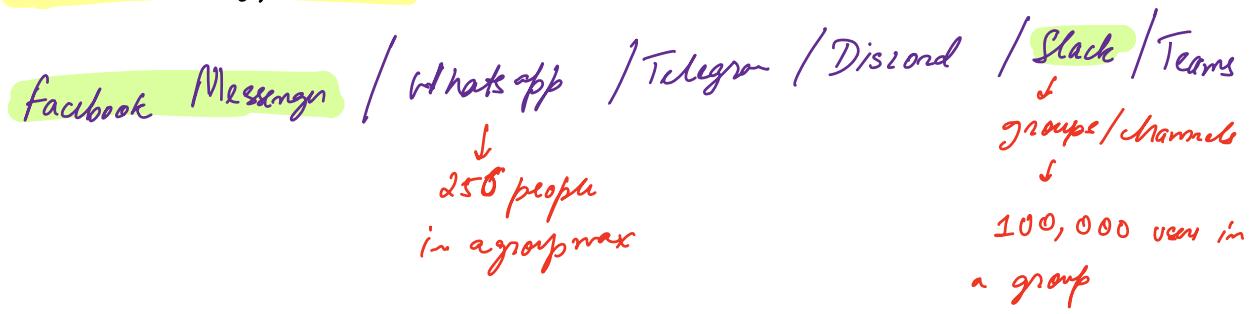


Agenda

- Problem Statement
 - discuss core features
 - do Not go on a fancy feature suggestion spree
- Minimal Viable Product (MVP)
 - Storage Requirements
 - # queries / sec
 - Sharding
 - Read / write heavy
- Scale Estimation
- Design Goals
 - Availability / Consistency
 - Eventual Consistency
 - Immediate
 - Consistency / Latency
 - Data Loss
 - no-consistency at all
 - Latency requirements
- API
- System Design
 - How is the system being used by external world
 - client / other services

Problem Statement



Minimal Viable Product ($MVP \rightarrow v0$)

- User registration/login → separate service
- Sender can send message to receiver
- attachments (Image / audio / video --) $v1$
 - discuss during initial
- online / offline $v1$
 - during class
∴ learning.
- read receipts / sent receipts $v1$
- group conversations $v2$
- recent conversations
- message history
 - ↳ all messages should be permanently stored on server side
- real time chat
- end-to-end encryption $v3$
- notification → separate service ($v1$)

Scale Estimation

QPS

— estimate yourself / ask interview

~~FB people is world~~ → ~~4B have access to FB~~ → ~~50% on DAU~~

~2B MAU at facebook → 50% on DAU → 1B DAU

each user sends 10 msg/day on average → 10¹⁰ message/day

Story

9 million messages on 1st day → Hackathon Message Service

On day 1 we expect to get ~10 billion messages / day

Facebook deals with > 70B message/day.

~~QPS~~
queries/sec

$$\frac{70 \text{ B}}{86,400} = \approx 1 \frac{7 \times 10^{10}}{86 \times 10^4} = 10^6 \text{ m/s} \Rightarrow 1 \text{ M msg/sec}$$

Storage requirements

Size of 1 message (on avg)

:+1:
:smile:
<3

content ≈ 100 letters (100 B)

timestamp = 8 B

receive-id = 8 B

sender-id = 8 B

attachment-path ≈ 50 B

→ large files/blobs → never stored
in DB (SS)

Total \approx 200B (size of each message)

70K messages / day

$$\begin{aligned}\text{Daily Storage requirem.} &\Rightarrow 200 \text{ bytes} * 70K \\ &= 14,000 \text{ B bytes} \\ &= 14 \text{ TB / day}\end{aligned}$$

$$\begin{aligned}20 \text{ years of data} &= 20 * 365 * 14 \text{ TB} \\ &\approx 100,000 \text{ TB} \\ &\approx 100 \text{ PB if data is total}\end{aligned}$$

Sharing

oh yes!!

