

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** → 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

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
 - DB 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

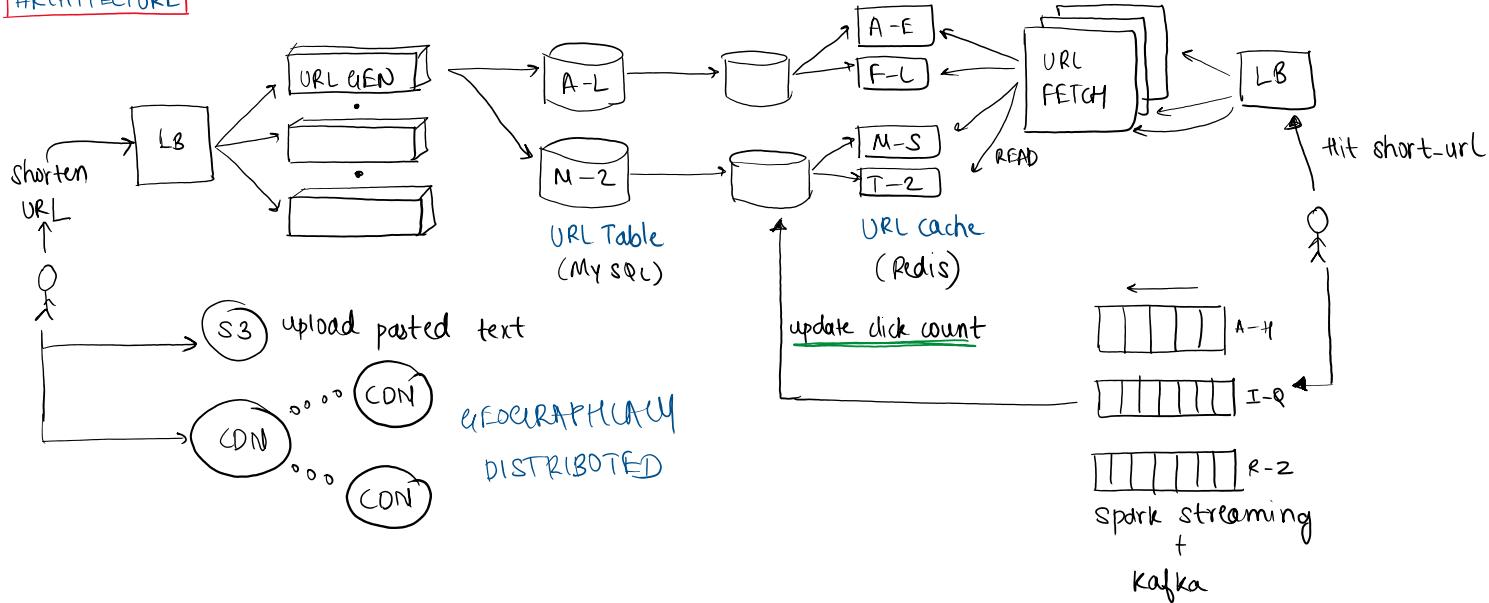
cheaper than HDFS / DBs

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

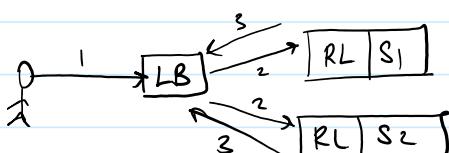
Unauthenticated service? → IP-address based

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

WHERE TO PLACE RL

Local to service

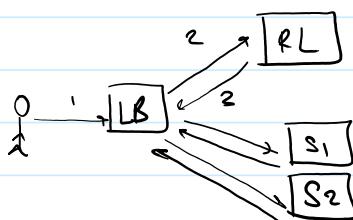
- + less components to manage
- + lesser n/w calls
- spammed requests may use up n/w bandwidth (payload size)
- tightly coupled with service



Distributed as service

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

THIS IS BETTER FOR US



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
→ (+) 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

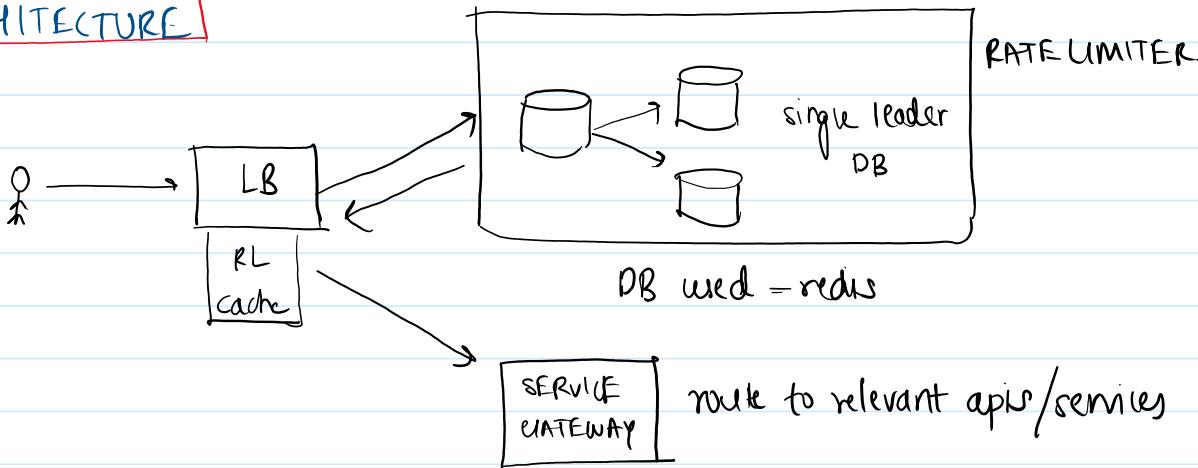
- + faster write throughput due to multi-parallel writes
- Counts may not be accurate since it needs to get propagated among replicas

single-leader write

- + Accurate counts since you write to a single leader
- slower writes since single write leader

