

Agenda



1. Invalidation & Eviction support@scaler.com
2. Local Caching - Case study 1 → Scaler's Code Judge
9-11:30 class
+10 minc
3. Global Caching - Case study 2
4. Case study 3 - Facebook News Feed doubt

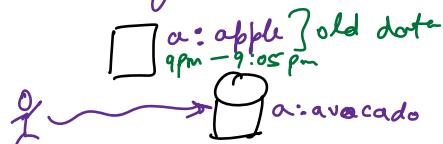
Most applications are read heavy → caching.



Challenges with Caching



- Cache is not the source of truth
↳ database
- Cache data can get stale



Cache Invalidation

↳ fix stale data



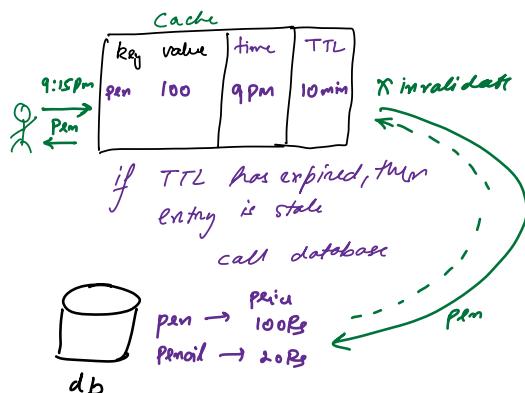
Time to Live (TTL)

Simple / Configurable TTL

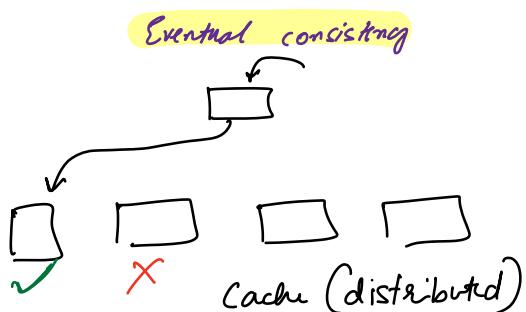
you will get slightly stale data

Write Through

- Small in size
- Cache miss → data wasn't there in cache
↳ lead to slowdown
- To load new data
Cache needs to remove existing data.
↳ what data to remove?



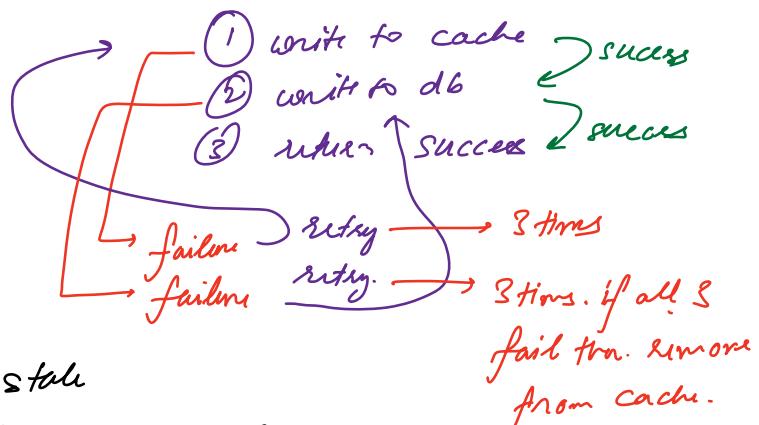
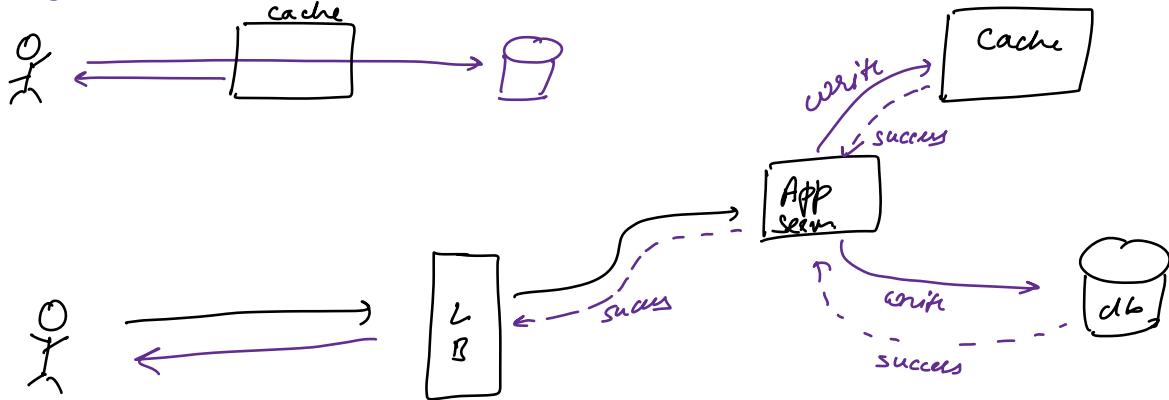
Write Back



Write Around

Write Through Cache

Any update must first be written to cache. } Synchronously (or after other) to both cache & db



? Can cache data be stale
in case of Write Through Cache?

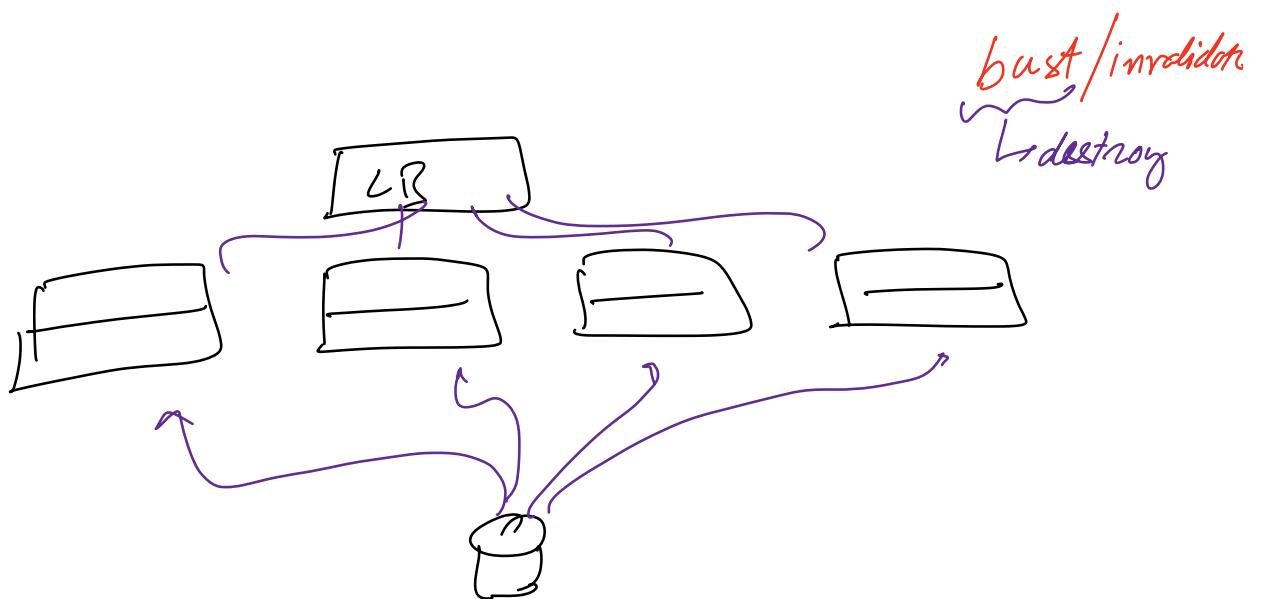
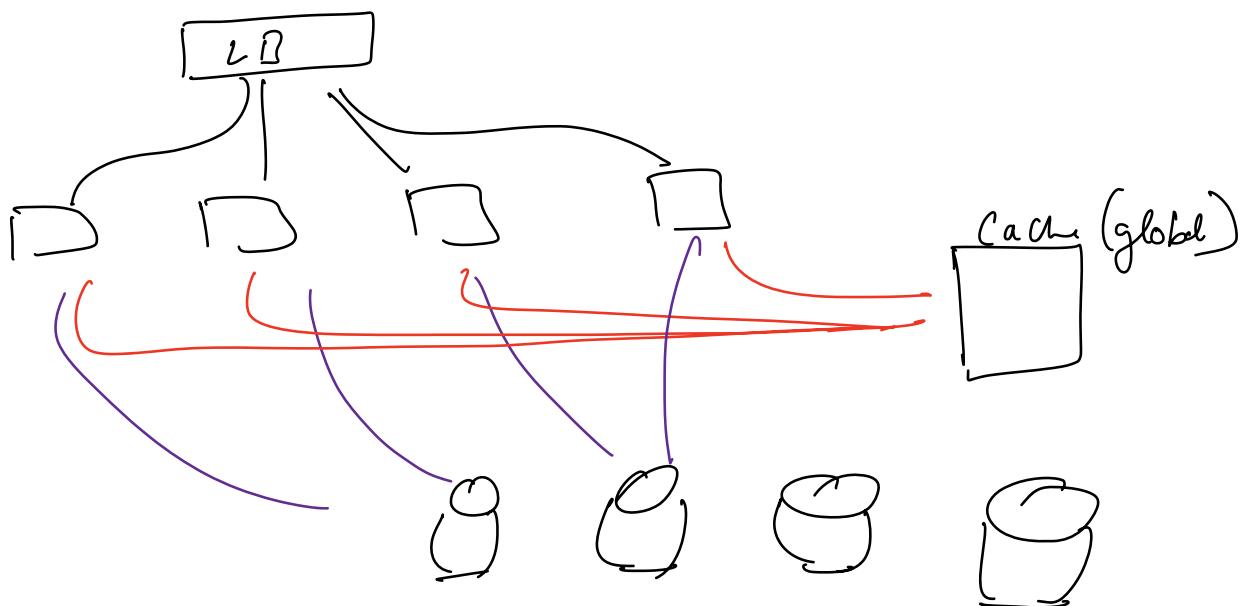
Never!!

