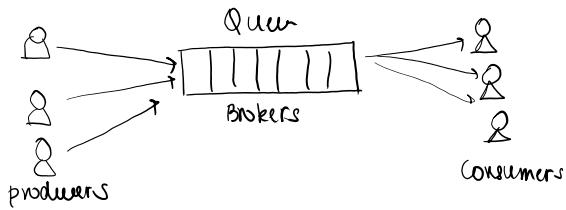
id	par_id	comment
A	-	good post
B	A	agree,+1
C	A	ur joking
E	C	+1
D	C	ur right

- + 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



MESSENGERS

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

joining streams and processing

Spark streaming

- micro batch processing

fault tolerance of streams
acts like consumers

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)
 - you can replay the events from your log-based queue
 - Done by a job manager/zookeeper

Architecture

04 April 2024 18:27

Tiny Url + Pastebin

16 December 2024 19:57

REQUIREMENTS

- ① Generate short url for each long url
- ② user must be able to see number of clicks for each url

TINY URL

PASTE BIN

PERFORMANCE REQS

- ① median url has 10k clicks. Most popular may have millions
 - ② At time, you may have 1 trillion urls
 - ③ Read speed > write speed (in terms of count also)
- avg = 1 kb → $1 \text{ Trillion} * 1 \text{ kb} = 10^{12} * 10^3 = 10^{15}$ (PB)

LINK GENERATION

DELETION / INVALIDATION

- hash (long-url, user ID, createTimestamp)
- combination minimizes collisions

→ simple CRON job to delete

* how many chars

will your tiny url contain? 8 chars → hash collisions can still happen

→ $(0-9, A-Z) = 36 \text{ chars} \Rightarrow 36^n \text{ possibilities}$

for $n=8, 36^8 = 2 \text{ Trillion}$ (which is more than enough)

DB REPLICATION (FASTER WRITES)

* Leaderless or multi-leader replication? NO
↓ USE

lets say you get hash of url, but that already exists in DB

→ you need to immediately find this (so that you can create another

→ In multi-leader the node you write to hash)

might not have hash, but another node might have it (got parallelly written)

SINGLE LEADER REPLICATION

+
Partitioning/sharding on HASH
(speeds up writes ready)

→ can support multiple writes since sharded across clusters

→ In case of conflict you write nearest next hash value

→ you can additionally use consistent hashing
↓ most probably will lie in same node

LOCKING

→ lets say 2 users try to write/update same row. You would need some locking mechanism

① Predicate Query

→ obtain locks on rows which don't even exist yet

→ This needs you to scan entire table for 'short-url'
can be sped-up by indexing on 'short-url'

② Materialize locks

→ make sure all possible short id rows are present in DB
 $2 \text{ Trillion} * 1 \text{ byte} * 8 \text{ chars} = 16 * 10^{12} = 16 \text{ TB}$

→ you can feasibly have a DB with size 16TB with all rows pre-computed/filled

DB Engine Impl

B-TREE

(can't store TB data in memory)

BEST SUITED FOR US

→ fast reads, slower writes
→ stored on disk

LSM + SS Tables

→ faster writes, slower reads
→ in memory + ss Tables stored on disks
↓ This makes retrieval/reads slower

CACHING HOT/POPULAR LINKS

* Can you use where DB loads certain keys to pre-populate cache? NO
→ you can't predict which url might become popular

MECHANISMS

EVICTION METHOD?

→ LRU

push cache keys to pre-populate cache? url might become popular

MECHANISMS

EVICTION METHOD?

→ LRU

- ① wire back cache (NO) → you can wait and later inform about conflicts
- ② wire through (NO) → overhead to write both cache + DB (you don't need it in cache immediately)
- ③ wire around (YES)
 - directly write/update in DB
 - first read will be cache miss, but we're fine with it

CLICK COUNTS

- ① Traditional DB (NO) → too many concurrent writes
- ② Atomic counters (NO) → too complex implementations
- ③ MESSAGE QUEUES
 - IN-MEMORY msg brokers (NO) → fast but not durable
 - topic-based msg brokers (YES) (use Kafka)

Partition brokers based on short-urls to make it faster

MESSAGE CONSUMPTION

- ① HDFS + Spark batch job (NO) → too infrequent updates
- ② Flink (NO) → it processes each event/msg at a time (might get crowded)
- ③ Spark streaming (YES) → supports 1 msg at a time processing
 - + mini batch processing in case of bottlenecks

use EXACTLY ONCE semantics

PASTE BIN - THREE PASTES

DATA STORAGE :- object based blob stores (S3)

only put and fetch

No updates or queries r done

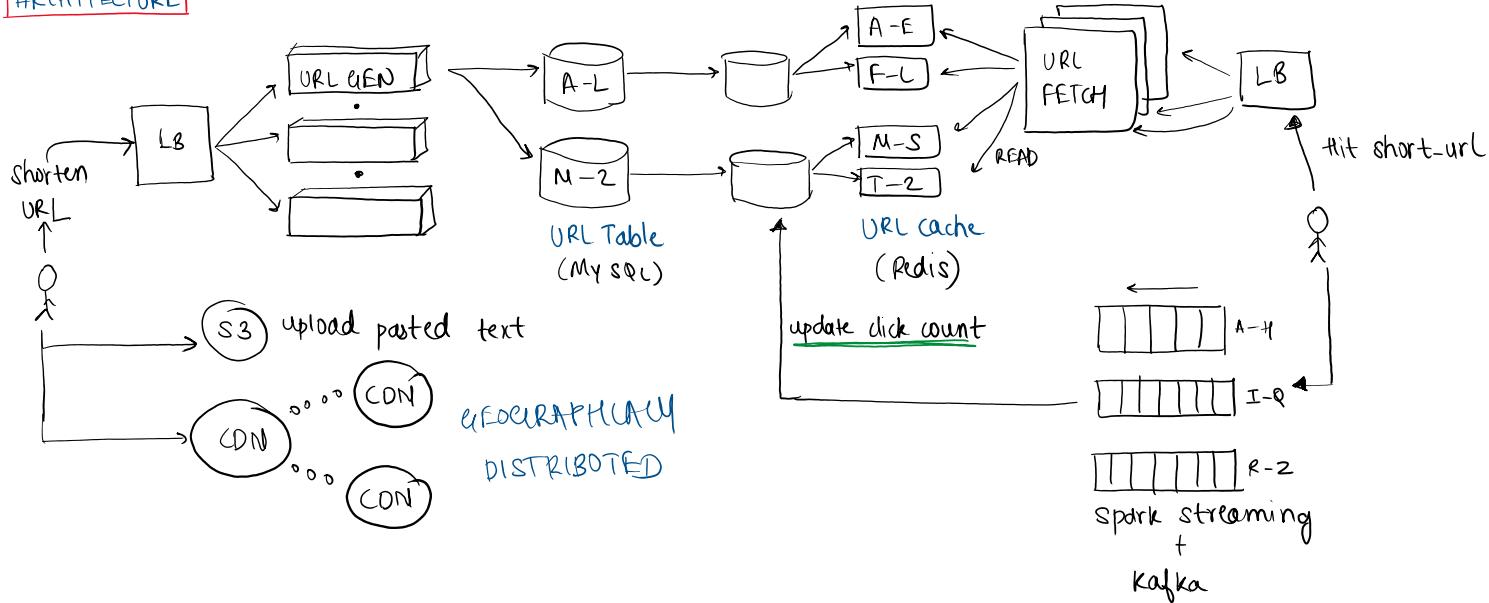
cheaper than HDFS / DBs

HUGE FILES → use a CDN → will distribute geographically

CACHING → write around? (NO) write is fast but first cache miss is expensive (huge file)

write through (YES) → slower writes is fine but no initial cache miss

ARCHITECTURE



Rate Limiter

17 December 2024 08:19 PM

CAPACITY EST

- ① 1B users
- ② 20 services to rate limit
- ③ fields → user-id, count
(8 bytes) (4 bytes)
- ④ Total storage needed = $10^{12} * 20 * 12 = 240 \text{ eB}$ → pretty small considerably

RATE LIMITING KEY

USER-ID

- + Easy to track and block users
- Dummy/fake accounts
- Not all services need user auth

IP-ADDRESS

- + No need for auth + user data
- + Multiple accounts handle
- Multiple users using same n/w (wifi)

choice? Authenticated service? → user-id based

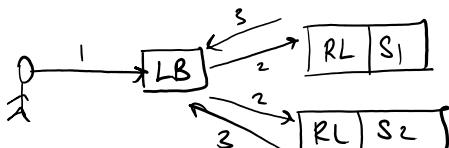
Unauthenticated service? → IP-address based

API format → bool ratelimit (user-id, ip-address, service-name, request TS)

WHERE TO PLACE RL

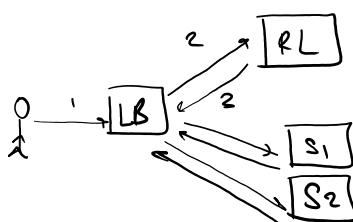
Local to service

- + less components to manage
- + lesser n/w calls
- spammed requests may use up n/w bandwidth (payload size)
- tightly coupled with service