Read/write heavy

every message is read

Read heavy

Both read & write heavy

Write heavy

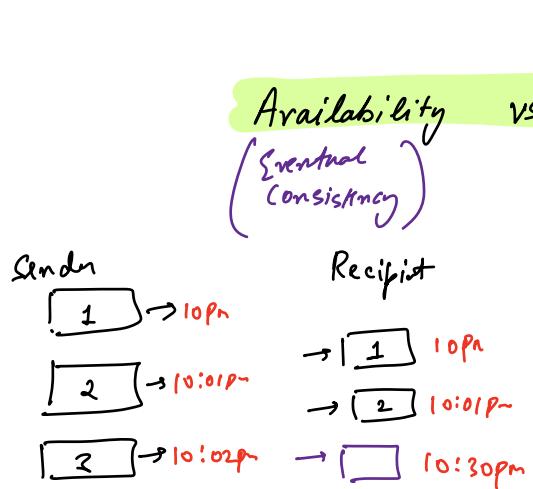
every message sent \rightarrow stored

Read heavy \rightarrow caching

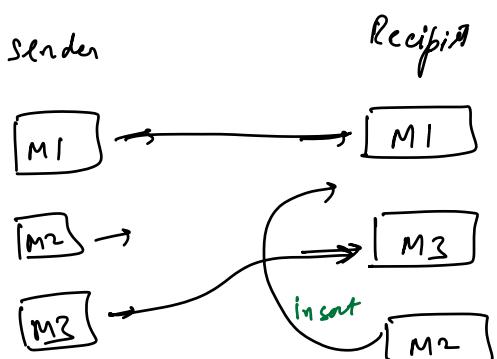
Write heavy \rightarrow Serialized data

\because Consistency was not important

Design Goals



messages can get delayed
Sender sees that M3 was
successfully sent at 10:02
but in reality it is
sent at 10:30



message received out of
order

Consistency → Communication on
billion of human
civilization

if a message is shown as
sent → then it must
actually be sent

Designing Data Intensive

Applications → Maotzi Klippen

Consistency vs Latency → Consistency

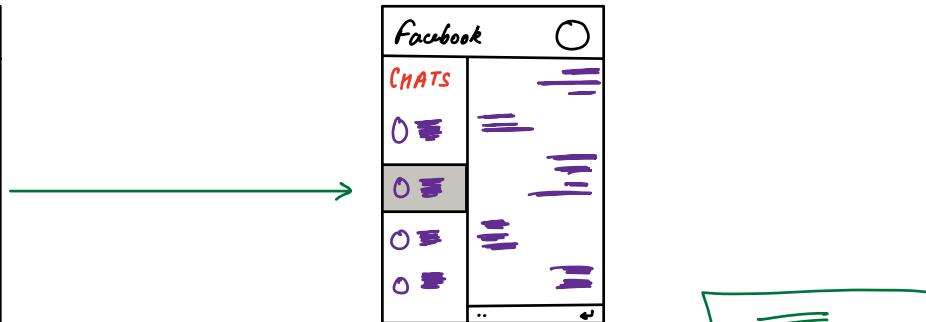
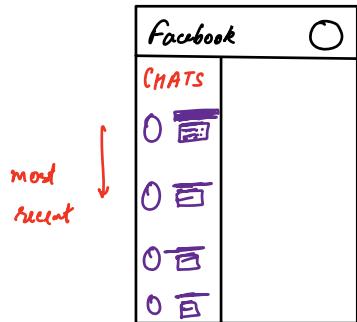
Data Loss → No!

Latency → pretty low

typical $< 50\text{ ms}$
100x
messenger $< 5\text{ sec}$

API → how will external world (client/other services) use our system

open APP

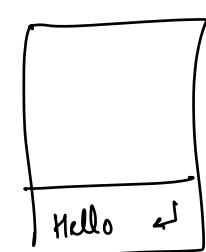


→ get Conversations (user-id, offset, limit)
sorted by timestamp
Pagination

- get Messages (conversation-id, offset, limit)
(convid, friend-id)
(group-id)

→ Send Message (conversation-id, message, void)
(convid, friend-id), message, void
(group-id)

Idempotence → $f(x) = f(f(x)) \Rightarrow$ idempotent function



{ sender-id : Akash
Rec id : Manya
Text : Hello }

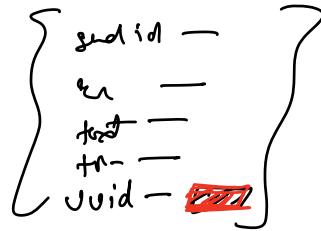
✓ stored

success

network issue

resent 'Hello'

We want to prevent two same messages from being stored twice, but still allow multiple messages with same text if user wants it



10:34 → 10:40

System Design

Sharding

only work for 1-1 conversations

userid → all conversations of 1 user must be present on 1 shard
any req for user will only read from 1 shard

get Conversations (userid, offset, limit) → easy → ∵ all the conversations of this user are in 1 shard

get Messages (conversation-id, offset, limit) → easy ∵ all conversations for this user

- ↳ (sender, receiver)
- ↳ (group-id)
- ↳ friend
- ↳ group

Send Message (conversation-id, message, uid)

- ↳ (sender, receiver) → store in both sends shard & receiver shard. (consistency)
- ↳ (group-id) → fan out write
| group | = 10,000 users → 10,000 writes
super expensive X

idea → each group behaves like a user.