- ① Always fresh data (in cache)
- ② Reads are super fast → every cache hit will return directly from cache
- ③ writes are much slower

Write through cache is good for a **read-heavy** system.

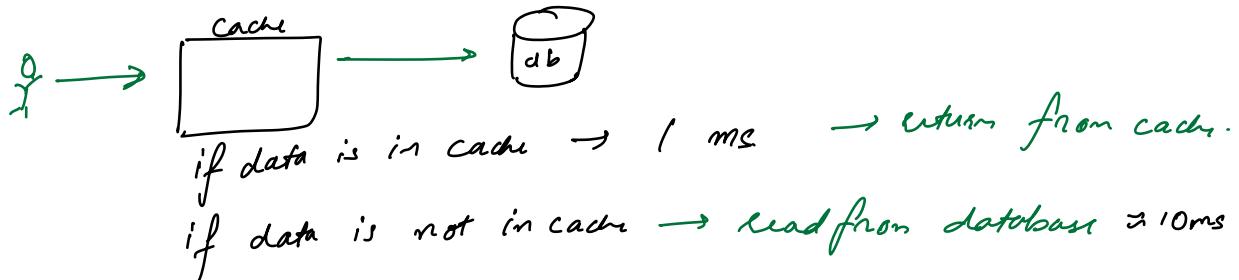
& strong consistency is needed.

most req. are reads & not writes
90% of applications



Reads are fast when using a cache

↳ ∵ of how cache works.

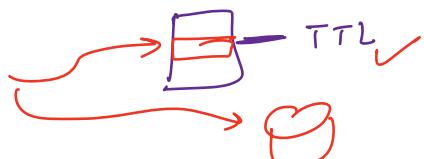


Hit → 1 ms → 99% cache hit

miss → 1 ms + 10 ms = 11 ms → 1% misses
Cache db

Avg. time per query.

$$\begin{aligned} & 0.99 * (1 \text{ ms}) + 0.01 * (11 \text{ ms}) \\ & = 0.99 \text{ ms} + 0.11 \text{ ms} \\ & = 1.1 \text{ ms} \end{aligned}$$



Write around cache

→ writes happens directly at the database

↳ the cache might be out of sync for a short while

→ cron job that queries recently updated data from db (10 mins) & updates in cache.

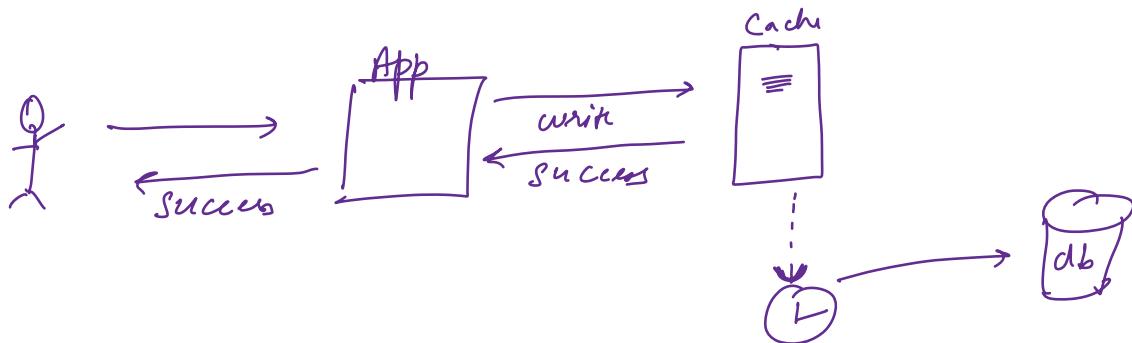
select * from product-price
where updated_at > (now() - 10m)
} → update in cache.

eventual consistency

Write-back cache (*inconsistent*)

↳ first write data in cache
↳ return success
even though data is not on db!!

↳ *chronical*
cron job that will periodically flush
this data on to the database.



Cons

if cache server crashes we lose some data!!

lose

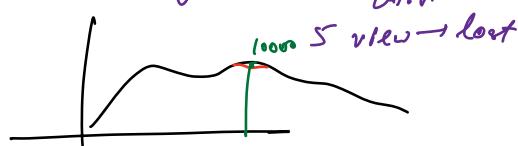
Pros

Extremely high ^{writes} throughput

Use write-back cache in write-heavy system
only if consistency is not important

Analytics / tracking clicks or ads

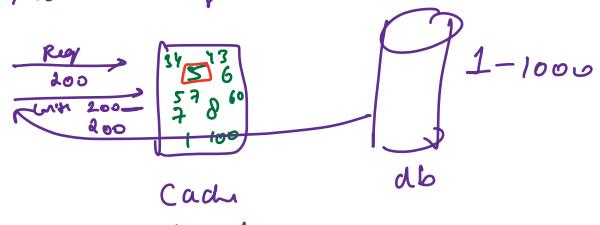
↓
trend based.
↳ 10,000 views
↳ only 9,995 views got credit
↳ 5 view → lost





Cache Eviction

Cache is limited in size



First in First Out

Simple
Simpler



Least Recently Used

time
↳ Hashmap + Heap
↳ Hashmap + DLL

if you see something, it is v. likely
that you will see it again in
near future.

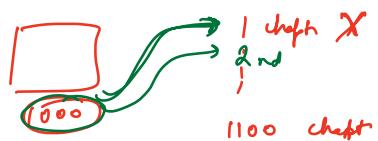
Last in First Out

↳ whatever you've read recently will get removed.
HW → find out circumstances when this is awful.

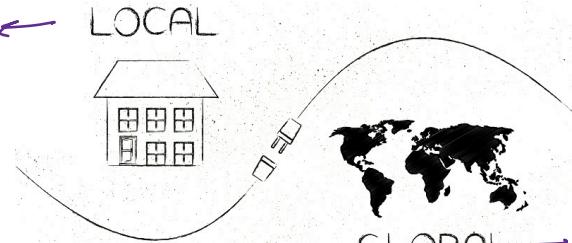
Least Frequently Used

↳ evicted.
read counts

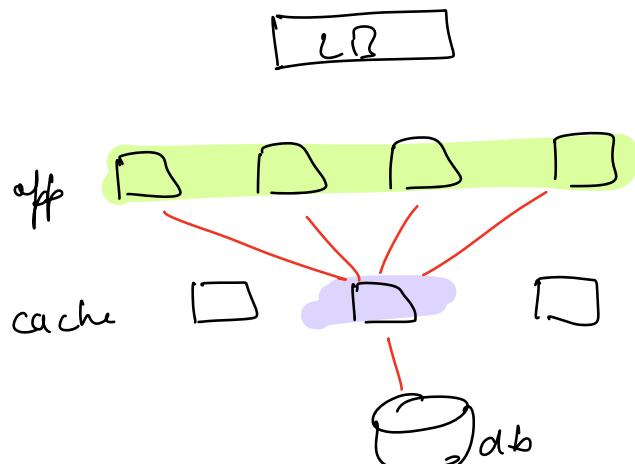
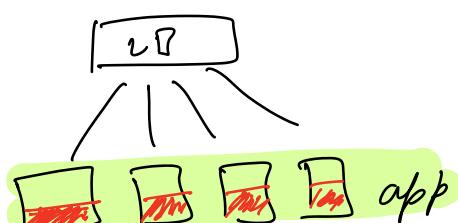
One Piece Meme: Reading groups.



used by only
1 appserver.