Distributed as service

- + Easily scale up/down
- + Handle spammed requests easily
- Added api calls (call RL → call service)



CACHING for LB

→ Best is **write back** → whenever you hit LB it will cache count of hits and return the redirected url
→ (+) increment will eventually get updated in RL

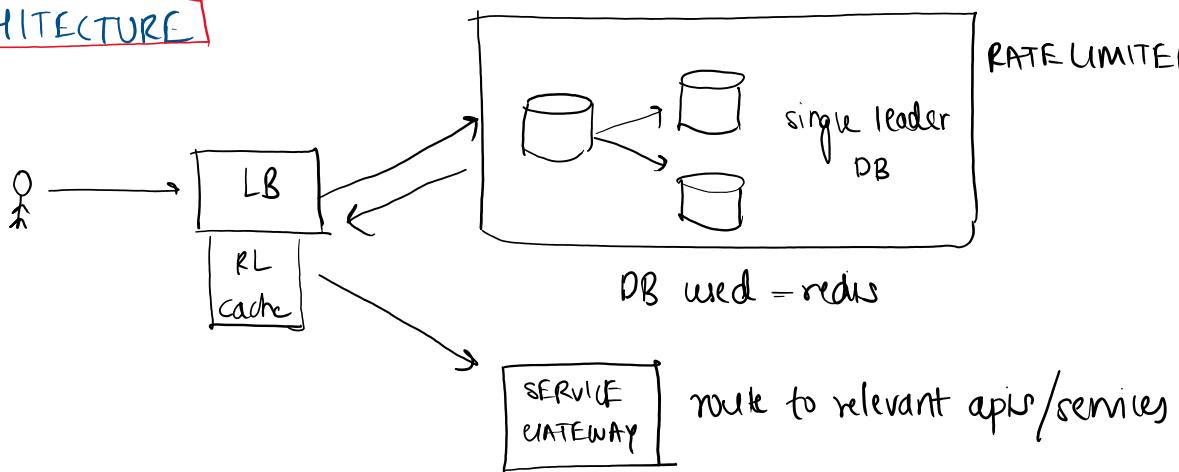
- and return the redirected url
 → (+1) increment will eventually get updated in (RL)
- which cache to use?
- ① Disk based (NO) → we want it to be quick + data here is smaller
 - ② In-memory based → (YES)
 (Redis)

DB REPLICATION

- LEADERLESS / multi-leader
- + faster write throughput due to multi-parallel writes
 - Counts may not be accurate since it needs to get propagated among replicas

- single-leader write → WE USE THIS AS ALREADY IS MORE IMP
- + Accurate counts since you write to a single leader
 - slower writes since single write leader

ARCHITECTURE



REQUIREMENTS

- Users must be able to edit documents simultaneously
- The updates need to be propagated to all users eventually

CAPACITY

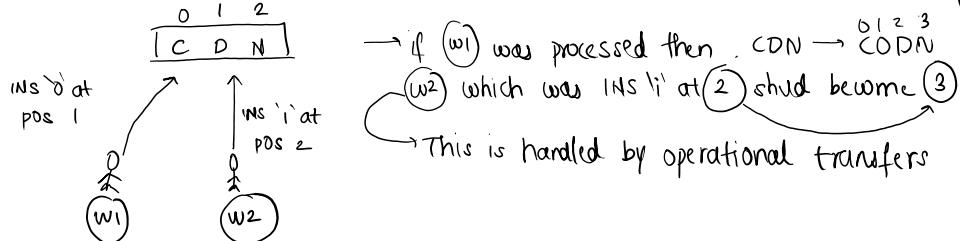
- Billions of users and documents
- Each doc can have 1000s of concurrent editors/viewers
- Doc size is in kbs

SOME BAD IDEAS

- ① using websockets
- propagate updates
 - propagate entire document
- ? Those r sent by each user to a centralized server
- CONS
- file size and num_users may be huge
 - propagate updates - you need handle conflicts

OPERATIONAL TRANSFORMS

- clients send writes to server as they please. Server handles them gracefully

**WORKING**

- Centralized server takes in all the writes and transforms them if needed
- Here order in which writes were processed matters
- Can be bottlenecked due to so many users + concurrent writes

CRDT

→ BEST DESIGN SOLUTION HERE

STATE CRDT

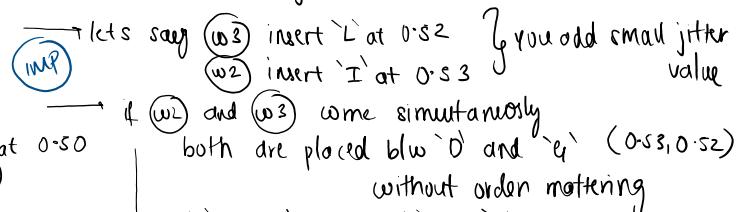
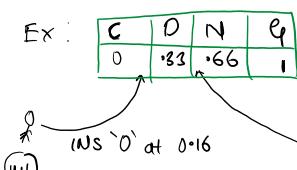
- Send entire document
 - Merge fn merges docs
- ISSUE doc size and writers are huge in number here

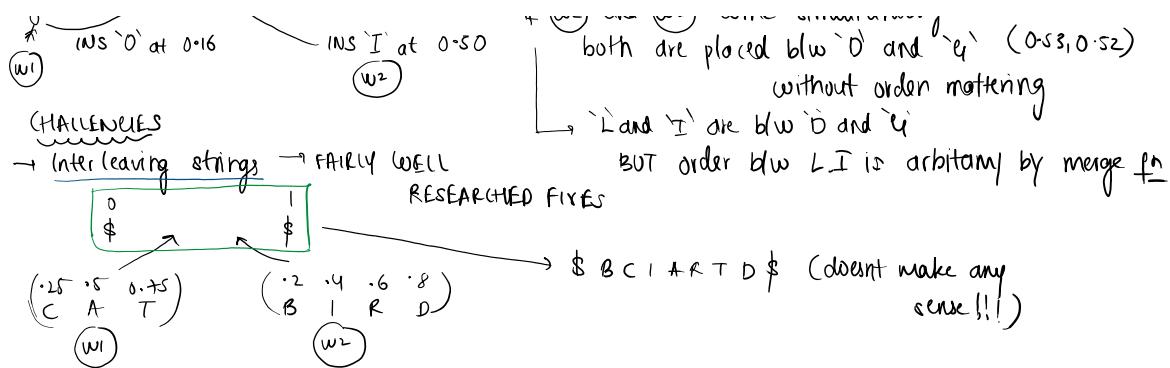
OPERATIONAL CRDT

- Just send op! (INS 'o' at pos 1)
- merge fn is commutative and associative (order won't matter)
- Idempotent? No (you will need state info for this impl) → using sort of a write id)
 - if you get same op twice it should be applied once

CRDT MERGE FN

- ① First char = 0 ② if you want to insert a char
Last char = 1 then pos = 0 + 1/2 = 0.5 + some random





Achieving Idempotence

- ways → UUID (extra space to be stored)
version vectors (we go with this)

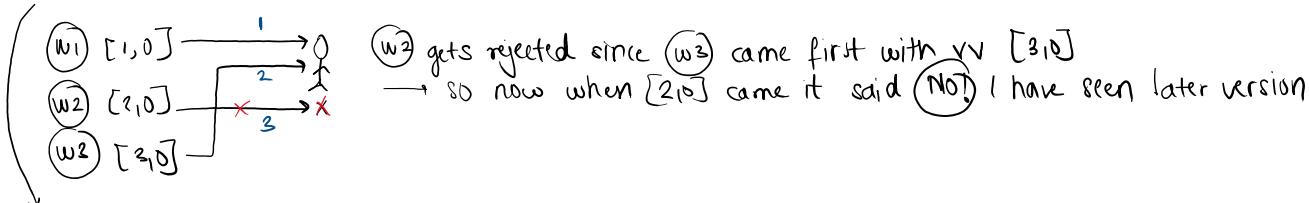
VERSION VECTORS

→ VV, which server made the write and how many writes that server did
 ↳ VERSION

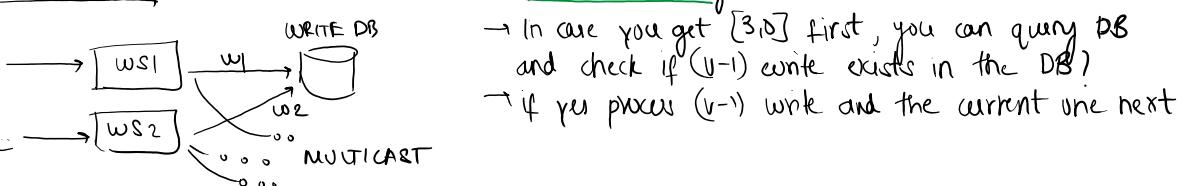
Ex - There are 2 write nodes/servers

(W1) [2,0] → if you have processed [2,0] now and later u get [1,0]
 ↳ No writes by s2
 2 writes by s1
 You ignore (u have already processed later version)