UI → $\xrightarrow{\text{send}} \underline{G1}$ → another in UI shard
& UI shard

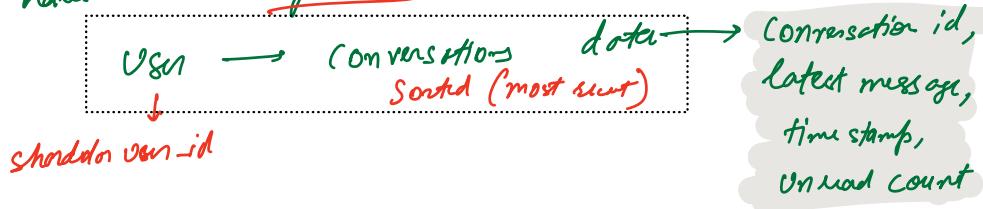
read?

Conversation id → (sender, receiver) 1:1
→ (group-id) groups

get Conversations (user id, offset, limit) → difficult
∴ to get conversation of a user we have to hit

diff shards

Alternate have another separate db that stores



get Messages (conversation-id, offset, limit) → easy
just hit that shard & get messages from the cor

Send Message (conversation-id, message, uid) → simple - just write to the conversation shard



secondary db stores (user id → conversation)
(sent by)

however if we have

when a message is sent to conversation → go to secondary db.

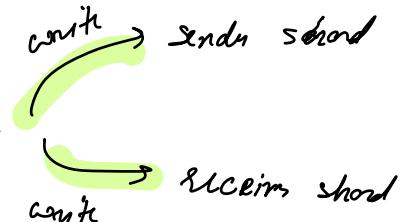
for each user that
is part of the conversation, mark the convos as most recent.

X expensive → apps like Slack → don't sort
most recent conversation.

Slack still has secondary db
to store user-id → list of conversations
↳ not sorted by
recency

Send Message Consistency

Handling by user-id → send message

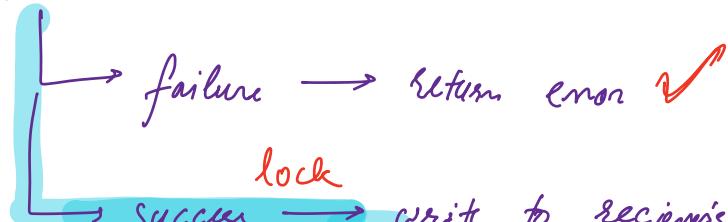


Problem → what to do if 1 write succeeds but
other fails

Sender → receiver
Sent X not recd

Case 1

first unit to sender's shard



latency is low

good case

Immediate retry
3 times ✓

~~Retry 10 min later~~

Eventual consistency
✗ not acceptable

DB is consistent ✓

Rollback &
return error ✓

Case 2

Sender → Receiver

find unit in receiver's shard

fail → ✓ return error

succeeds → write to sender shard

success ✓

fails ↗ retry with delay

low latency!

note → when sender refreshes app (browser)
then sender will not sustain sent
message

↳ still get info from app cache ✓

choice of database

→ both read & write heavy.

↳ requires a lot of sharding.

↳ NoSQL

No DB in the world optimized for both reads &
writes

Concur to either read heavy or write heavy.