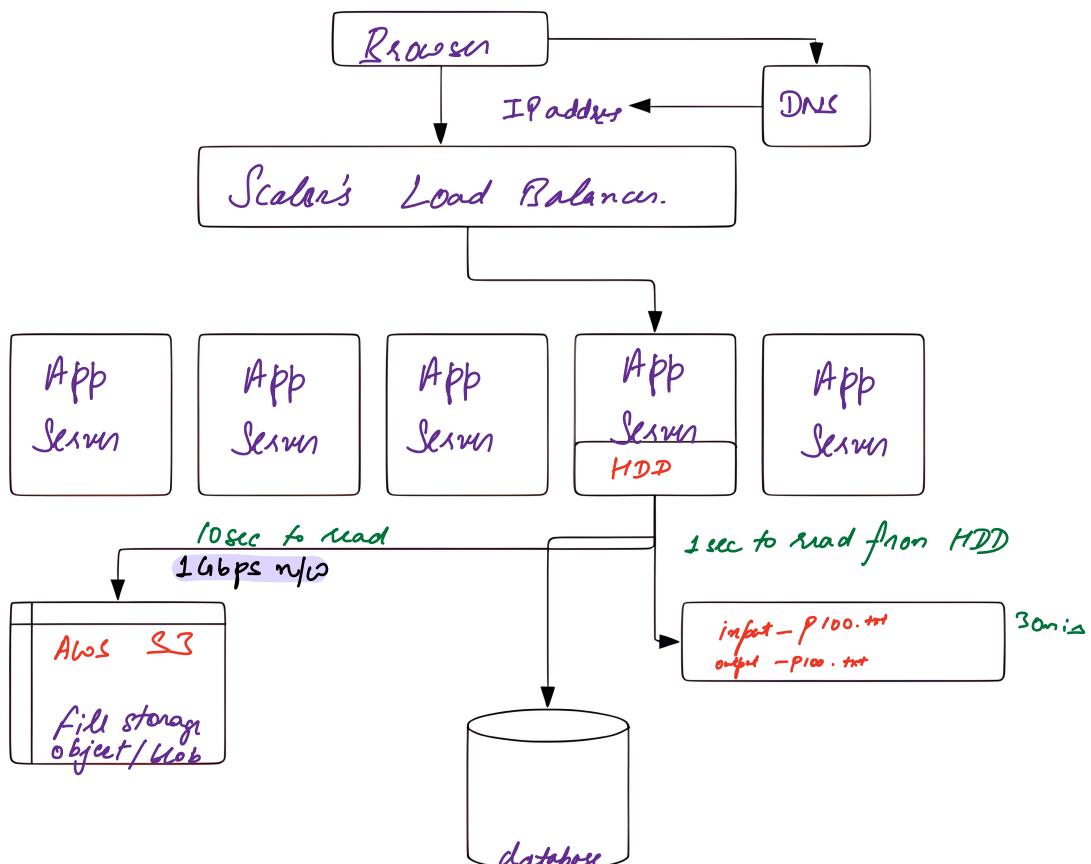
→ shared by
all appserver



CASE STUDY



Online Code Judge



<https://gist.github.com/jboner/2841832>

? What does the AppServer need to evaluate a solution?

problem-id problem-metadata

user-id solved-already score

code language. input testcases & output testcases

? How large can this data be?

input
output } 500 MB each → 1 GB of testdata
for a single problem!

given array → sort it
 $O(n \lg n)$

$\approx 10^6$ size

40-50 testcases

500MB input.
500MB output



Change to the testcases

if a test case is changed → how do we update it in the local cache?

- at every write, push change to all app servers.

↳ write through cache



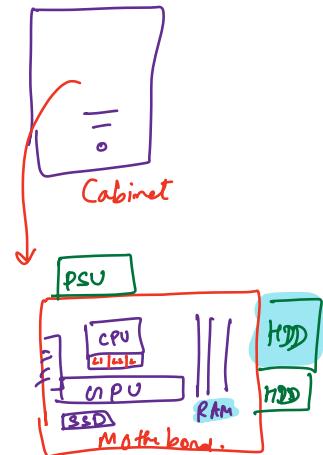
→ every app server has its own cache.

→ update every app server

→ wasteful

↳ write back cache & TTL

some period during which data will be stale.



CPU runs at $\approx 2.6 \text{ GHz}$
read/write
8 cores

$3 \times 8 = 24$ billion op/s
0.3 ns

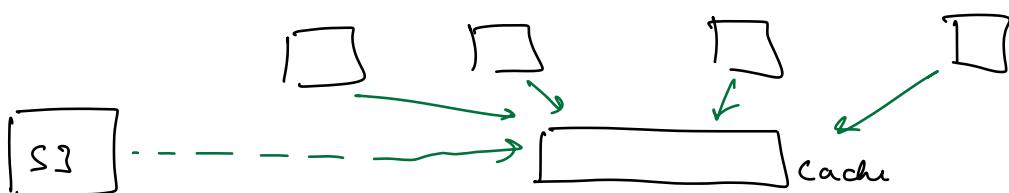
- TTL / write back cache.

bad solution → we want immediate effect

10mins TTL → entire content gets over &
test case are not updated

30sec TTL → most req will have to go to S3
no benefit of cache

- Global Cache

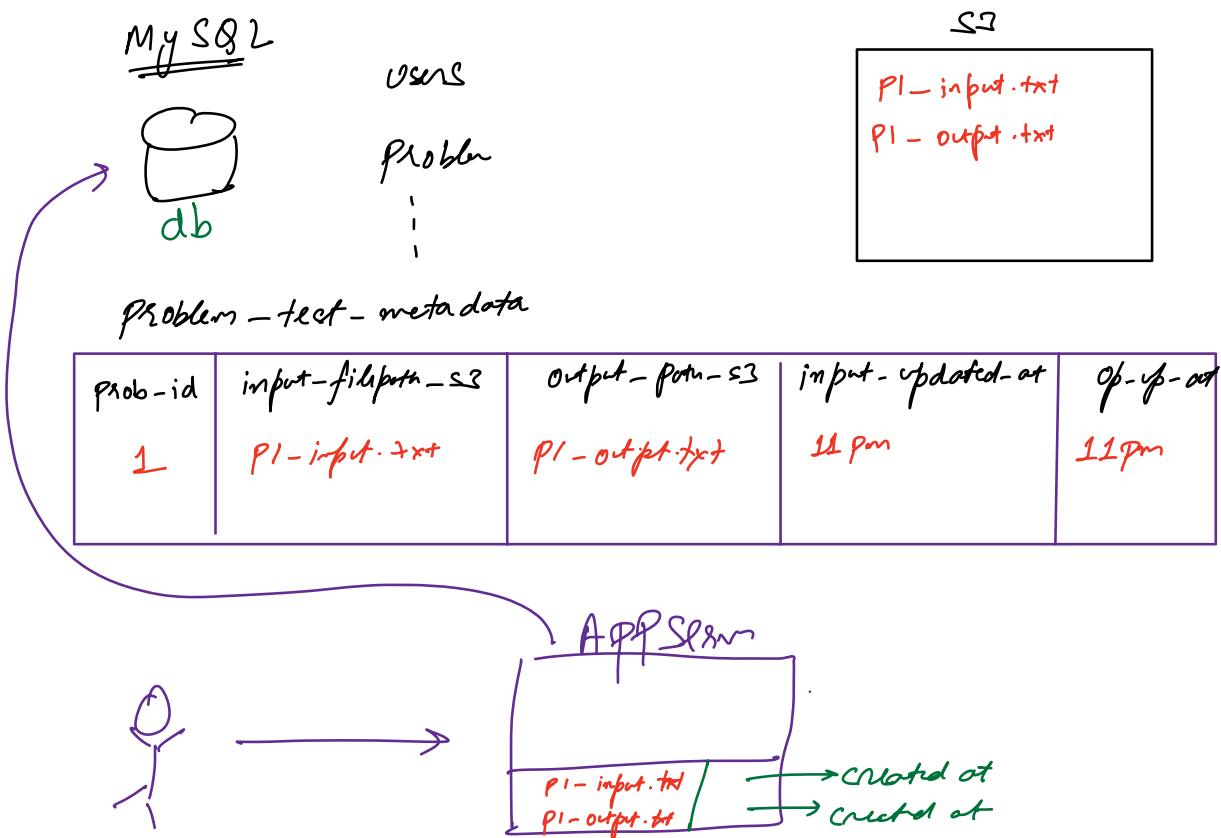


- ① a lot of m/w traffic will come to cache
- ② if cache seem has 32GB RAM
↳ 32 files

a lot of cache misses.