WE USE THIS AS ALREADY IS MORE IMP

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

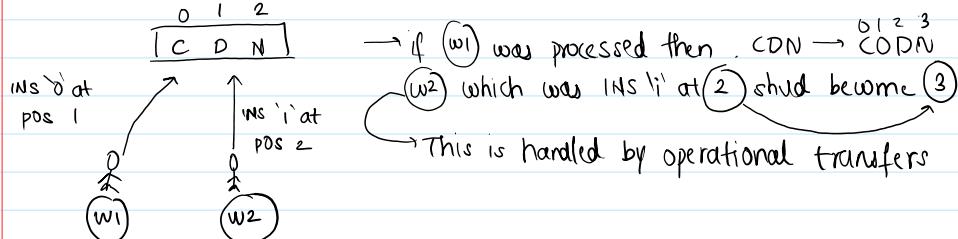
- ① using websockets
- propagate updates
 - propagate entire document
- These are sent by each user to a centralized server

CONS

- File size and num_users may be huge
- Propagate updates - you need to 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^t (INS \otimes 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^t 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 = $0 + 1/2 = 0.5 + \text{some random}$

Ex:

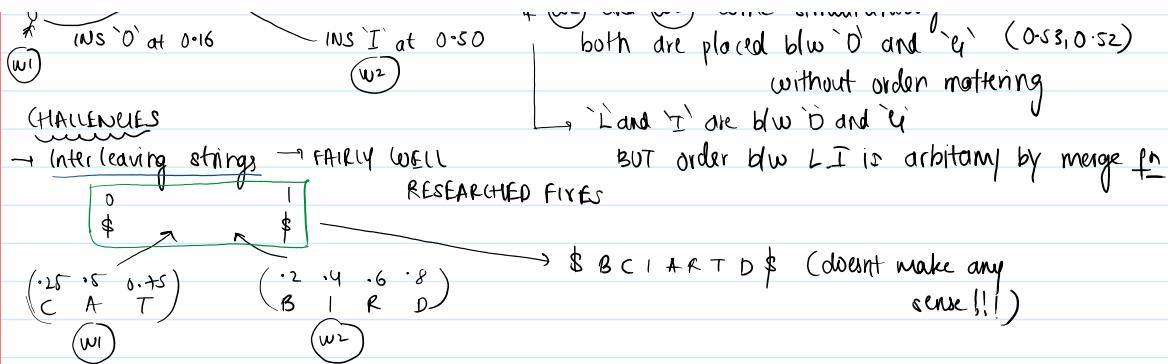
C	D	N	e
0	•33	•66	1

w_1 INS '0' at 0.16

w_2 INS 'I' at 0.50

jitter

→ Let's say w_3 inserts 'L' at 0.52
 w_2 inserts 'I' at 0.53
→ If w_2 and w_3 come simultaneously
both are placed b/w '0' and 'e' (0.33, 0.52)
without order mattering