VERSION SKIP ISSUE



SOLUTION

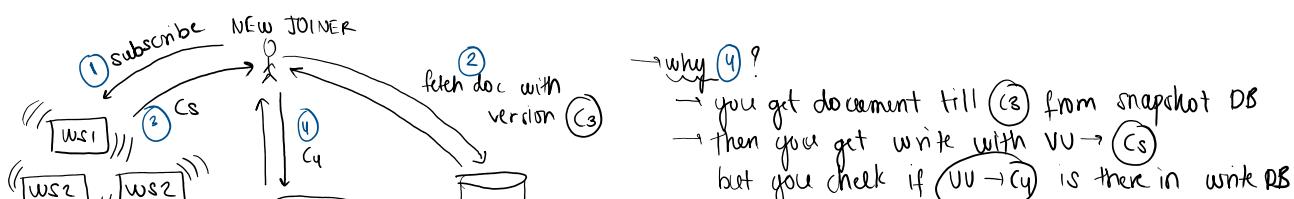
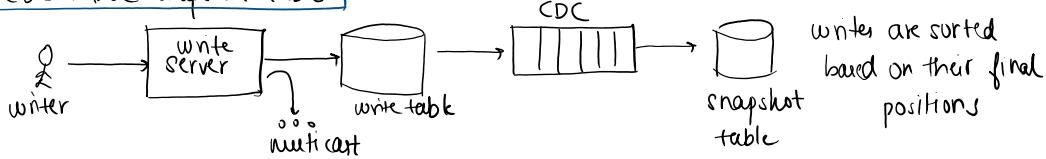


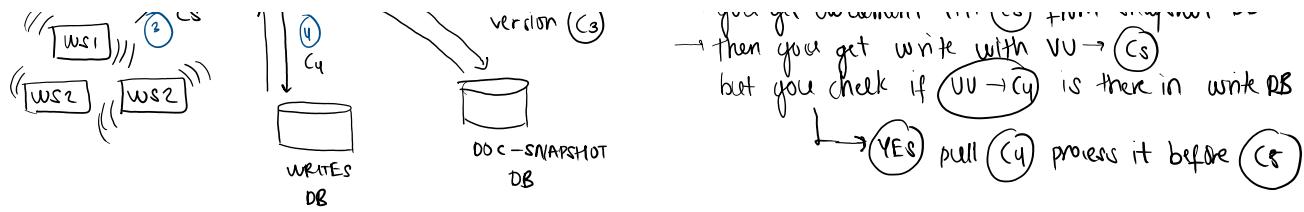
NEW JOINERS

↳ Basic: fetch the base document & apply writes

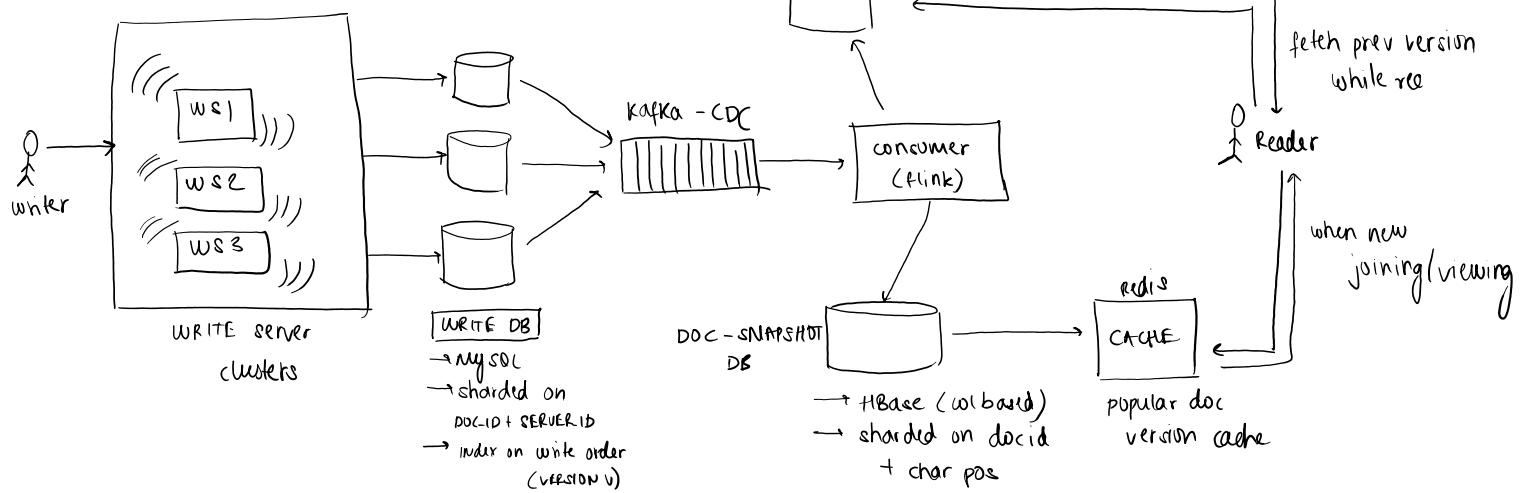
✓ BETTER

CDC + Doc snapshot table





ARCHITECTURE



REQUIREMENTS

- Users can post, watch, search videos *by name*
- You can comment on videos (single level of comments)

CAPACITY ESTIMATES

- 1 billion users
- Avg views on video in 1000's (some may have millions views)
- Avg video size = 100 Mb
- 1 million new videos each day

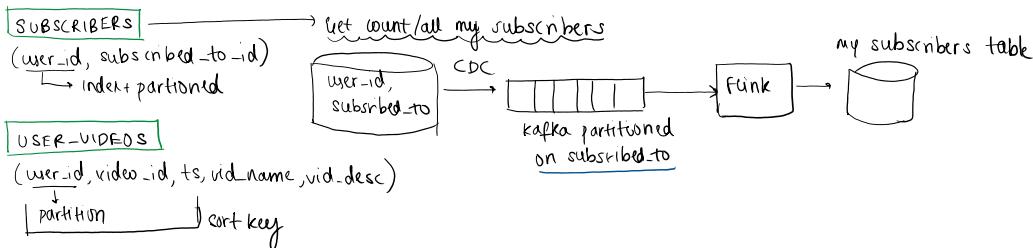
$$\begin{aligned} 100 \text{ MB} \times 1 \text{ M users} \times 400 \text{ days} &= 40 \text{ PB/year storage} \\ (100 \times 10^6) \times (10^6) \times 400 &= 4 \times 10^2 \times 10^6 = 4 \times 10^{16} \\ &= 40 \times 10^5 = 40 \text{ PB} \end{aligned}$$

VIDEO STREAMING

- ① Support multiple types of devices
 - ② Support multiple different n/w speeds
- G Basically [chunking] + [diff video resolutions]*

CHUNK + MULTIPLE RESOLUTIONS FOR EACH CHUNK**PROS**

- Parallel uploads (each chunk uploaded parallelly)
 - Lower barrier to start video (fetch first few chunks instead of whole video)
 - Grab next chunk based on n/w speeds of whole video
- You somehow maintain cur chunk + fetch lower/higher res of next chunk based on n/w speed

DATA BASE TABLES**VIDEO COMMENTS**

(video_id, ts, user_id, comment)

You can also use a channel_id + video_id

VIDEO CHUNKS

(video_id, encoding, resolution, chunk_order, hash, url)

- partition on video_id : all chunks must be on same node
- sort by video_id, encoding, resolution, chunk_order

If your cur is (f164+360p) → next chunk will likely be the same

WHICH DB TO USE?

- We have 1000x reads than writes : Btrees based ones are fine
- lots of video comments
 - * Use CASSANDRA → supports high write throughput (multi leader rep/?)
 - write conflicts are rare (we don't have nested comments)

VIDEO UPLOAD

- This is not very frequent and can be async
- We want this to be cheap + correct
- for each video → split into chunks → multiple resolution
 - (we don't care about order of processing here)

STEPS FOR UPLOAD

- ① Chunk video
- ② Upload chunks to object storage
- ③ Let a queue process chunks
- ④ monitor till all chunks are processed

which msg Broker to use?

In-memory (RabbitMQ)
log based (Kafka)

- why?
 1) Order doesn't matter
 2) We don't want control over delivery semantics
 3) Replicating missed logs is scarce
 4) No state being stored here

which object store?