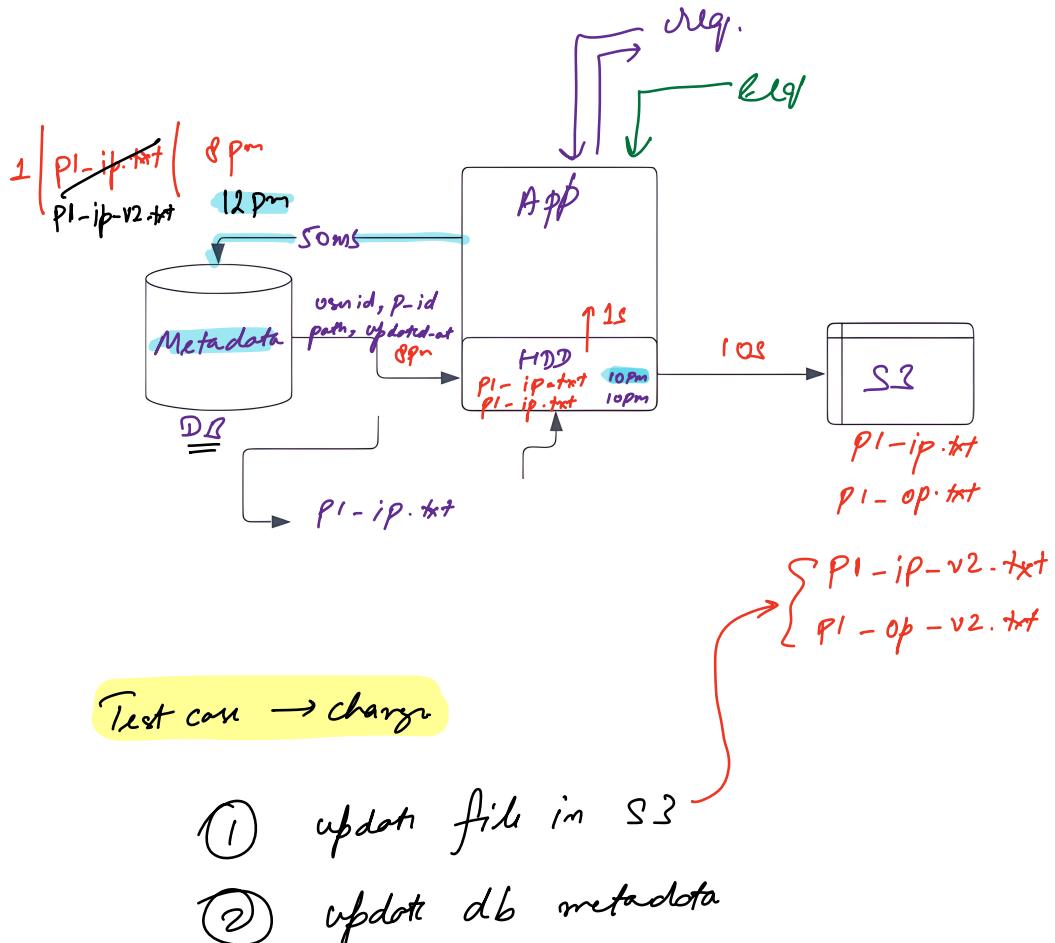
- ③ store in HDD → SSD ↳ $\sim 300 \text{ MBps}$
- HDD

File Metadata in DB





Solution



CASE STUDY



Contest Leaderboard

a lot of traffic
hundreds of submissions
every minute

? How can we compute the ranklist?

User-id, name,
Problem-id, score

User-problem
U-id, prob-id, status ✓ TLE
WA

? How frequently does the ranklist change?

true ranklist update
after every submission.

} several times per second

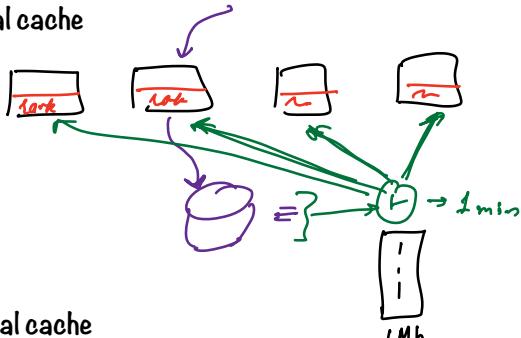
? Is immediate consistency important?

not important

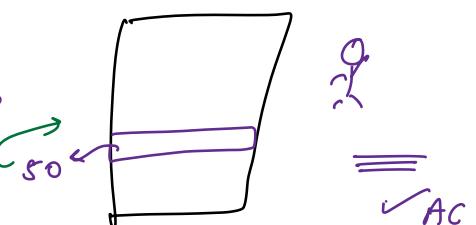
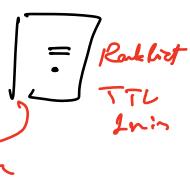
ranklist update every 1 min



Local cache

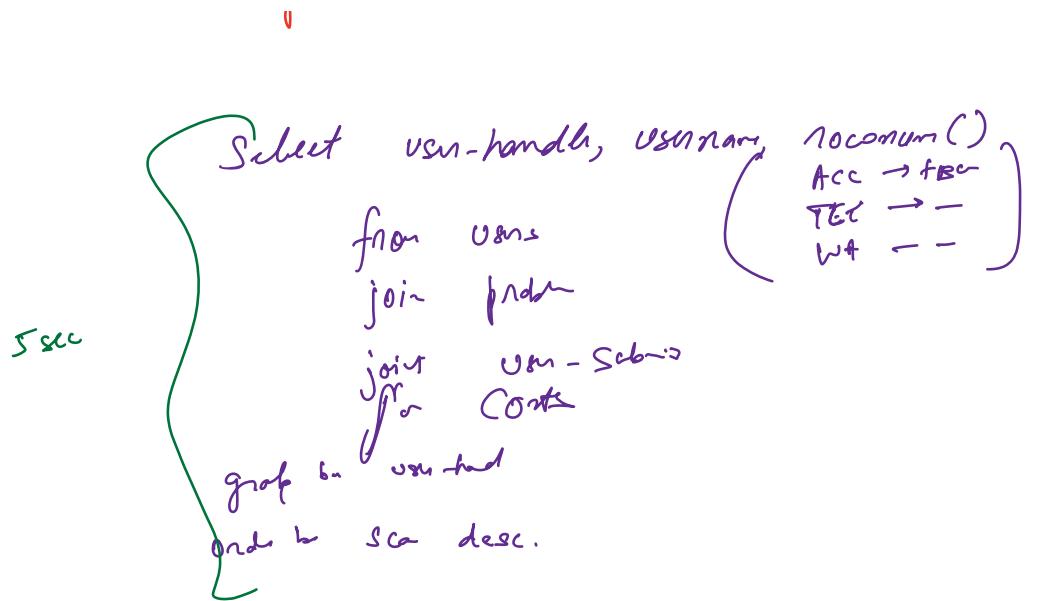


Global cache



TTL/with around
Cache
write-around (push)
re-computing ranklist once
every minutes Not bad ✓

TTL (pull)
every 1 min each appsum
will go to db & re-compute
ranklist 100 hrs each
minute X bad

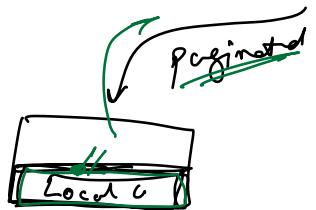


Client - around	global
Client job that recomputes even 1 min ✓	Client job except even 1 min ✓
Cache is local to each appserver	Send data to global cache
Send full entire ranklist to each appserver. 100 msg/min	queries to read from cache or paginated \rightarrow small data. ✓

appserver can serve
data from local cache ✓

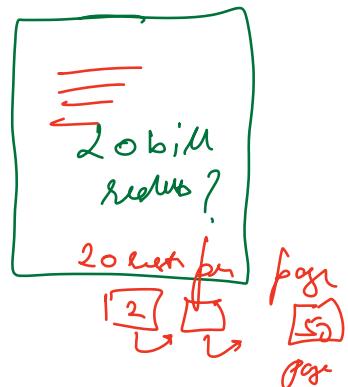
User-handle \rightarrow max 32 char string } 36 bytes
rank \rightarrow 4 byte intign } 100,000 users

2,600,000 R
3.6 MB



google

20 billion results in 0.1sec



(O cont)





Redis

C

try.redis.io
redis.io/docs/data-types
university.redis.com

H/w

single-threaded key-value store
(in-memory)
rich data types

String / Integers / Booleans
Lists / Sets / Sorted Sets

Page = 1000
20 rows
per page

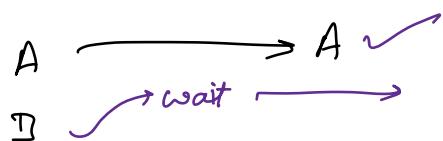
2000
1020

- ① give me ranklist for page X
- ② given a user-handle, find rank

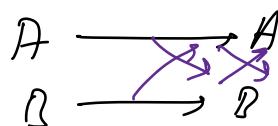
Key : Contest_id - user_handle

Value : Score, Rank

SLSVM is 1-threaded



multithread



10:30 → 10:35

CASE STUDY



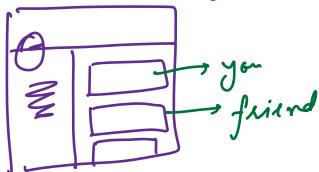
Facebook News Feed

? How can we compute the News Feed?