Achieving Idempotence

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

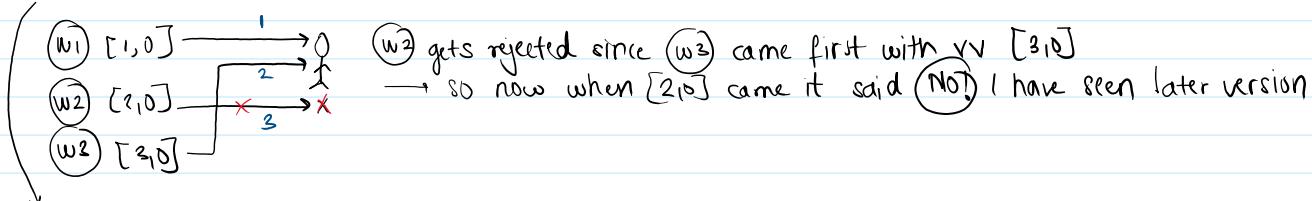
VERSION VECTORS

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

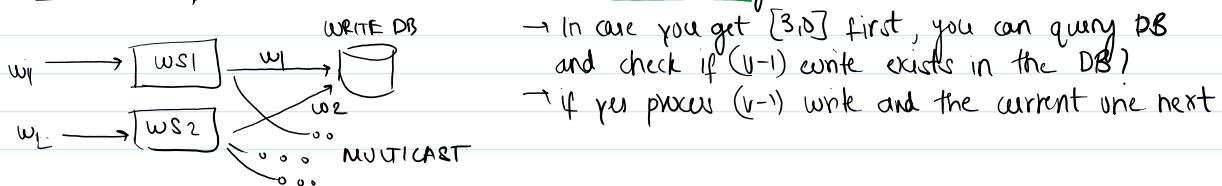
Ex - There are 2 write nodes/servers

- (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

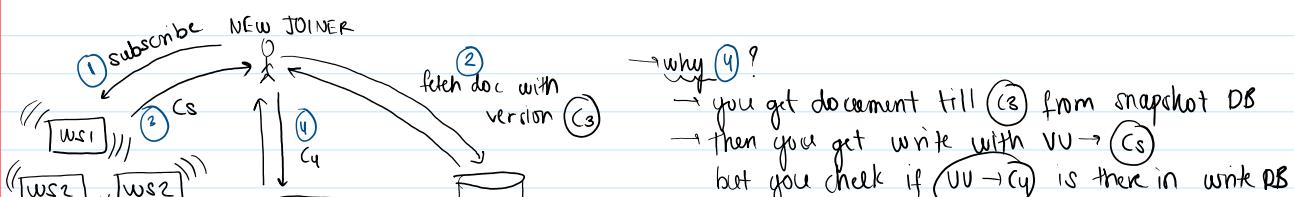
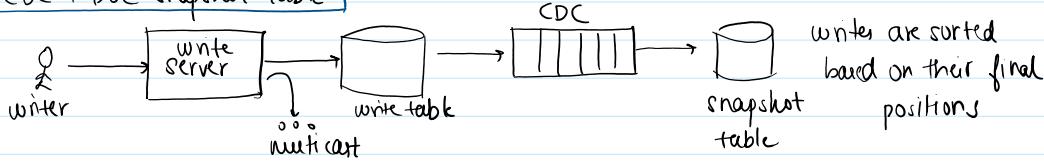


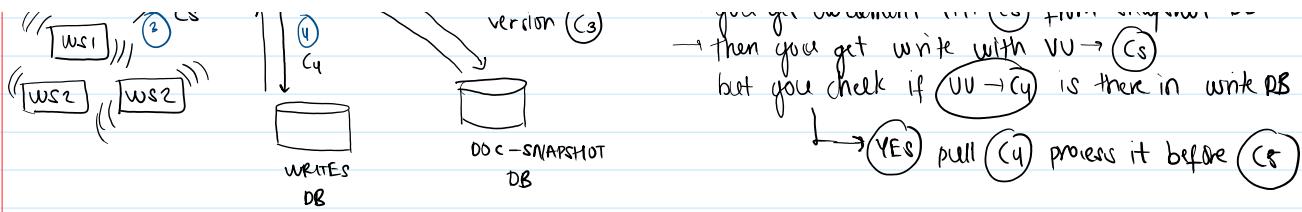
NEW JOINERS

↳ Basic: fetch the base document & apply writes

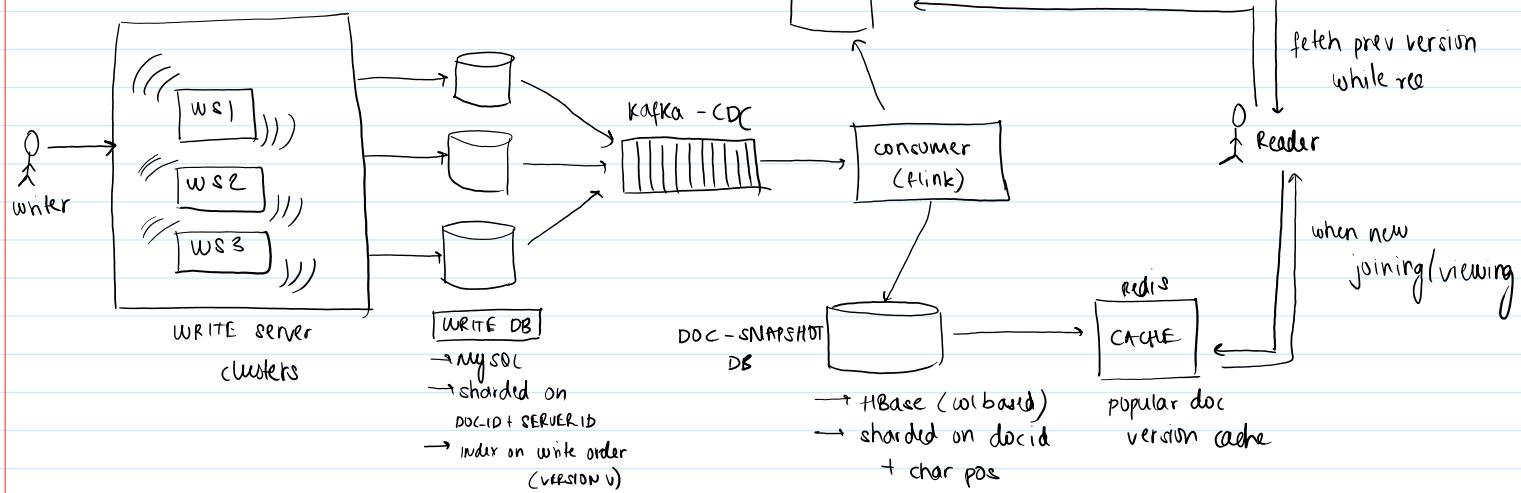
✓ BETTER

CDC + Doc snapshot table





ARCHITECTURE



REQUIREMENTS

- Users can post, watch, search videos
- 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