HDFS

- + Data locality (huge cluster)
- + High processing power

- Costly due to CPU needed

* we pick S3

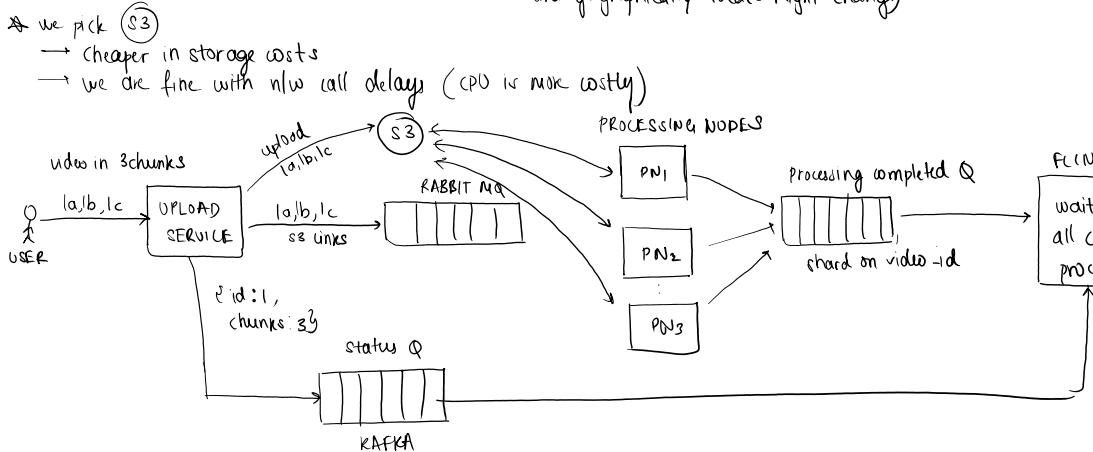
→ cheaper in storage costs

→ we are fine with inflw call delays (CPU is more costly)

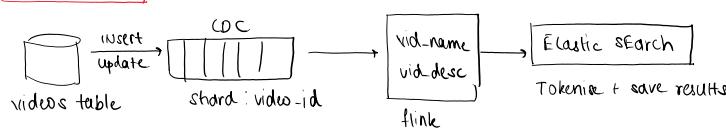
OUR CHOICE S3

- + Cheaper storage costs
- + No CPU costs → just upload n forget

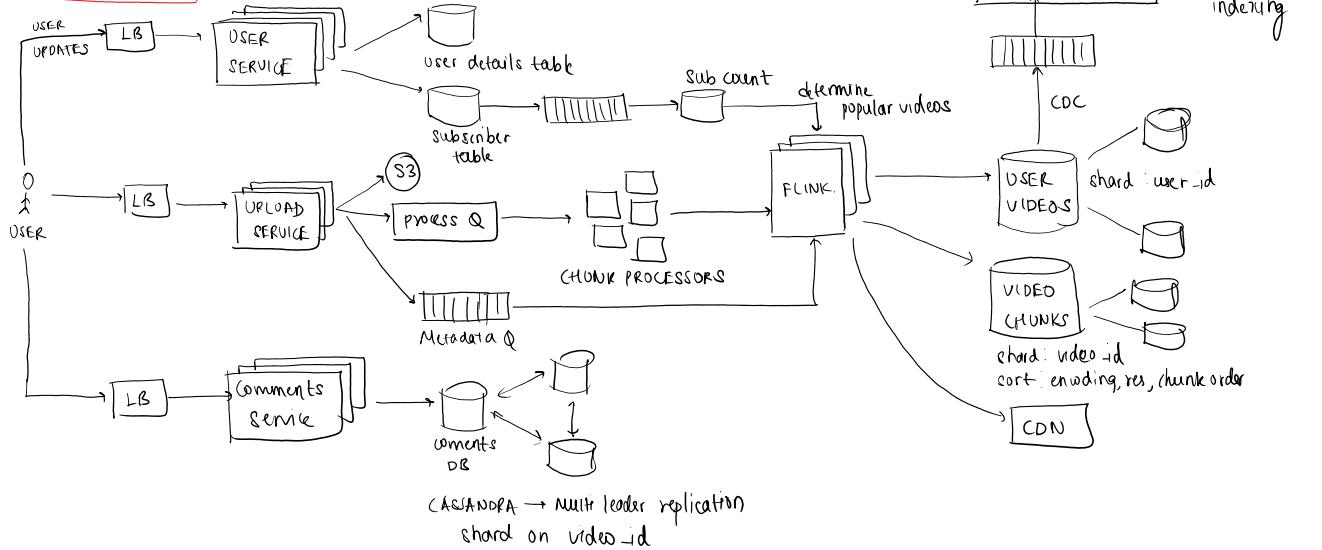
- inflw calls vary (each chunk needs a call
and geographically locat! might change)



VIDEO SEARCH → Elastic Search



ARCHITECTURE



Type-Ahead System

27 December 2024 01:58 PM

Resources:

- [System design - Design Autocomplete or Typeahead Suggestions for Google search](#)

REQUIREMENTS

- When typing in textbox, quickly load the suggestions.
- Top keywords are updated regularly (every 30-60 mins)
- Include words + entire sentences in suggestions
- Reads need to be almost instantaneous

CAPACITY ESTIMATES

→ Assume word suggestions (NOT SENTENCE)

- + word predict 3 more words
- 200K English words, avg len = 5 chars
- + word you need to store result of all its prefixes

apple = "a", ap, app, appl, apple

$$200K * 5 = 1 \times 10^6 \text{ entries}$$

$$1M \times (\text{1B word} + \text{1B for suggestions}) \\ = 1M * 20B = 200 \text{ MB}$$

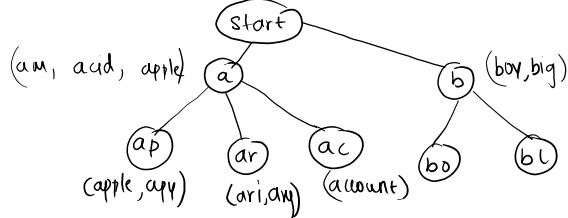
can be stored locally
on client device

Query Based

- There are 1B unique search queries each day \Rightarrow storing counts of each = 4B bytes = 4e9B → can't be stored locally
- Assume avg search term is 20 chars
 - ① 1B entries * 20 prefixes each = 20B entries
 - ② 10 suggestions for each word = $10 * 20 = 200 \text{ b/entry}$
- $20B * 200b = 4 \times 10^{12} = 4 \text{ TB storage each day}$

TRIE

- Obviously we will be using TRIE for storing our recommendations online



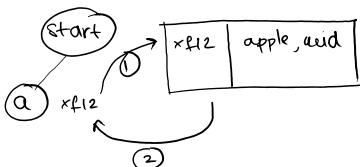
- while typing TC for getting next char suggestions is O(1)
 - "a" - am, acid, apple
 - "ap" = apple, apply

- TIME COMPLEXITY
- ① Insert word = O(len(word))
 - ② Delete word
 - ③ Search word
 - ④ Search prefix = O(P+k)

len of word
with prefix P

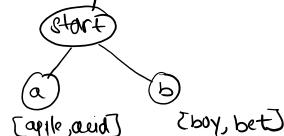
MEMORY EFFICIENT

- Store the references of suggestions at each node.
- You have to fetch it from another place - DB, disk etc
- + LESS MEMORY → only store addresses
- MORE TIME → fetch takes time



TIME EFFICIENT

- Store entire suggestions directly as list of words + node
- No need to fetch other places to get data
- + LESS TIME
- MORE MEMORY



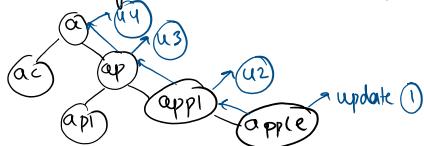
UPDATING KEYS

- Let's say for "apple" the recommendations have changed
 - apple: apple buy, apple big → apple pie, apple phone

You may need O(len(word)) operations. You update each prefix of word

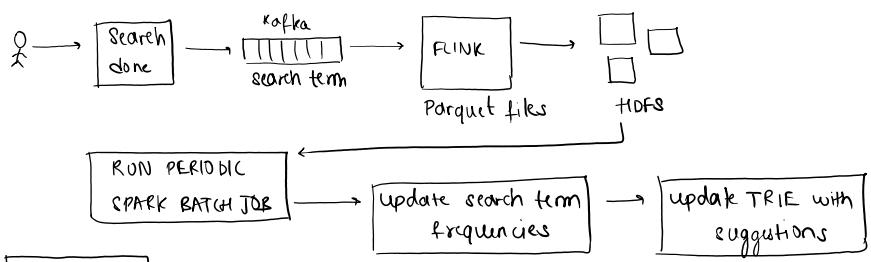
apple: "apple buy, apple big" → apple pie, apple phone "

You may need $O(\text{len}(\text{word}))$ operations. You update each prefix of word

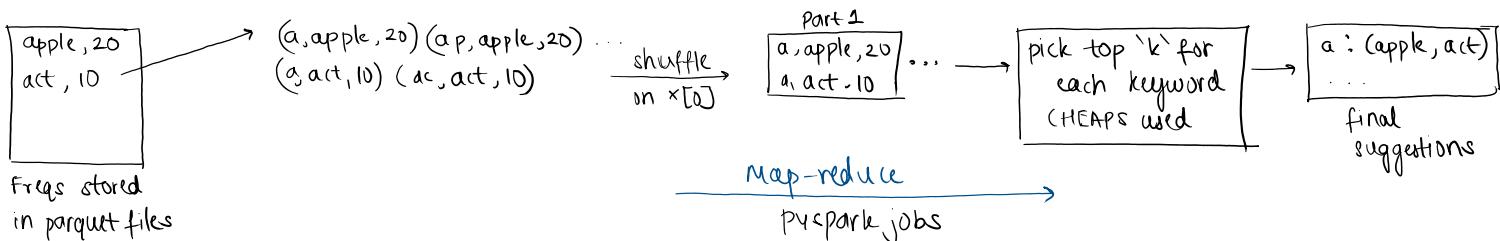


→ You might need to acquire locks for this. (In case u update 2 words with common prefixes)