↓
needed ∵ we need
immediate consistency

→ ensure that most reads never go to db

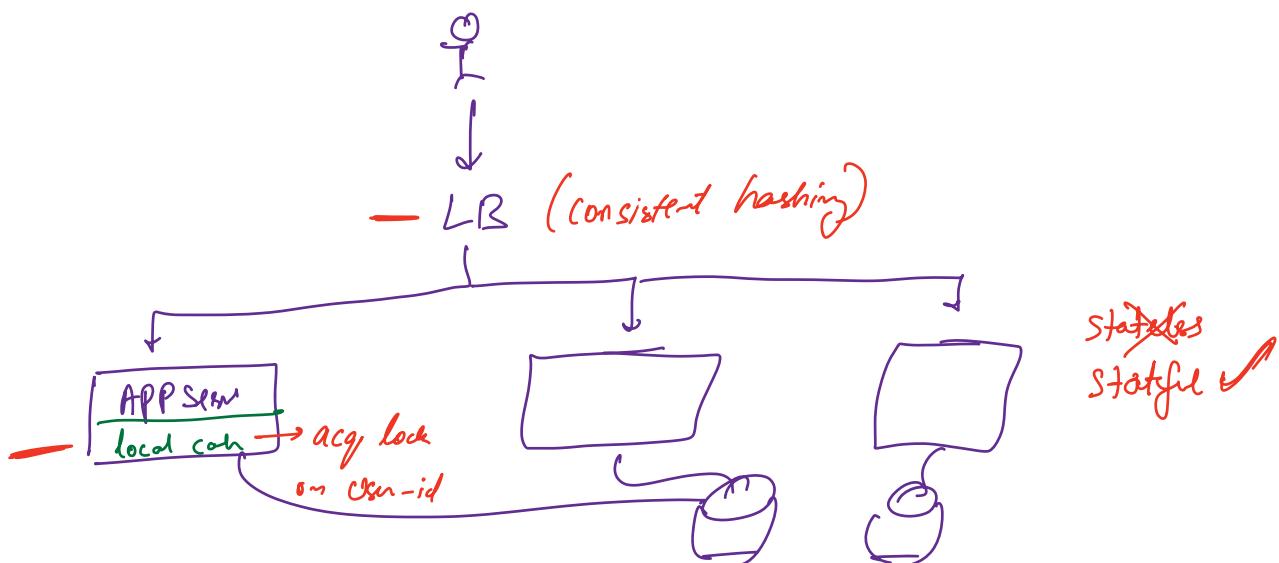
↳ heavy caching. → Sharded

↳ Cache is consistent with DB

↳ write-through cache

→ make sure that we handle high
concurrency in cache.

Don't use global cache → local caching.



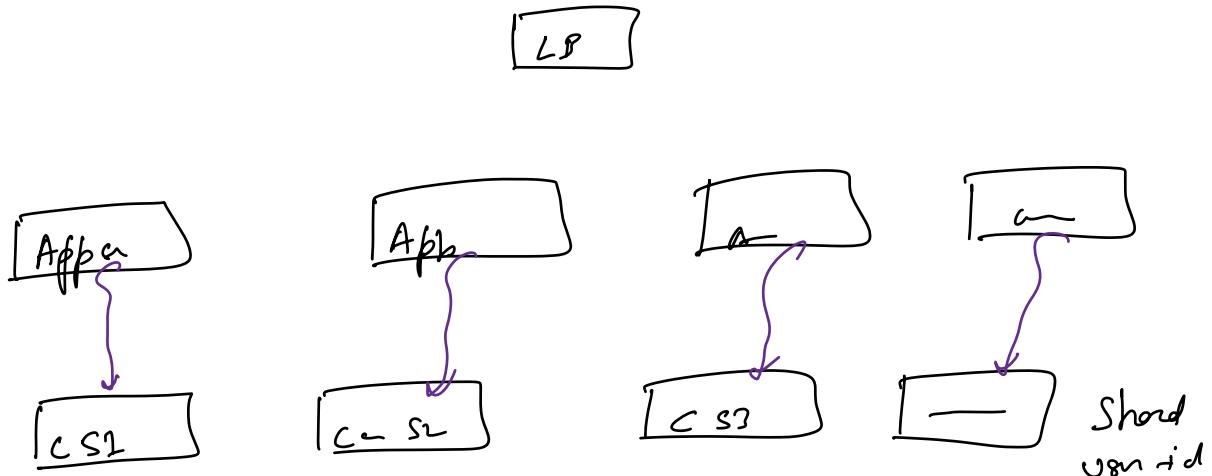
- Pros → ① Scale horizontally
② Race condition on handled-

Cons → ① If app session goes down. → Until user
is assigned a new session
Slave will be unavailable
to that user

{ } cold start of cache for new
app session.