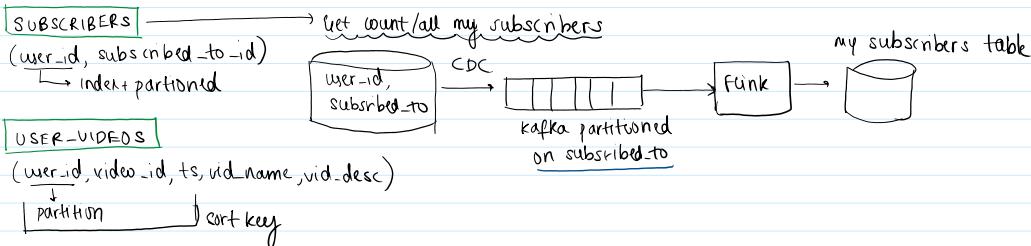
$$\begin{aligned} 100 \text{ MB} \times 1 \text{ M users} \times 400 \text{ days} &= 40 \text{ PB/year storage} \\ (100 \times 10^6) \times (10^6) \times 400 &= 4 \times 10^2 \times 10^6 = 4 \times 10^{16} \\ &= 40 \times 10^5 = 40 \text{ PB} \end{aligned}$$

VIDEO STREAMING

- ① Support multiple types of devices
 - ② Support multiple different n/w speeds
- ↳ 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 (Hd4+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

OUR CHOICE S3

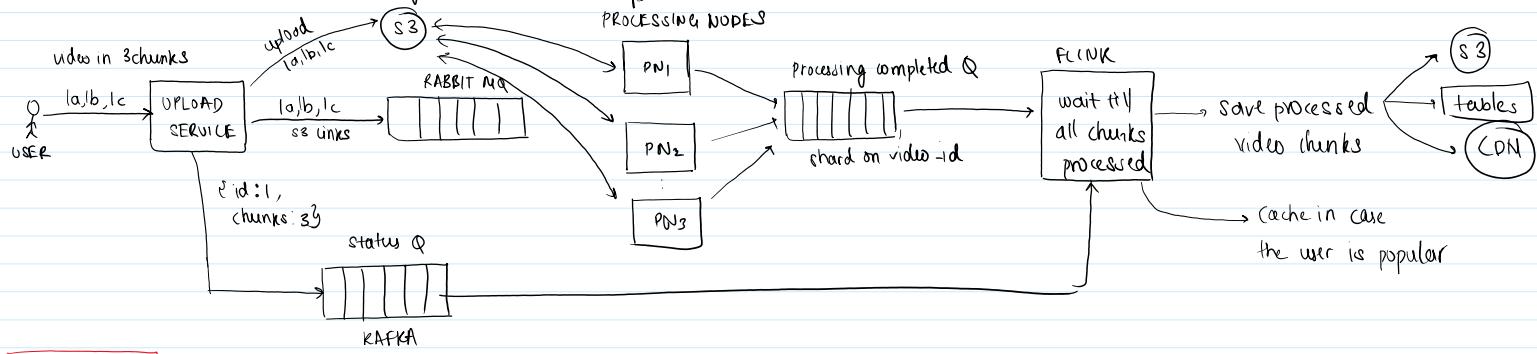
- + Cheaper storage costs
- + No CPU costs → just upload & forget

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

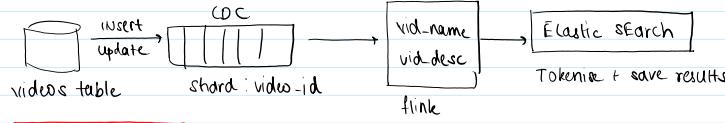
* we pick S3

→ cheaper in storage costs

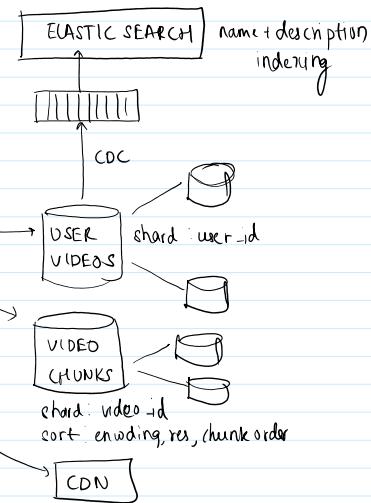
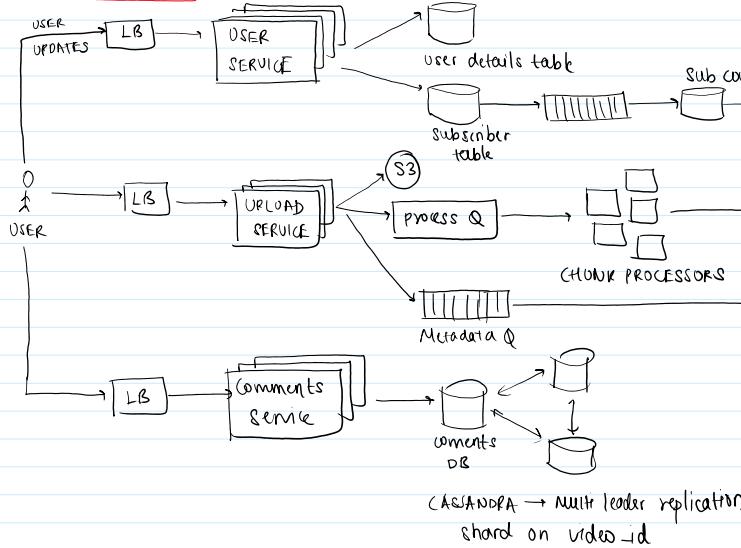
→ we are fine with NW call delays (CPU is more costly)