UPDATING FREQUENTLY HIT TERMS ACQUISITION SERVICE



BATCH JOB



How users will READ

① API calls - BAD CHOICE

→ You will make 1 api call for each character being typed
Too much overhead due to headers + payload to fn

② Long polling (same issues as api over HTTP)

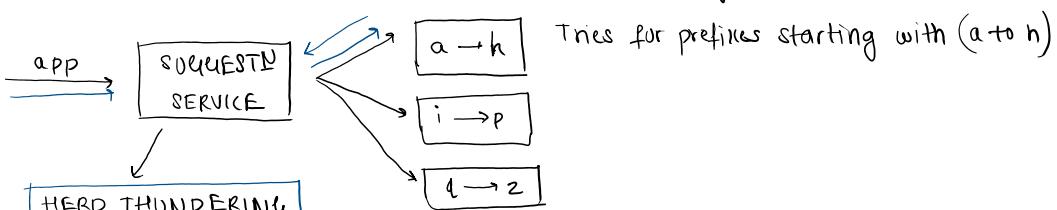
③ SSE (server sent events) → NO (you need bidirectional flow)

SOLUTION **WEB SOCKETS** → Bidirectional
No round of headers
You need a websocket manager
→ Zookeeper + consistent hashing + Load balancer

SCALING - DISTRIBUTED TRIE

→ You can't store all keywords in single Trie (TOO LARGE to fit on disk)

SOL: distribute tries across machines using range based partition.



HERD THUNDERING

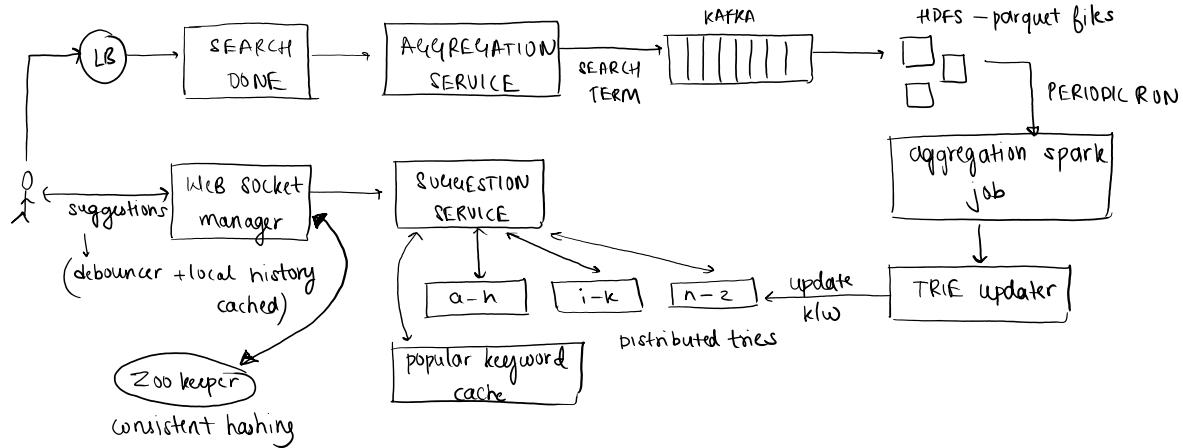
what if you get lot of requests at a time for certain keywords

① Dynamic Repartition

: Reduce or increase range partitions pre-emptively
→ will take some time to redistribute and recalculate
done in bg

- ① Dynamic Repartition : Reduce or increase range partitions pre-emptively
 → will take some time to redistribute and recalculate
 (done in bg)
- ② Cache popular keywords

ARCHITECTURE



REQUIREMENTS

- 1) Users can see who they follow + their followers instantly.
- 2) Low latency news feed for each user
- 3) Posts can have configurable privacy types
- 4) Users can comment on posts (Normal vs nested both)

capacity estimates READS >> WRITES

- 1) 100 b for each tweet text + 100b for metadata
- 2) 1 Billion daily posts : $1B \times 200b \times 365 = 7.3TB/\text{year}$
- 3) On avg user has 100 followers. 'Verified' users have millions of followers
- 4) Comments also are 200 b (100b content + 100b metadata)
- 5) Max comments in a popular post = 1M = $10^6 \times 200b = 200\text{ MB}$

FETCHING FOLLOWERS/FOLLOWING

- ① get_followers (user-id) → all ppl who follow user
- ② get_following (user-id) → all ppl user follows

FOLLOW (N)Y TABLE

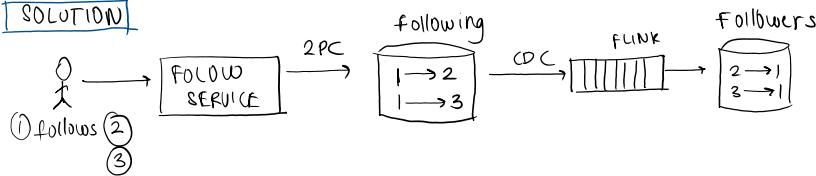
user_id	following
1	2
	3

FOLLOWER TABLE

user_id	follower
1	2
	3

→ get following (user) is efficient (indexed) → REVERSED
 → get followers (users) is inefficient

→ To make both efficient you need both tables
 BOTH writing to both tables in extra time

SOLUTION**WHICH DB?** → CASSANDRA

- Multiple concurrent writes happen.
- No need to handle write conflict/duplicates
- 1 follows 2] just deduplicate, nothing complex

→ Partition key: user_id, sort key: follower/following_id

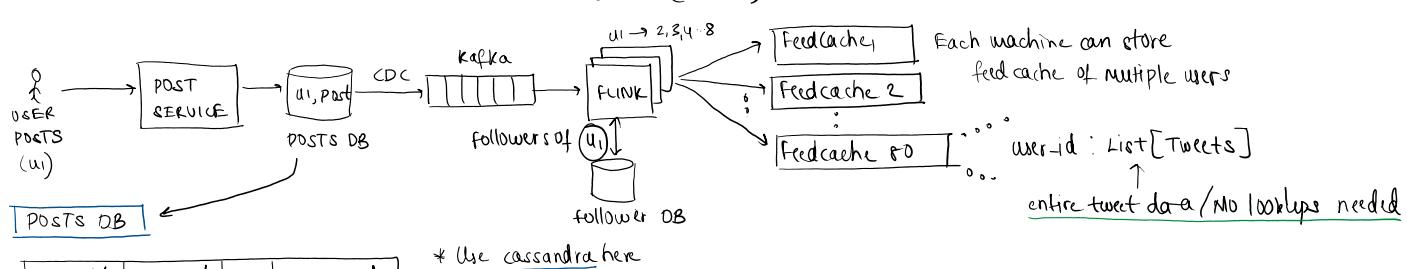
NEWS FEED**Naive + Bad**

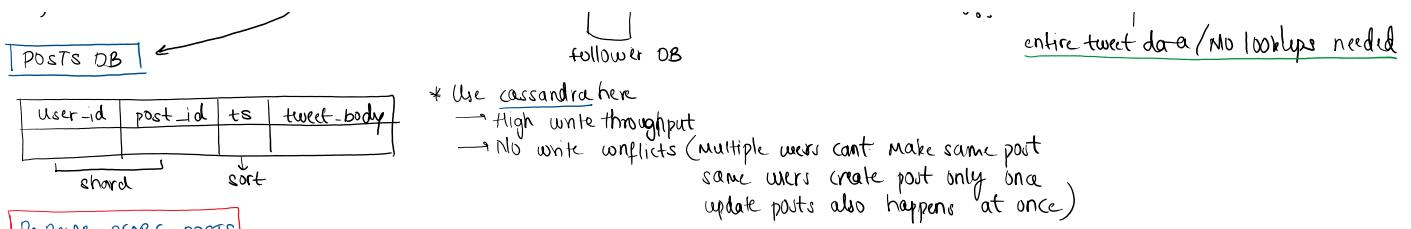
→ getFeed (user) → fetch followers → fetch all posts made by ppl followed

PRE COMPUTE FEEDS FOR EACH USER? → is this possible?