Posts made by me
Posts made by friends,
" " people / Pages you follow
" " groups
Posts by other X that your friends have interacted with.

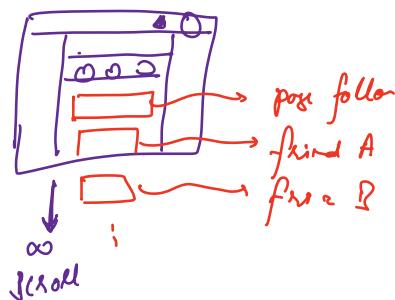
Recent data

Profile Page ✓

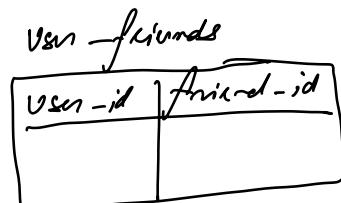
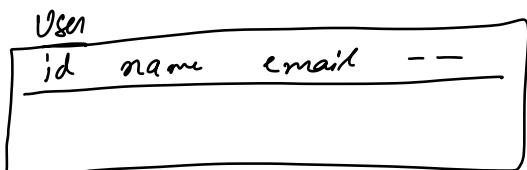
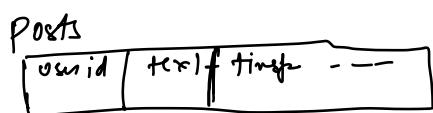


select * from posts
where user-id = _____
order by time desc
limit x offset y

News Feed X



select * from posts
join user
join user-friend
join user
order by time desc
limit x offset y

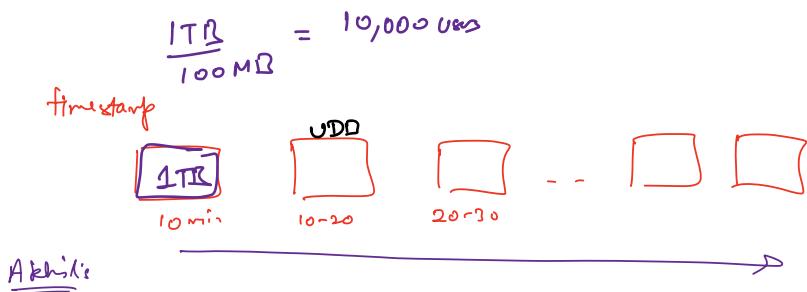


? Does all the user data fit on a single machine?

No. shard data

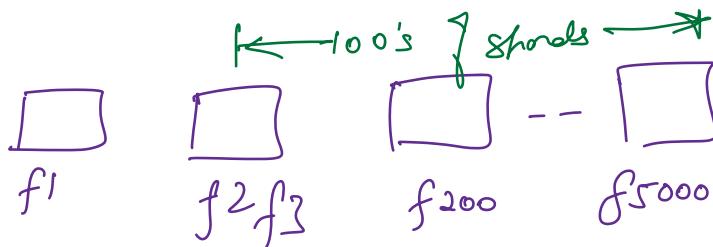
? What is my sharding key?

~~Post-id~~



Shard by user-id

? How do I retrieve data for the news feed?



UDN
↳ user database

more size of 1 user data

→ posting once every hour
→ all year for past 20 years

⇒ $12 \times 865 \times 20$

$\ll 100,000 \Rightarrow 1KB \Rightarrow 100MB$

X not work

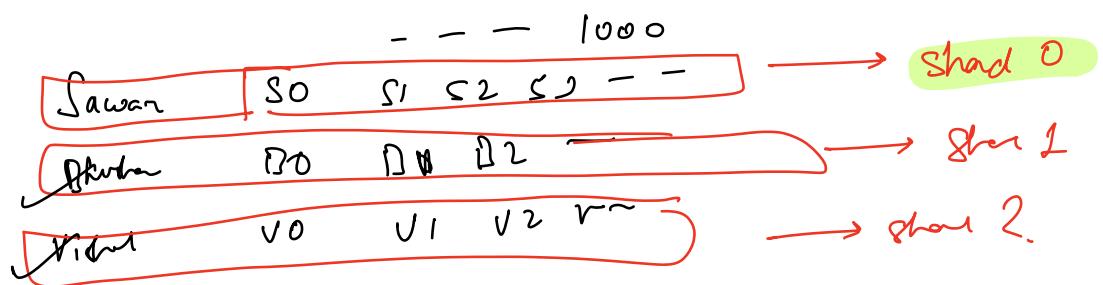
? How can you optimize the construction of the News Feed?

① Even if we cache it → first issue is how to compute it

to compute a single news feed we need to hit 100 shards

1 → store posts in sharded dbs
e → for news feed only store post-id

expensive

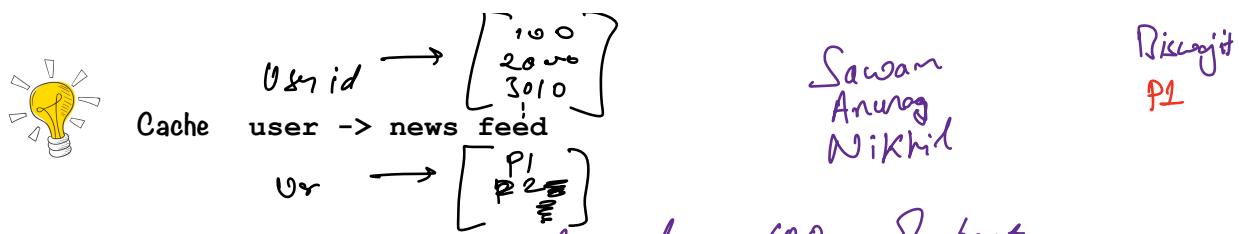


,

,

,

+



Sawan
Anurag
Nikhil

Ricard
PL

1. Billions of users, news feed has 100s of posts and each post is replicated in many feeds
2. Fan out updates: every post has to
3. update in news feed to be copied to all feeds X
algorithm → v. hard



Back of the envelope estimates

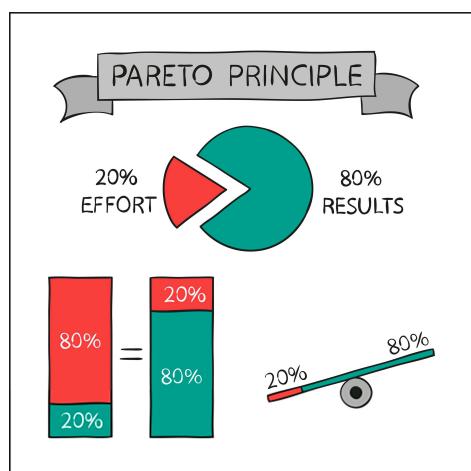
Facebook → 1+ billion MAU

→ 800 million DAU

posts made per day. $\Rightarrow 800 \text{ M} * 1\% * 10 = 80 \text{ M}$

posts per day by avg user / 10 posts
(over estimate)

how many people would post
on a day?



80:20:1 principle
↓ ↓ ↓
passive interacting general content
reading (like (posts))

? Daily post data

$$80 \text{ M} * 1 \text{ KB} = 80 \text{ GB}$$

posts						
P-id	U-id	Content	Location	Media Path	Timestamp	
8B	8B	250	8B	250	8B	1 KB