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 ⇒ storing counts of each = 4B bytes = 4 GB → can't be stored locally

→ Assume avg search term is 20 chars

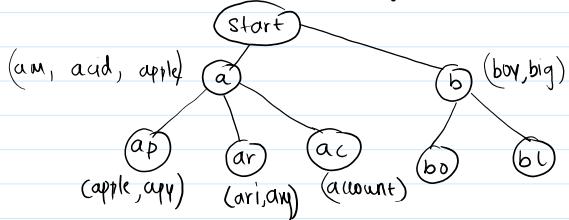
① 1B entries * 20 prefixes each = 20 B entries

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



TIME COMPLEXITY

① Insert word = $O(\text{len(word)})$

② Delete word

③ Search word

④ Search prefix = $O(P+k)$

len of
prefix

len of word
with prefix P

→ while typing TC for getting next

char suggestions is $O(1)$

"a" - am, acid, apple

"ap" = apple, apply

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

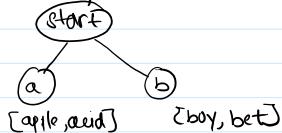
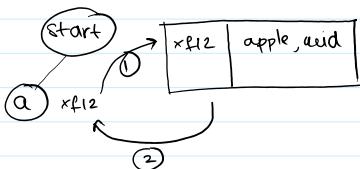
TIME EFFICIENT

→ store entire suggestions directly as list of words + node

→ No need to fetch other places to get data

+ LESS TIME

- MORE MEMORY



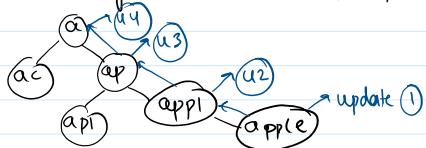
UPDATING KEYS

→ lets say for "apple" the recommendations have changed
(apple, apple buy, apple big → apple pie, apple phone)

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

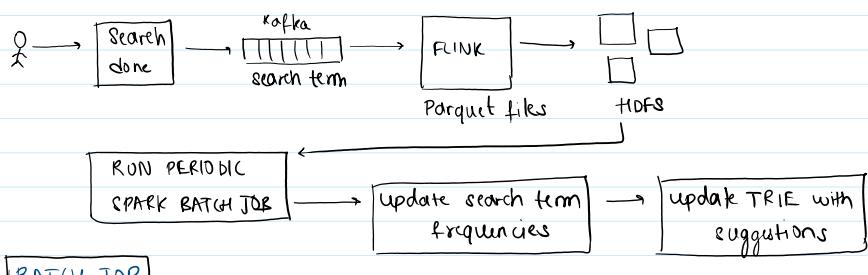
apple: "apple buy, apple big" → apple pie, apple phone "

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

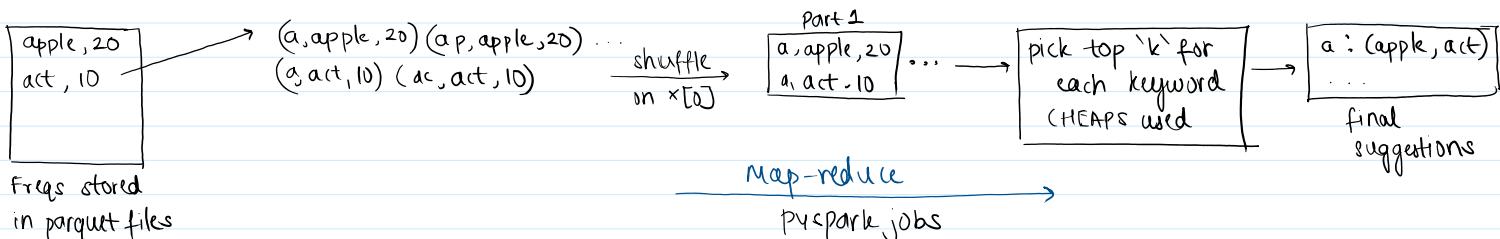


→ You might need to acquire locks for this. (In case u update 2 words with common prefixes)

UPDATING FREQUENTLY HIT TERMS | AGGREGATION 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)

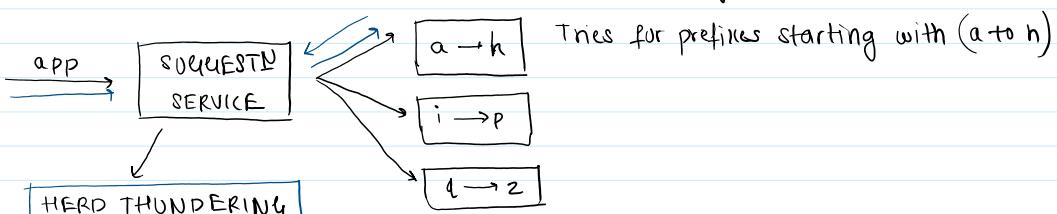
③ SSE (server sent events) → NO (you need bidirectional flow)

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 (Trie must fit on disk)

Soln: distribute tries across machines using range based partition.