- Somehow distribute tweets to feeds they belong to
- Assume a tweet belongs to 100 feeds → reasonable
 $(1B \text{ tweets/day}) \cdot (200b \text{ each tweet}) \cdot (100 \text{ copies}) = 2 \cdot 10^{13} = 20\text{ TB tweets/day}$

20TB distributed into 256 4GB caches = 80 such machine caches (DOABLE)



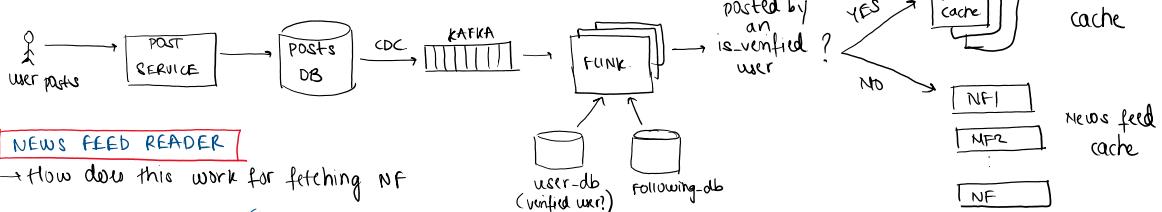


POPULAR PEOPLE POSTS

→ lets say a popular person posts (has millions of followers OR verified)

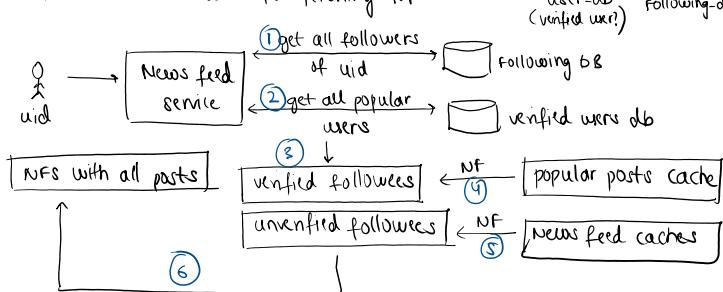
In this case you will copy that post to 1M News feed caches

→ very time taking → slow → post cache for popular posts
and more intensive



NEWS FEED READER

→ How does this work for fetching NF

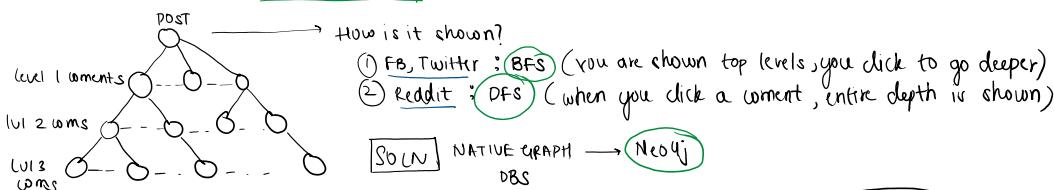


COMMENTS

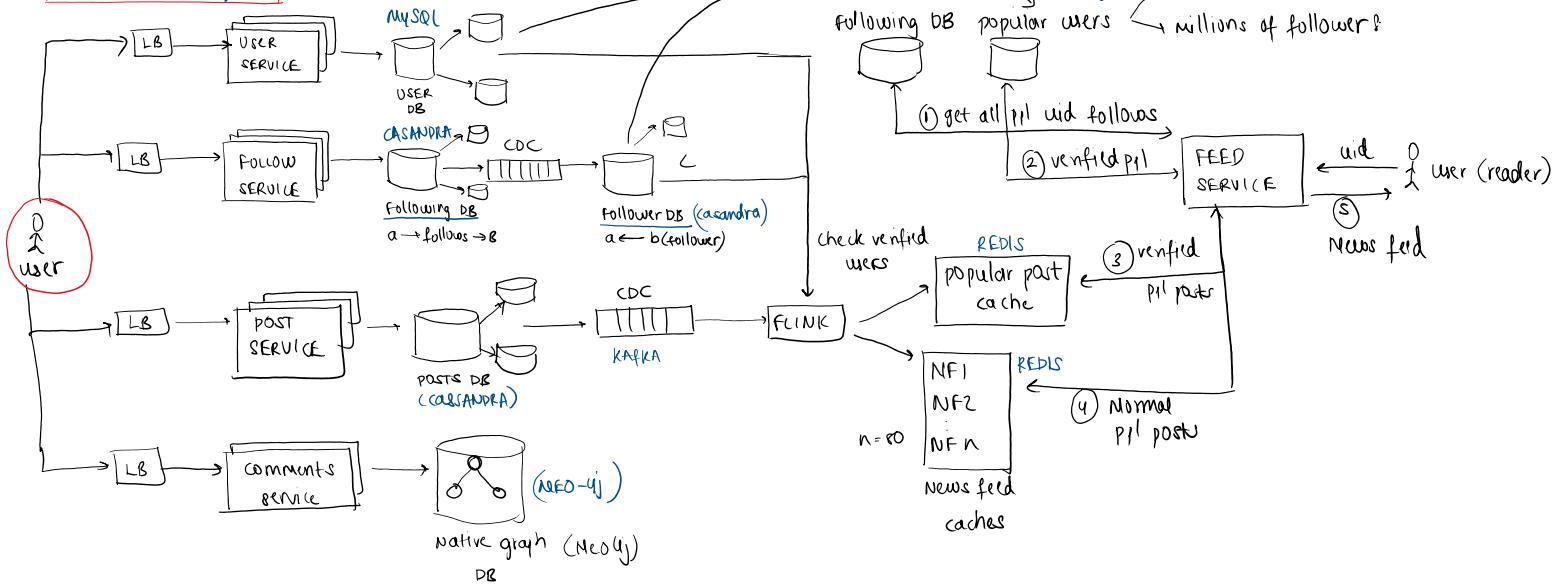
① Only one level → CASSANDRA

(high write throughput / no conflicts)

② Nested comments → GRAPH DB needed



ARCHITECTURE DIAGRAM



Whatsapp

28 December 2024 12:57

REQUIREMENTS

- Group chat with 10 members
- Sending + receiving messages in realtime
- Persist messages to read later (store conversations)

CAPACITY ESTIMATES

- 1 Billion users, each sends 100 messages/day
- Each message is $100\text{b} = (1 \times 10^9) \times (100) \times (100) = 10^{13}$ = $10 \text{ TB/day} \times 400 = 4 \text{ PB/year}$

CHAT MEMBERS DB

- Find all chats user is part of (left pane in whatsapp web)

user_id	chat_id
1	12
1	15
2	3
3	10

user_id → Index + shard

- multiple writes may happen: single leader replication (MySQL)
- can think of SET CRDT (too complex for this)

why MySQL? users don't frequent start chat/leave or join groups
→ performance impact significantly less

USERS DB

(user_id, email, pwd-hash, wr-metadata) ⇒ MySQL Infrequent change to this table
insert + partition

MESSAGES TABLE

partition

SERVER BASED TIMESTAMP

chat_id	ts	message	metadata	uid

which DB to choose?

- Reads and writes are both imp + frequent
- Column oriented storage? → usually you read only part of columns
(Metadata is usually not read)

(1) High writes → Cassandra? dealing with conflicts

(2) HBase → single leader (NO need of conflicts)

column based DB

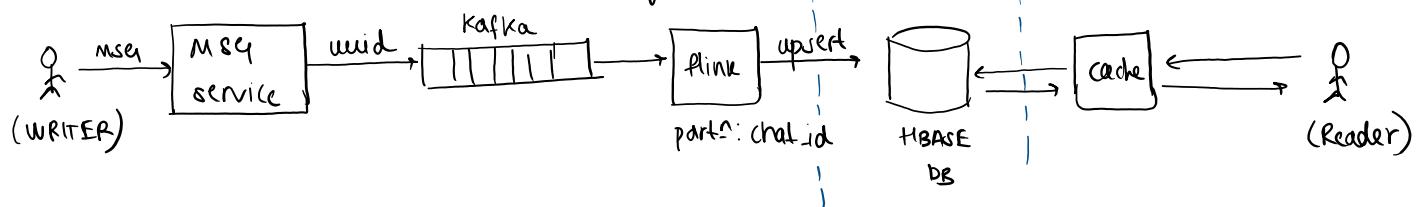
LSM Tree indexes → faster writes

T

column based DB
LSM Tree Indexes → faster writes

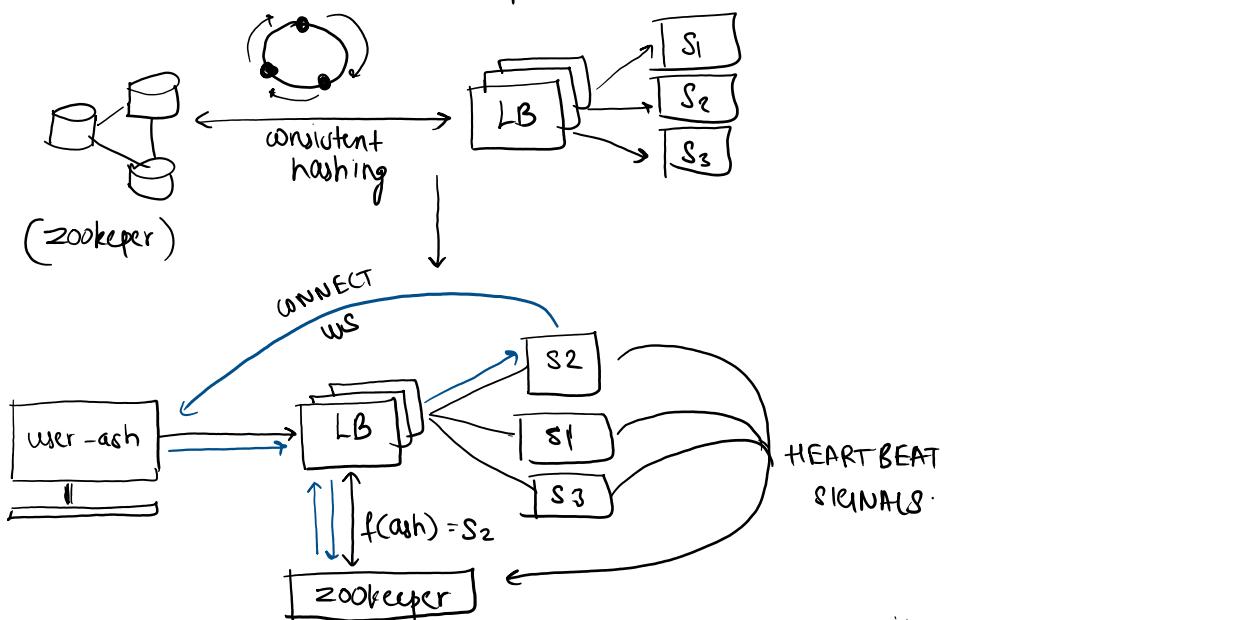
Tuning for faster writes

- ① Read directly from HBase
- ② Writes happen via stream processing (uid + msg handles idempotence)