Column Family \rightarrow ∵ Select messages / common

Write optimized \rightarrow Hbase



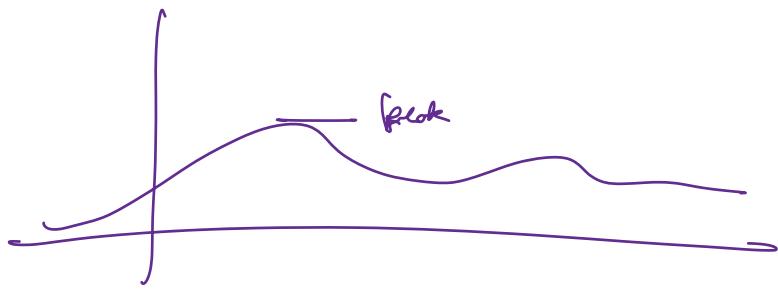
1 Million messages / sec $\xrightarrow{\text{avg traffic}}$ Write $\rightarrow 1$
Read $\rightarrow 3$

1 app can handle

$\xrightarrow{}$ 1000 writes/sec
 $\xrightarrow{}$ 1000 reads/sec

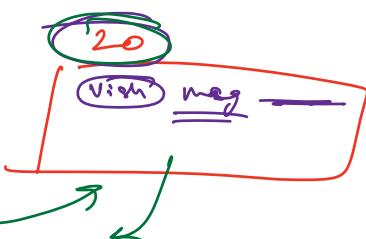
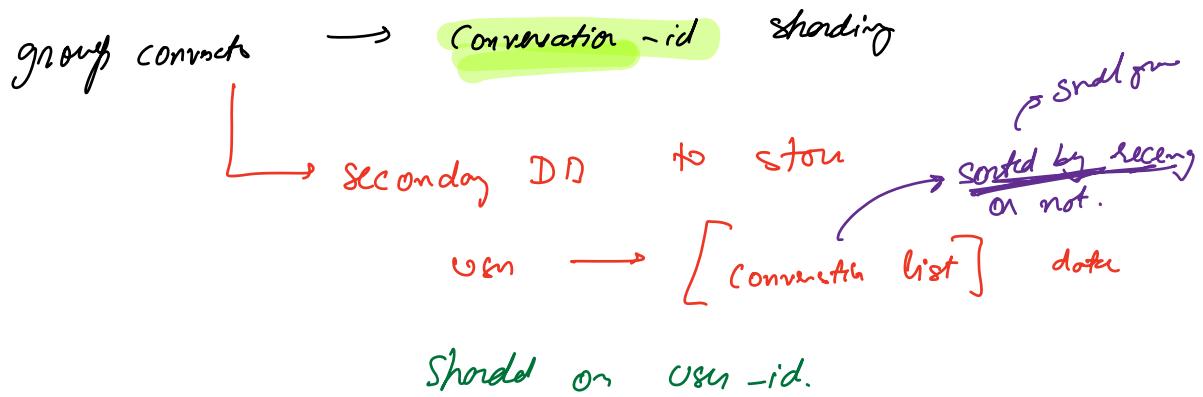
(If cache available
is for a diff user)

~ 1000 app servers.



WebSockets → persistent, two way connection
 HTTP/HTTPs

1:1 connection → osnid sharding



$$R + \omega \approx X$$

$\uparrow \uparrow$

Cache \rightarrow 10 year old msg ∇

\hookrightarrow last 10 days of data

$$70.1 \text{ msg/day} \times 200 \text{ R} \times 10$$

$$= 1.4 \times 10^5 \text{ R bytes}$$

$$= \underline{\underline{120 \text{ TB}}} \text{ data in cache}$$

Idempotence

$$\text{for all } x \quad f(x) = f(f(x))$$

if you apply the function multiple times
the result doesn't change

\neq fixed point
means some x

$$\text{abs}(x)$$

$$\text{abs}(-10) = 10$$

$$\text{abs}(\text{abs}(-10)) = 10$$

$$\begin{array}{l} \text{floor}(x) \\ \text{ceil}(x) \end{array}$$

$$\begin{array}{l} \left(\frac{d}{dx} \right) X \\ \left(\frac{d}{dx} c^x \right) \end{array}$$

API Idempotent

GET (read)