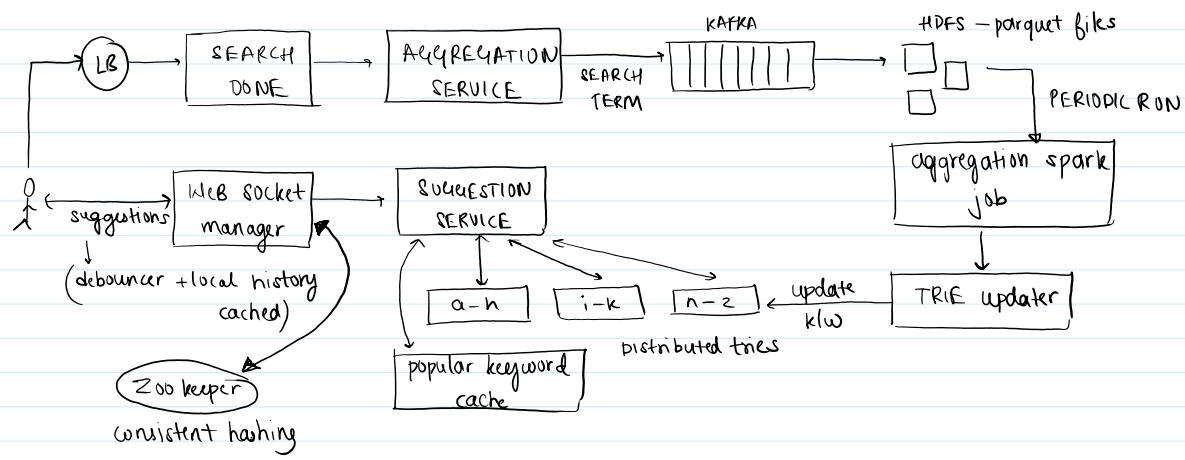


what if you get lot of requests at a time for certain keywords

→ 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 who user follows.

FOLLOW(MY) TABLE

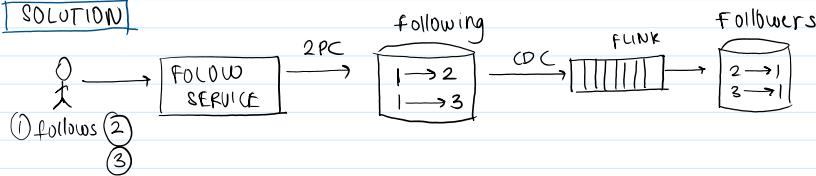
user_id	following
1	2
	3

FOLLOWER TABLE

user_id	follower
1	2
1	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

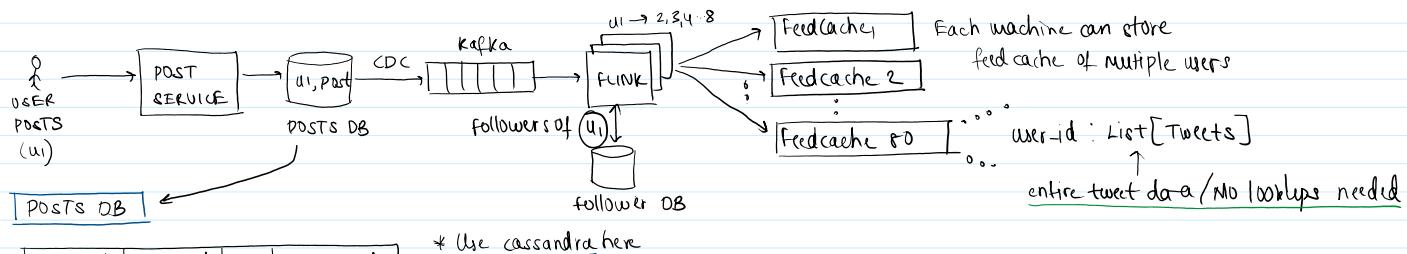
NEWS FEED**Naive + Bad**

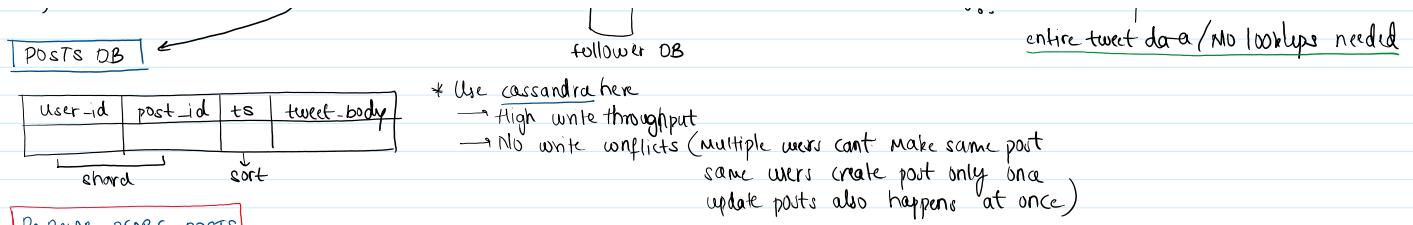
→ getFeed (user) → fetch followers → fetch all posts made by ppl followed