LOAD BALANCING

- WEBSOCKETS (Allow bidirectional data transfer)
- Load Balance on user-id: requests from user must go to same websocket
 - ↳ will use zookeeper to keep track of active servers
 - ↳ internally use consistent hashing for dynamic mgmt
 - ↳ LB will listen to zookeeper



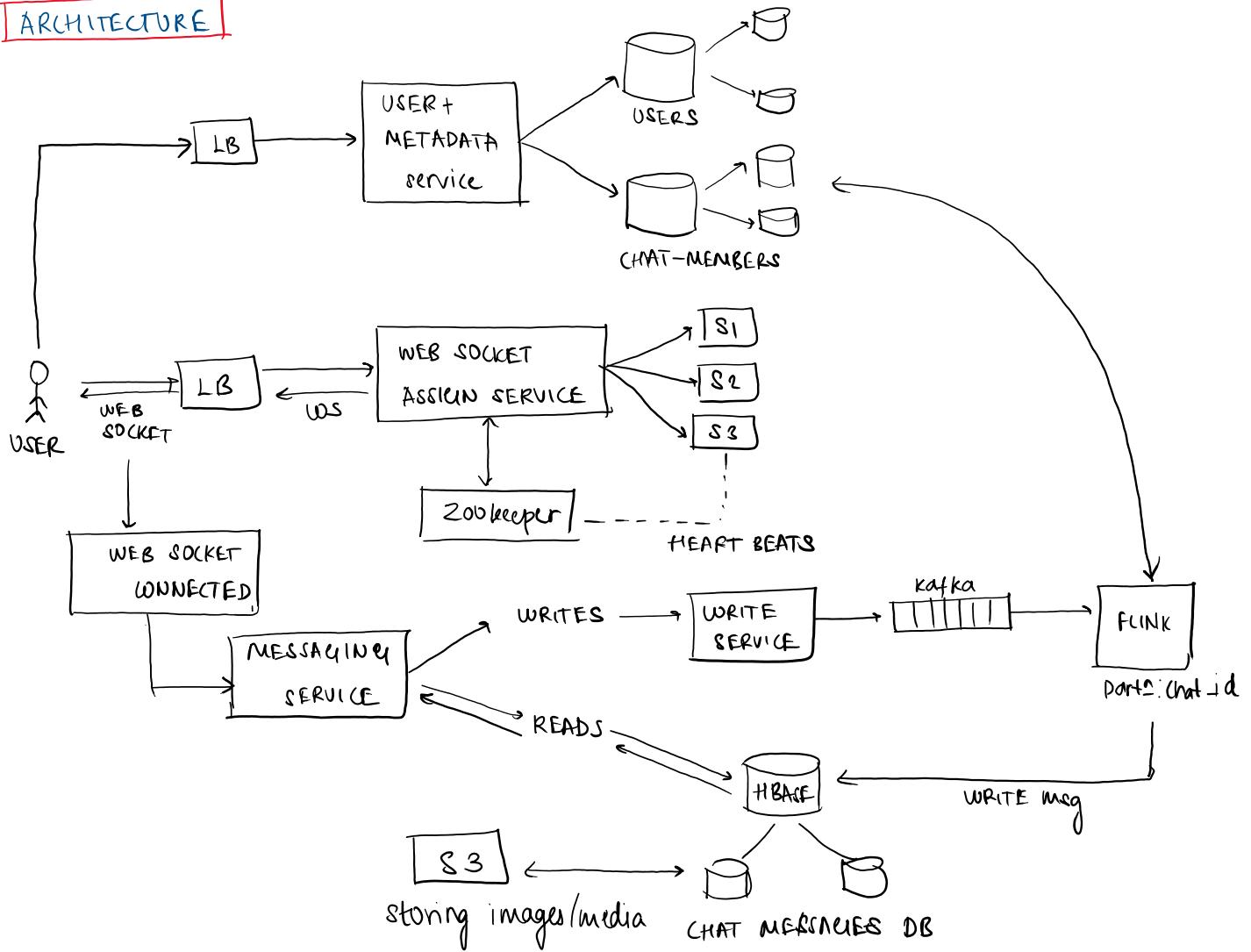
- ① user hits LB
- ② LB hits zookeeper with user-id
- ③ Zookeeper updates LB with which server to assign (S2)
 - All servers send heartbeat to Zookeeper at freq times
- ④ LB now connects client ← (S3)

IN CASE (S2) goes down

- LB asks zookeeper again (after some interval)
- Zookeeper knows (S2) is dead (didn't get heartbeat)
- It assigns another server to user (ash) using consistent hashing