? Monthly post data

$$80^{\text{th}} \text{ mo} = 2400 \text{ GB}$$

= 2.4 TB

1 byte → ≈ 1 billion values

Sharding → Partitioning across Shards

will improve read/write speeds [not true]
true purpose → reduce data size

split data so that each partition has different data

Replicate data → multiple servers having

same copy of data

→ Improve read speeds

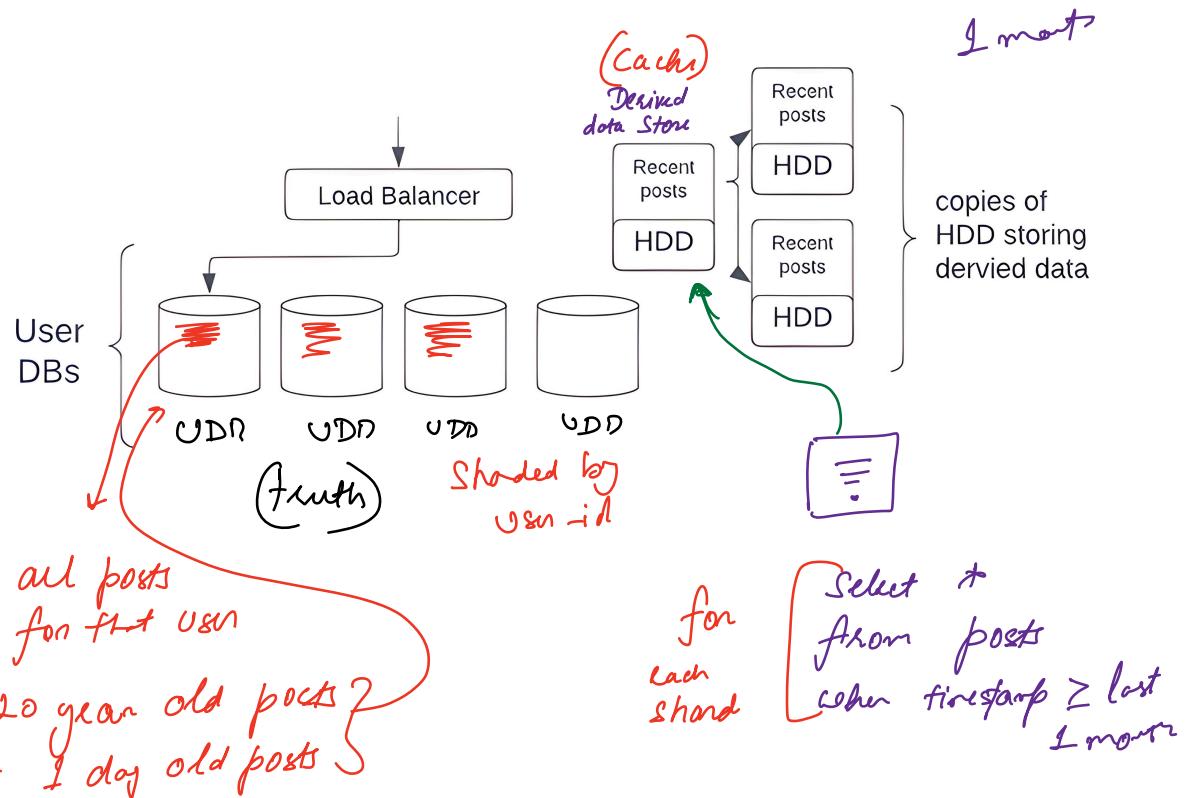
→ availability

Ctrl + C

Ctrl + V

○ ○ } 1 hour old
○ ○ }

Singh sevon!!



Akhil → news feed

① fetch all friend ids of Akhil

② select recent posts made by user_id

Select * from
derived on cache all_recent_posts
where user_id in (friends of Akhil)

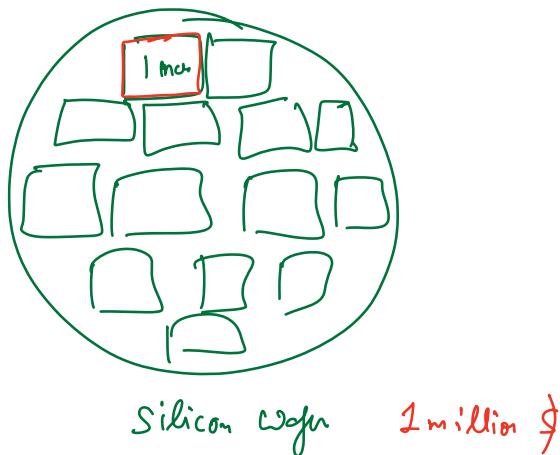
- ① Update data \Rightarrow cron for each db shard that updates last 30 min data into the cache
- ② Remove old data. \Rightarrow every night at 6 am (most traffic time)
delete * from all_recent_posts

When timing $>$ 1 order old.



CPU \rightarrow Intel
clean rooms

ULS I
Ultra large scale integration
3nm
5nm
V.V.V. expensive



flip-flop \rightarrow registers 50Mb
 \rightarrow cheap but slow
L1
 \rightarrow clean but slow
L2
;

-
- ① write through \rightarrow no stale data
 - ② write around \rightarrow cron job that updates cache using db
 \approx TTL
stale for some time
 - ③ write back \rightarrow inconsistent data

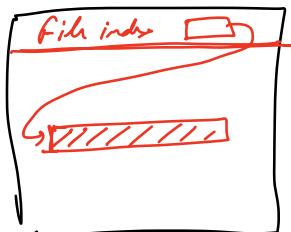
meta data

→ s3 → same of that? did we have stale data?
→ db metadata

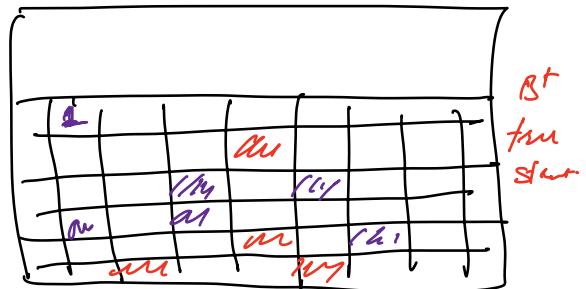
No stale data!



S3
↓
NTFS / CxTz
fill storage

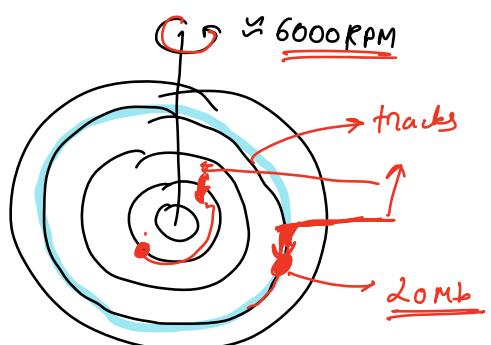
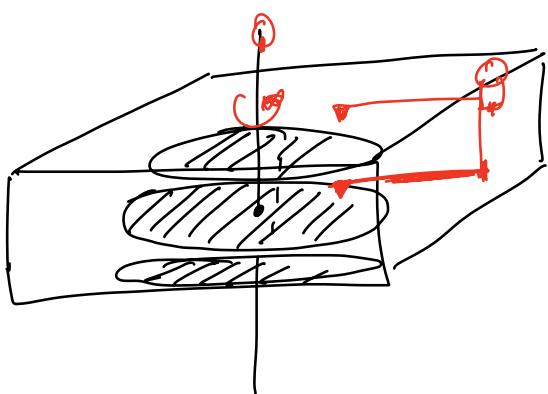


DB querying data



continuous data access

P-id | input soon / out soon
random data access



$$\text{disk seek} \Rightarrow 100\text{ms} \quad \frac{6000}{60\text{s}} = 10\text{ms}$$

$TTL \rightarrow$ stale data

global \rightarrow high m/o usage

\hookrightarrow HDD \rightarrow slow

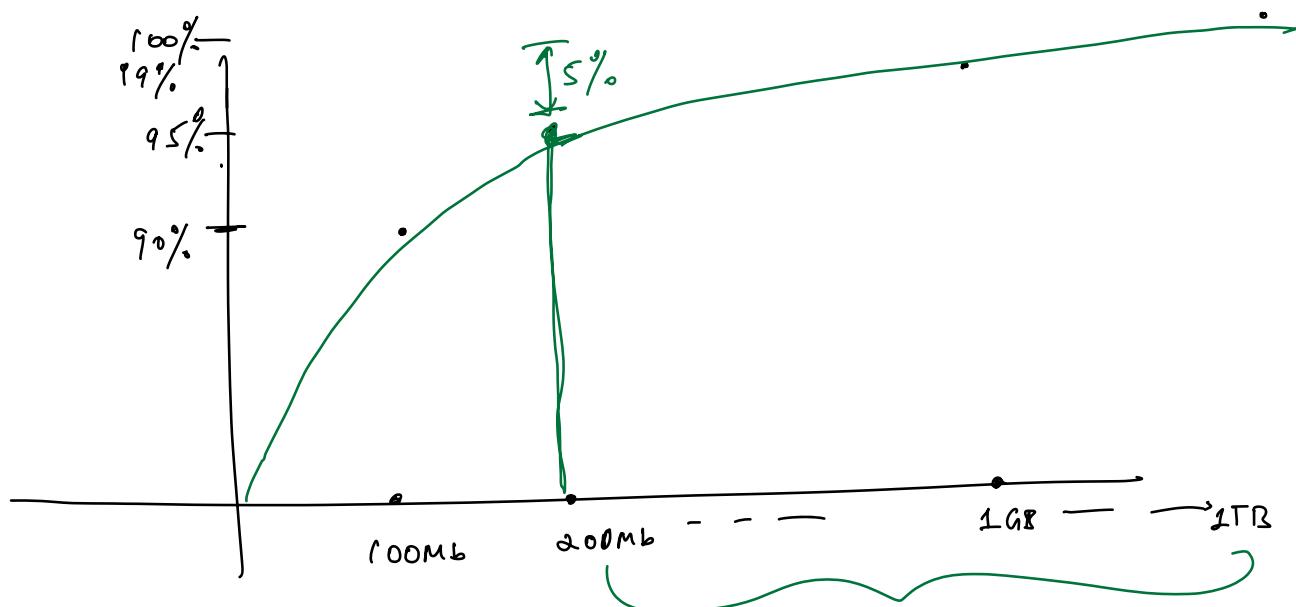
\hookrightarrow RAM \rightarrow small

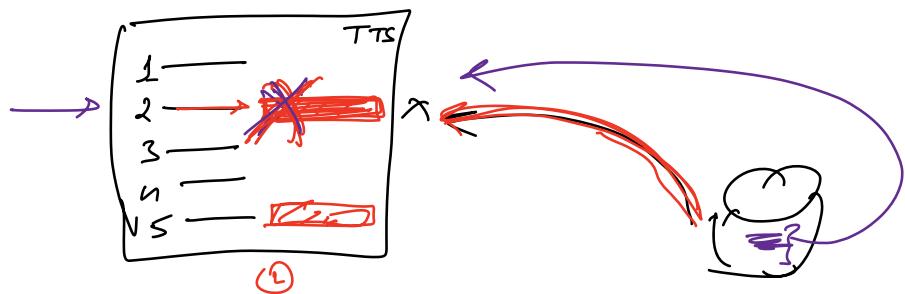
written back $\times \rightarrow$ inconsistent