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)



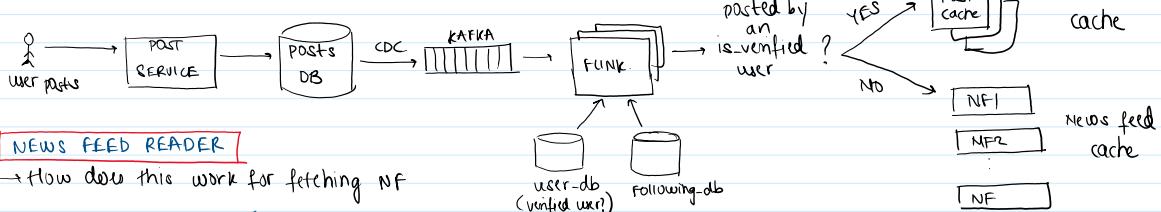


POPULAR PEOPLE POSTS

→ let's say a popular person posts (has millions of followers OR verified)

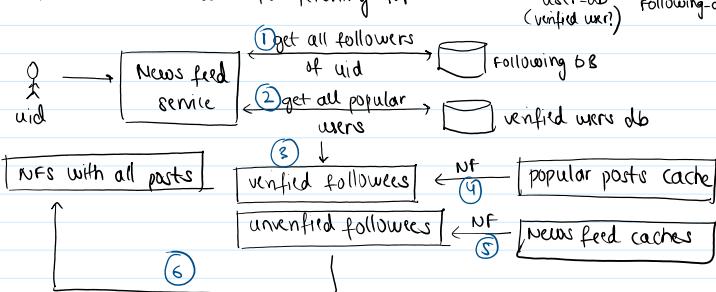
In this case you will copy that post to 1M News feed caches