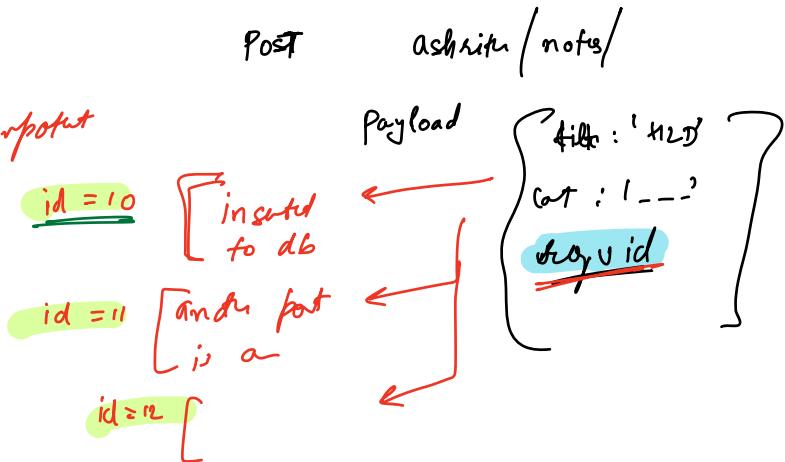
GET /ashrik/messages/
← msg

GET /ashrik/messages/
← no msg

GET /ashrik/messages/

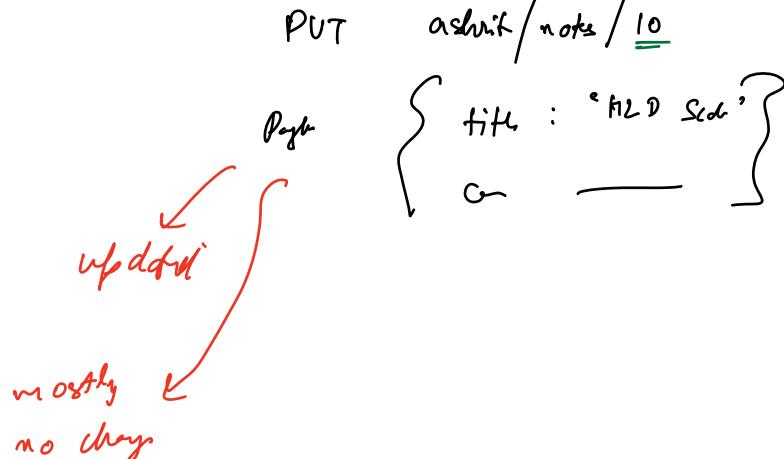
POST → insert

↳ usually not idempotent

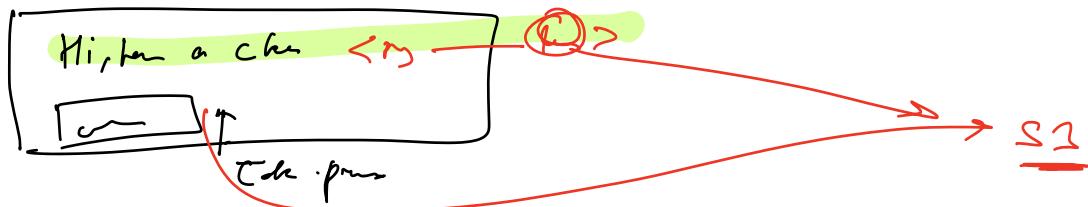
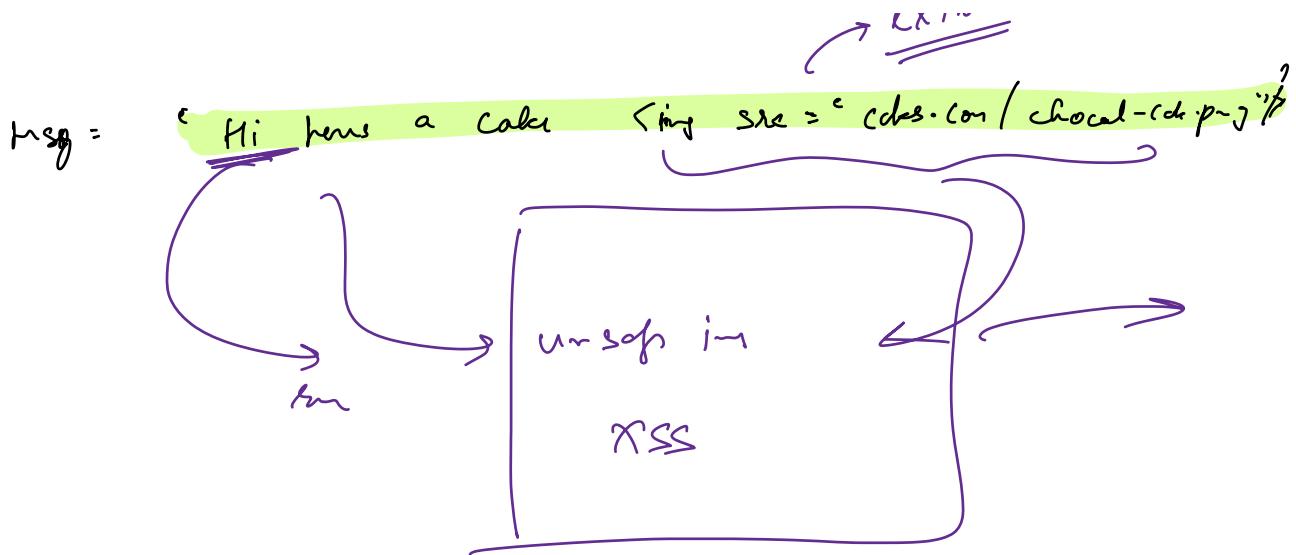