it just has a map
user → server

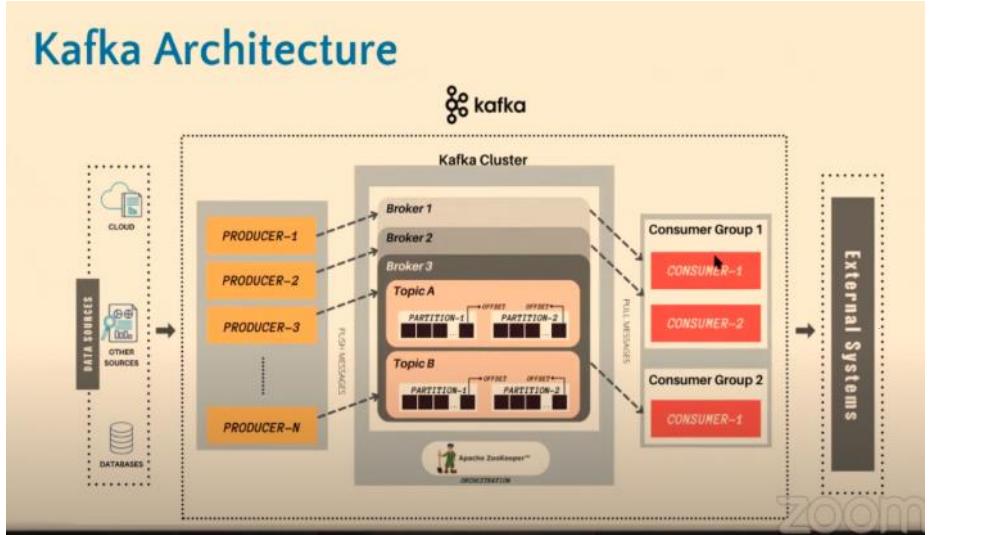
ARCHITECTURE



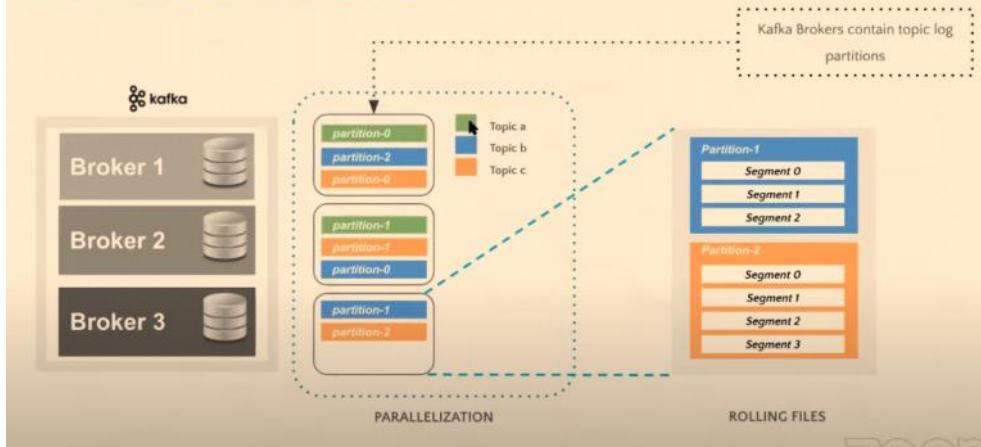
Kafka Architecture

08 April 2021 14:03

KAFKA ARCHITECTURE BASED



Partitions in Kafka



PRODUCER DETAILS

Every msg you push you need the following detail

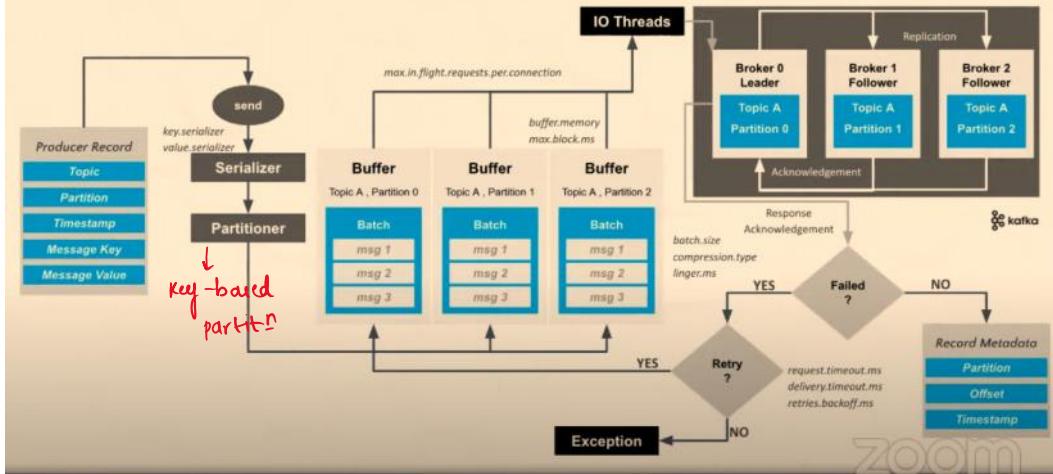
Is:

- a) key (optional)
- b) value
- c) topic
- d) group_id
- e) partition_num (optional)

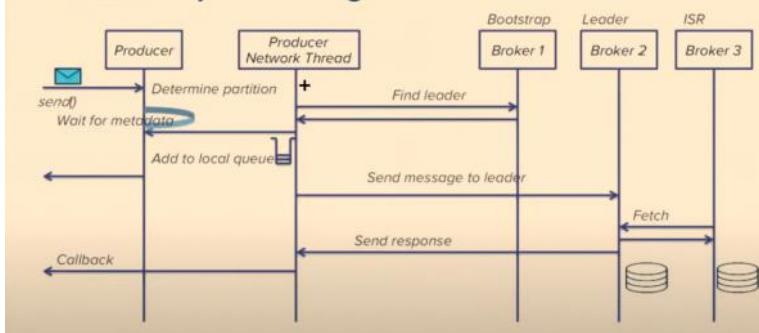
How to decide which partition to write to?

1. key is passed: Use $\text{hashCode}(\text{key}) \% \text{n_partitions}$
2. key is not passed: simple round robin
3. You explicitly mention which partition to write to
4. You define custom Hashing function

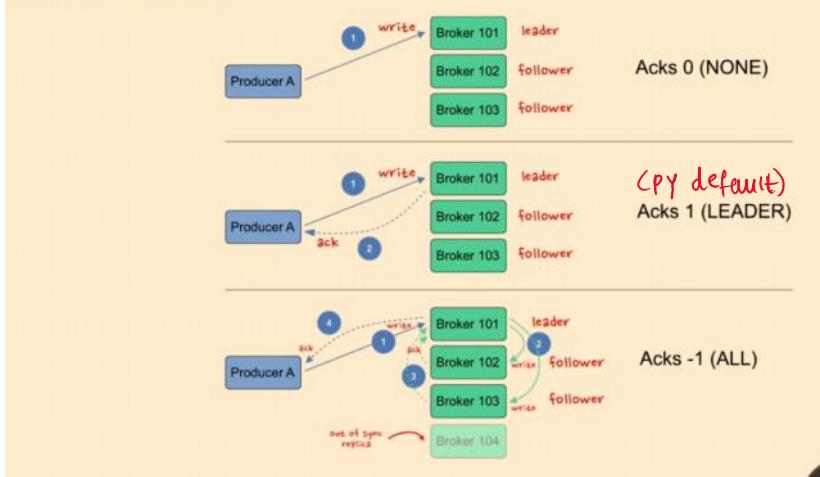
Producer Side Architecture



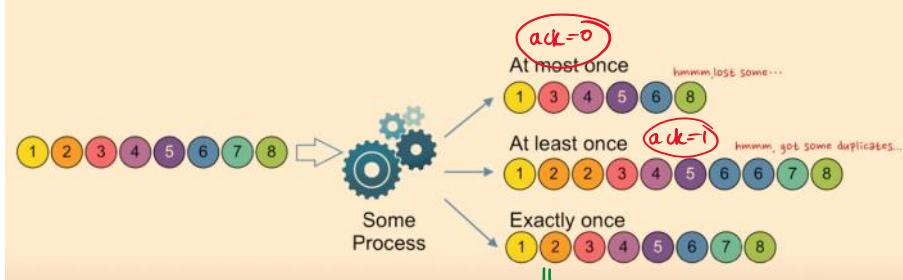
Producer sequence diagram



Producer Guarantees



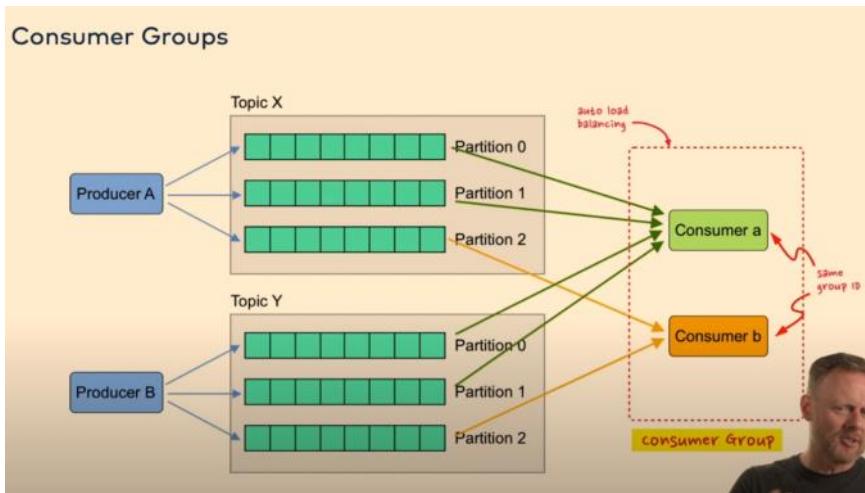
Delivery Guarantees



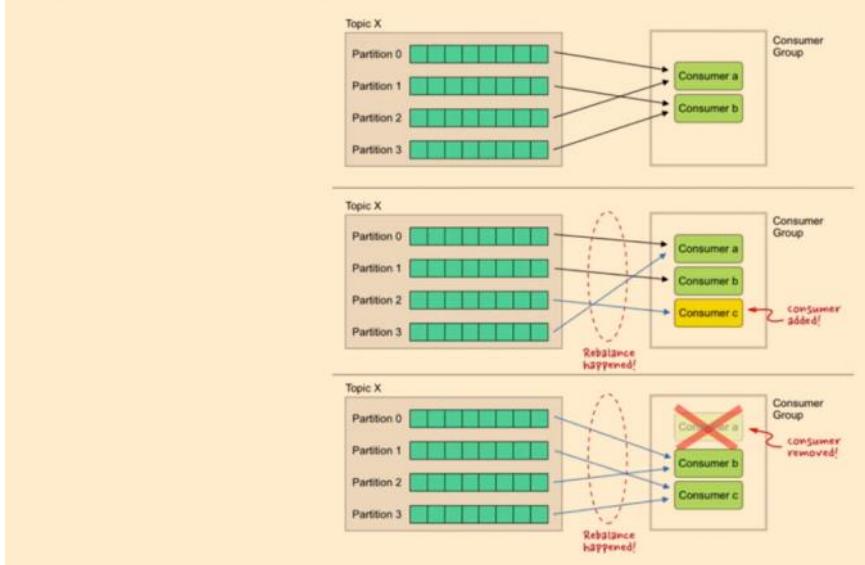
* Hard to implement

* You need to sync both (prod cons) → you need something as a transactional api server

CONSUMER DETAILS



Consumer Rebalances



OFFSETS

- What if your Consumer loses connection? or new consumer group gets added
- Kafka Offsets Committed vs Current vs Earliest vs Latest Differences | Kafka Interview Questions

Types of Offsets

- **Current offset:** Kafka maintains this offset of current message being read by the Consumer

Types of Offsets

- **Current offset:** Kafka maintains this offset of current message being read by the Consumer
- **Committed offset:** Consumer after successfully reading message will commit this (Auto in python) indicating that till "i" all messages have been read

Auto.offset.reset = "earliest | latest"

- earliest: start reading again from start of the partition
- latest: start reading from "latest_committed_offset" of last read

producer commit

→ Default: when u read it auto commits

→ LNC soln

- (1) Read msg (dont commit)
 - (2) process
 - ↳ write to op topic
 - (3) Now commit read to broker
- In case your pgm crashes you havent yet updated cur as read
Next consumer will pick it up again

$$7 - 2 = 1$$

DOCS

- Best Intro: [4. How Kafka Works | Apache Kafka Fundamentals](#)
- Confluent Playlist: [Getting Started with Apache Kafka®](#)
- Delivery Semantics: [Message Delivery Semantics in Distributed Systems: Kafka included](#)