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



NEWS FEED READER

→ How does this work for fetching NF

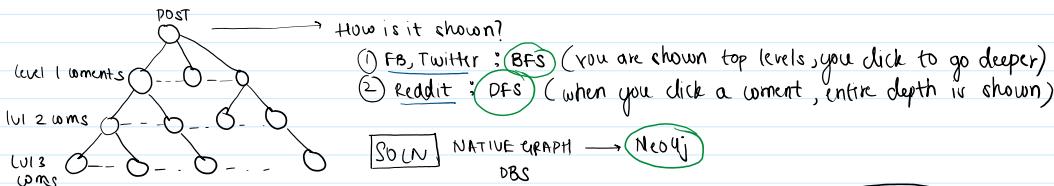


COMMENTS

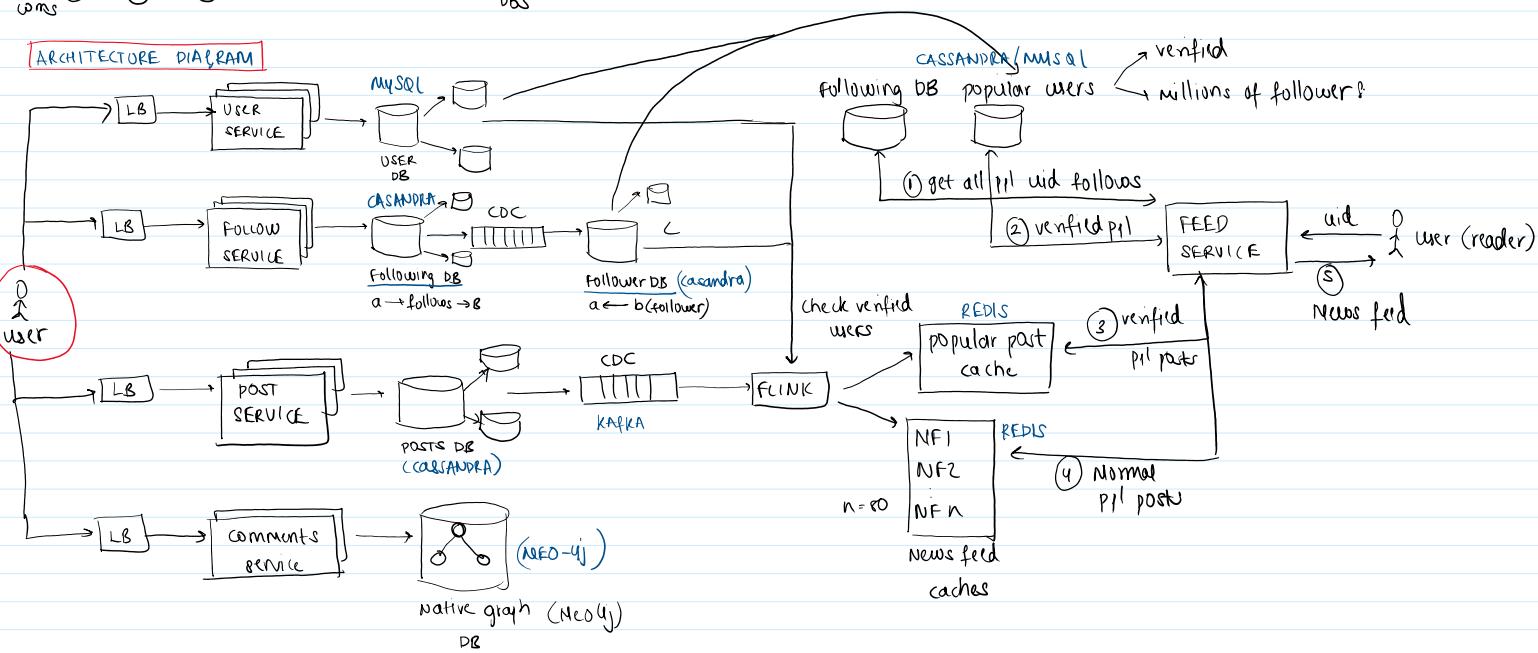
(1) only one level → CASSANDRA

(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 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 changes 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

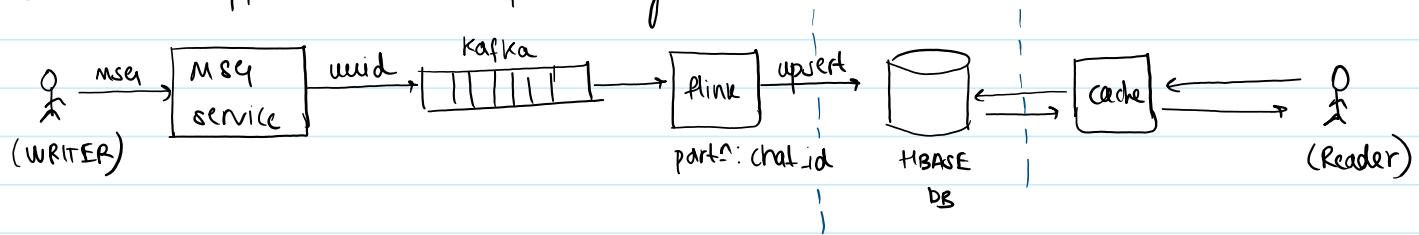
LSM Tree indexes → faster writes

T

column based DB
LSM Tree Indexes → faster writes

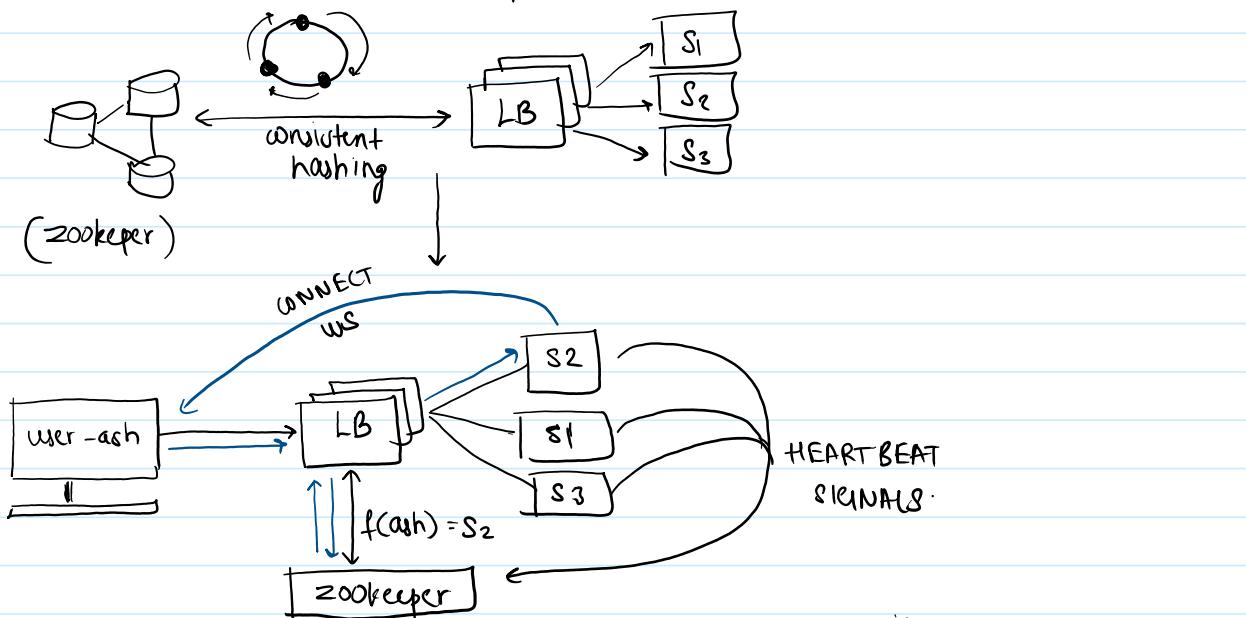
Tuning for faster writes

- 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 of 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 (S3) using consistent hashing

ARCHITECTURE