written around \rightarrow similar to TTL

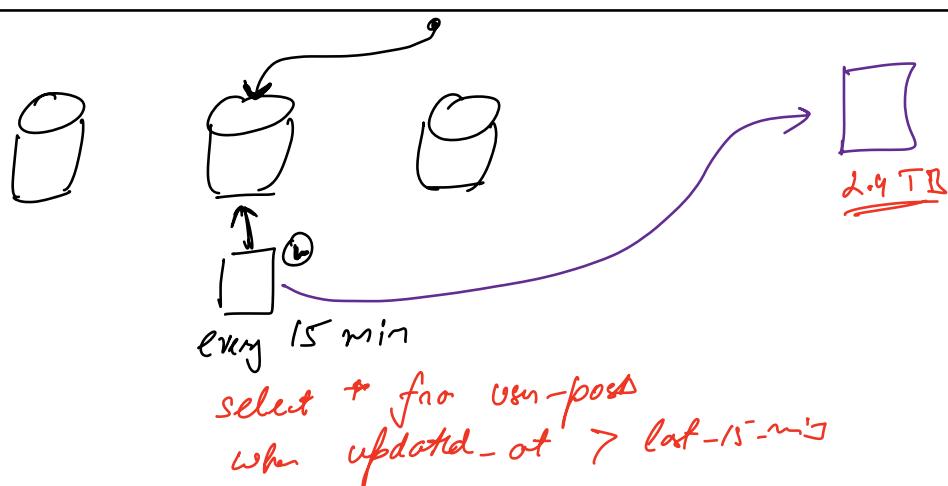
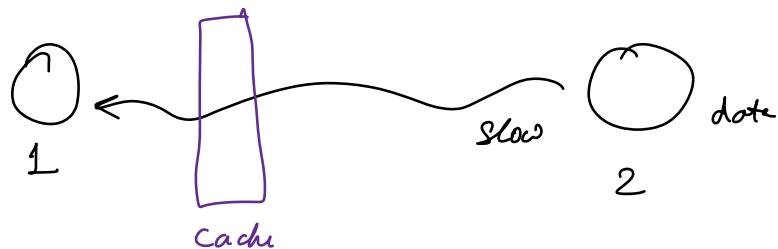
DB \rightarrow 1 TB

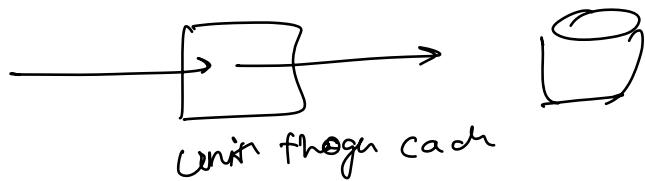
Cache \rightarrow 100 Mb





- ① wait for req (2) → expired → invalidate
 - ② Redis detected that (2) is still → remove it
 - ③ cron → stale data → refetch from db to update in cache
-

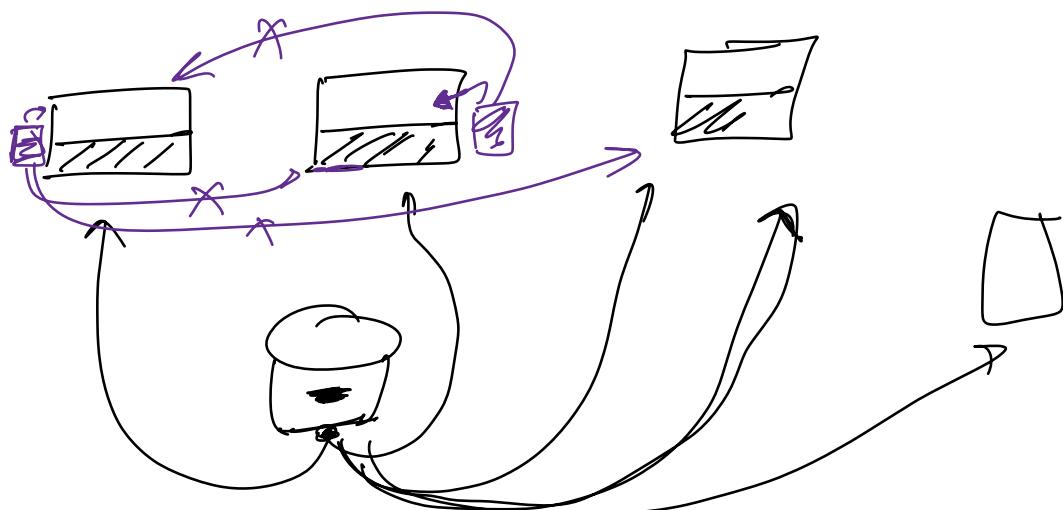




Cache & db must always be in sync.

TTL

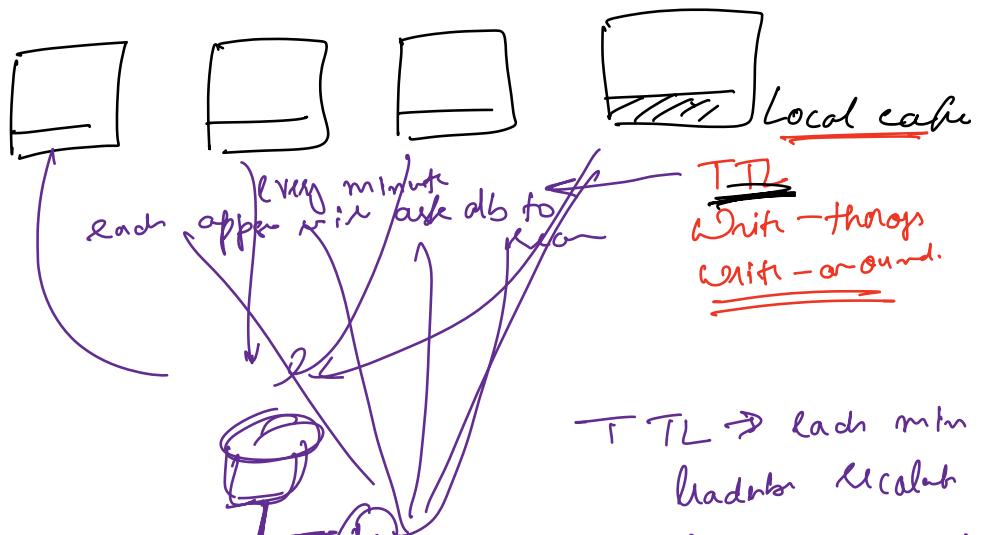
Write around cache



CDN → Cache for infrequently changing data
 (mostly static) time/sec
 audio/video/img
 geo-located so that user don't have to hit server far away.