PUT → update

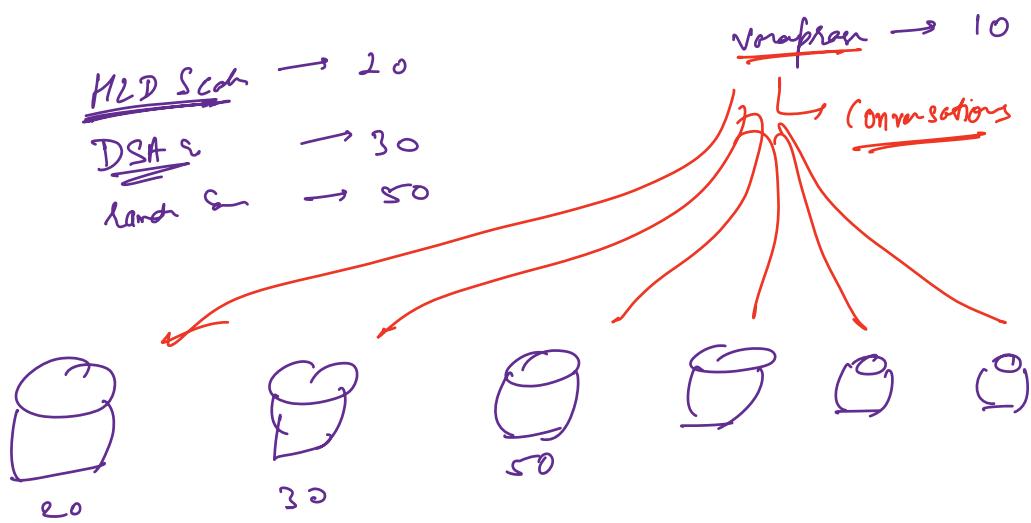
↳ always idempotent

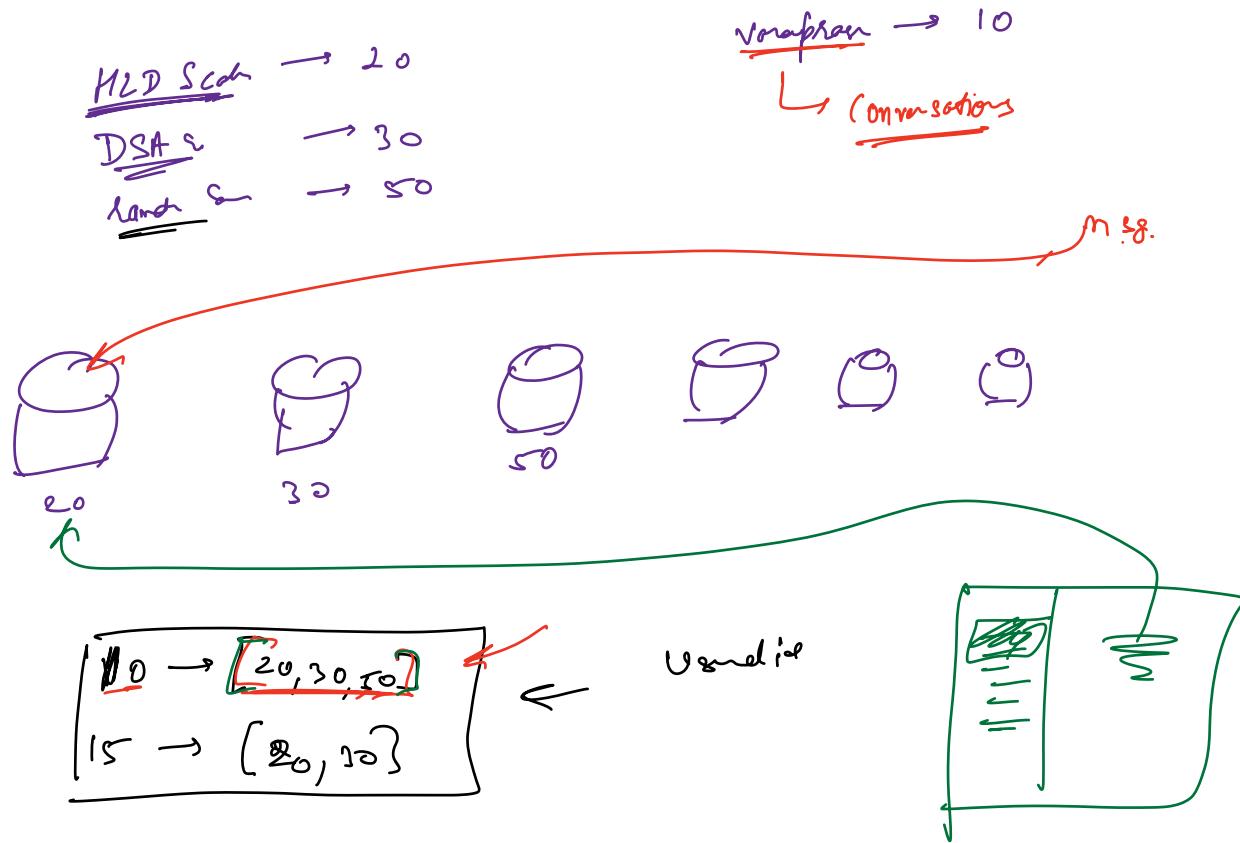


... instead.



Conversation-id [group id / (rec/sen)]



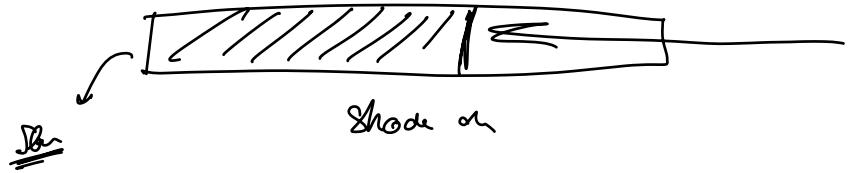


3×10^5 MAU for 800k
 ors 10 purchase each user
 & max 10 purchase each day