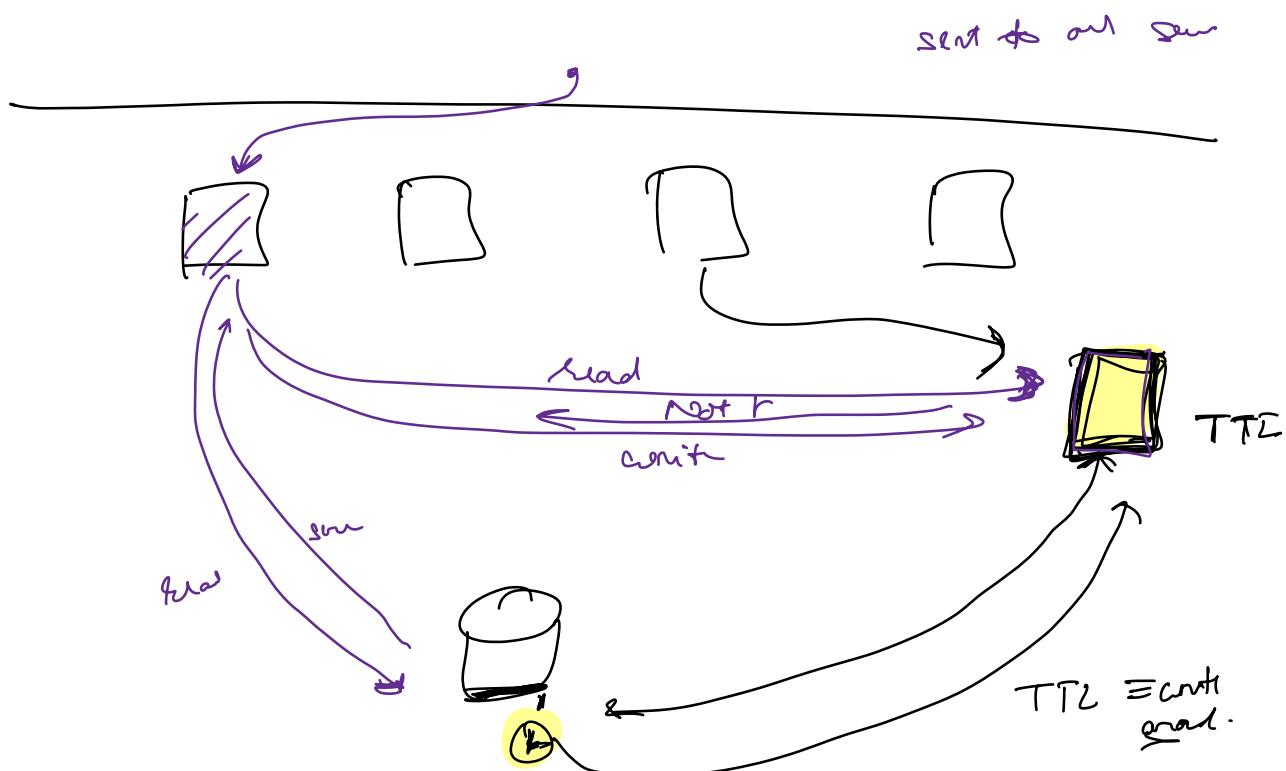
S3 → source of fast to store files
 db

Leader board



TTL → each min
leaderboard calculate positions
& send one w/u
to server

Write one

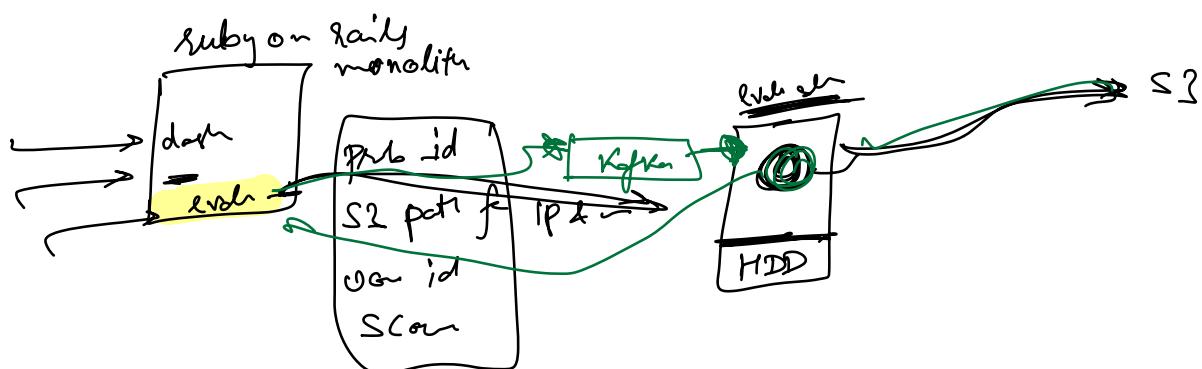


1. Latency req/sec

consistency

fast responses

- ① service | proba etc
- ② Read how/with how



USM - Post - Views

→ write back Cache



1 Billion

100 Billion post

1B + 50,000
views



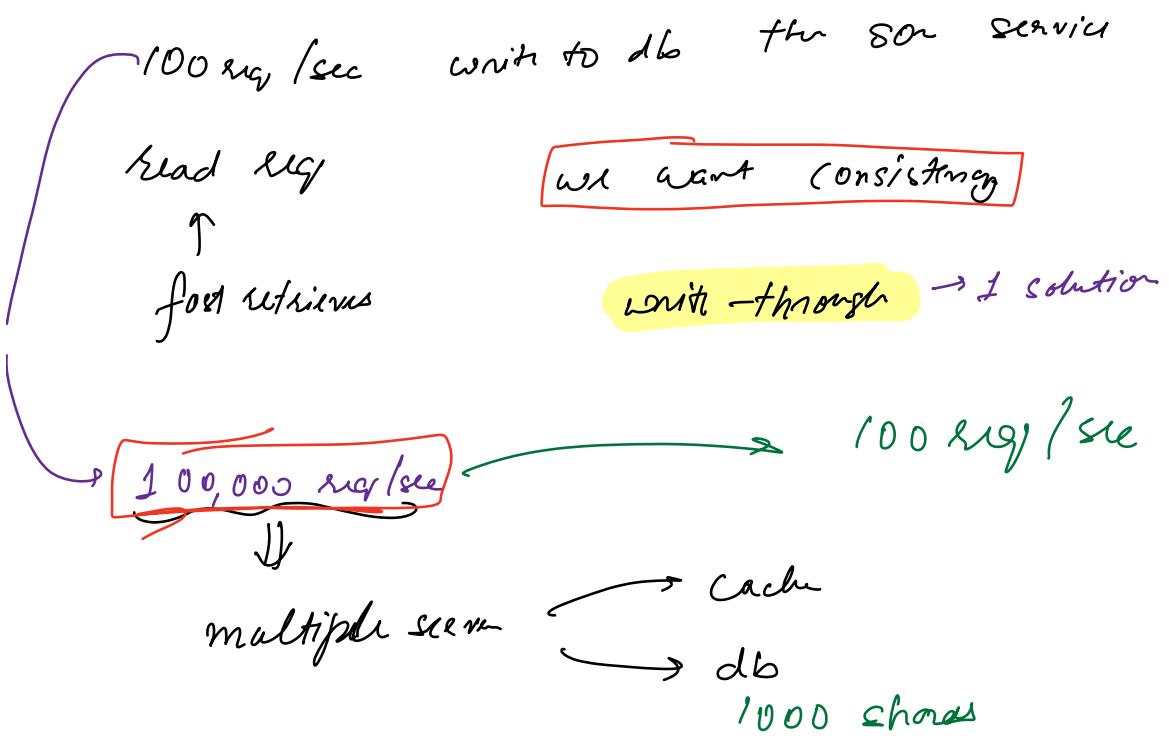
Page → 20 each

5322/20

→ int div.

Range

5322



1.4 billion people

50% of the use UPI

0.7 B people

avg user will use UPI how many day.

10 req per day

7B req per day.

$$\frac{7 \times 10^9}{86400} = 1 \times 10^5 \text{ req/sec}$$

1+ Lakh per day
25
Bachelor
Software

for 1% people
of India.

# of users	MAU	2 Billion
	DAU	800 million
↳ activity is the user daily		80:20:1 rule 80:20 rule

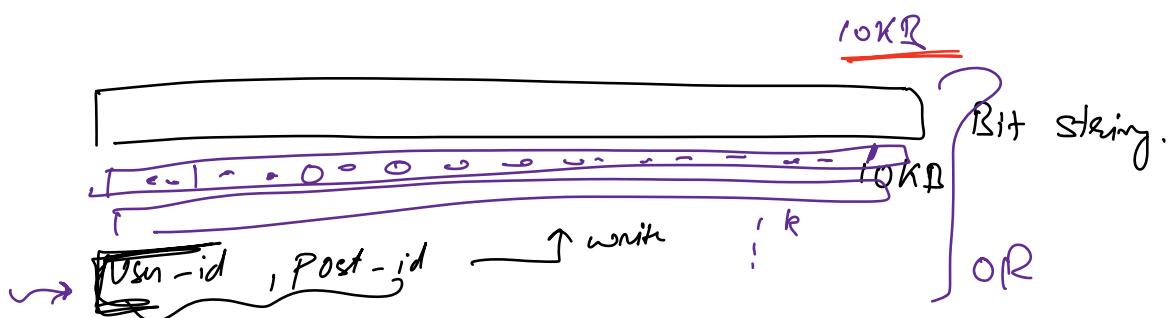
of activities per day → per second
86900

27,000 bookings per min March 2020

$$\frac{27000}{60} \approx 450 \text{ rev/sec}$$

100,000 cycles
in sum!

S. Sudarshan → DRMC
Head of CSE → IIT R
DR of IRCTC



Hash(User-id, Post-id) - - - for hash.

 New value
No kid