$3 \times 10^9 \rightarrow 10^5$
 10^4

10^6 purchases per day

$$\frac{10^6}{86,400} \approx \frac{10^6}{10^5} \approx 10 \text{ req/sec}$$

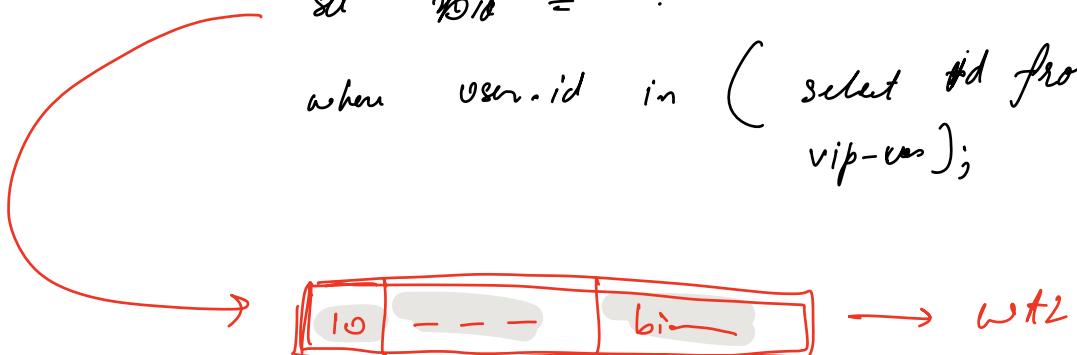


WA2

insert into table users

set ~~id~~ = '---'

where user.id in (select id from
vip-users);



SSTable → disk

Sparse Index → memory.