

Agenda



1. Invalidation & Eviction Recap
2. Local Caching - Case study 1
3. Global Caching - Case study 2
4. Case study 3 - Facebook News Feed

support@scaler.com



Challenges with Caching



Data can become stale
∴ DB is source of truth
↓
Invalidation

Small size — cache can get full
Eviction



Cache Invalidation



TTL

- Eventual Consistency
- fast Reads (10K/s)
- fast Writes (100/s) write → cache miss / TTL expired

Immediate Consistency

Eventual

No Consistency (data loss)

Write Through

- Immediate Consistency → write to DB & cache in atomic manner
- fast Read (10K/s)
- slow Write (10/s) → 2 Phase Commit (Roll back / sync)

Write Back

- No Consistency → Data Loss!!
- fast Reads (10K/s)
- ultra fast Writes (10K/s) → ∵ only writes to in-memory Cache

Write Around

- Eventual Consistency
- fast Reads (10K/s)
- fast Writes (100/s) → writes via cron job



Cache Eviction

CDNs are only for client facing data

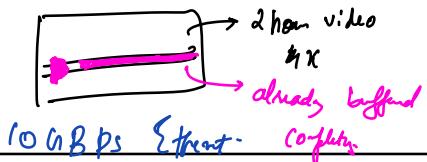
First in First Out → 1.9%

Least Recently Used → 89%

Last in First Out → 0.1%

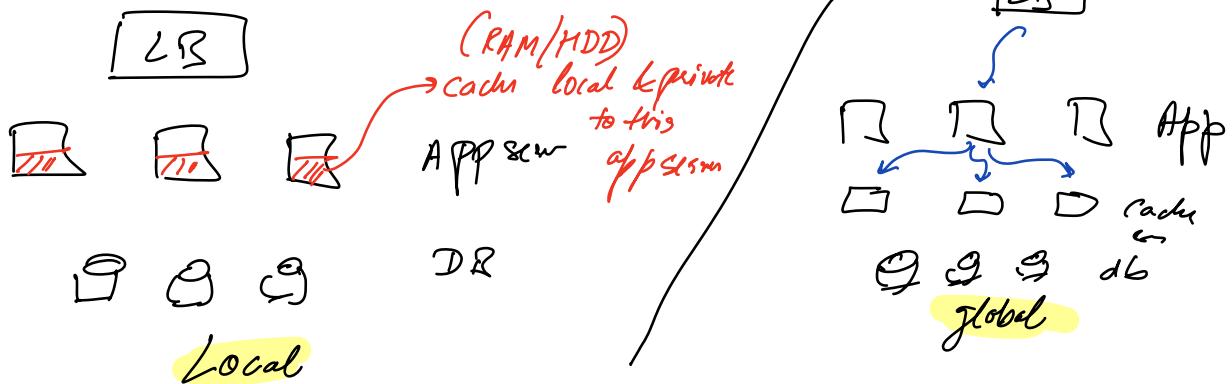
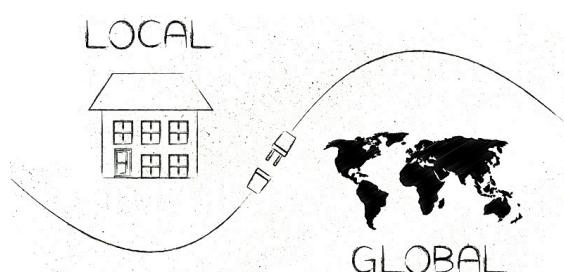
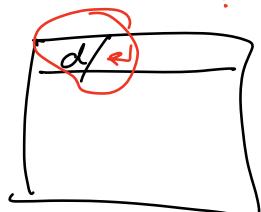
Least Frequently Used → 9%

CDN → Large companies like Akamai
Org might have their own private CDNs.



netflix → global cache for all orgs
per user cache
predict what user might watch
GoT → Ep 2 watched
likely to watch Ep 3 next

If I watched Ep 3 last, am I likely to watch it again in near future? X
No!





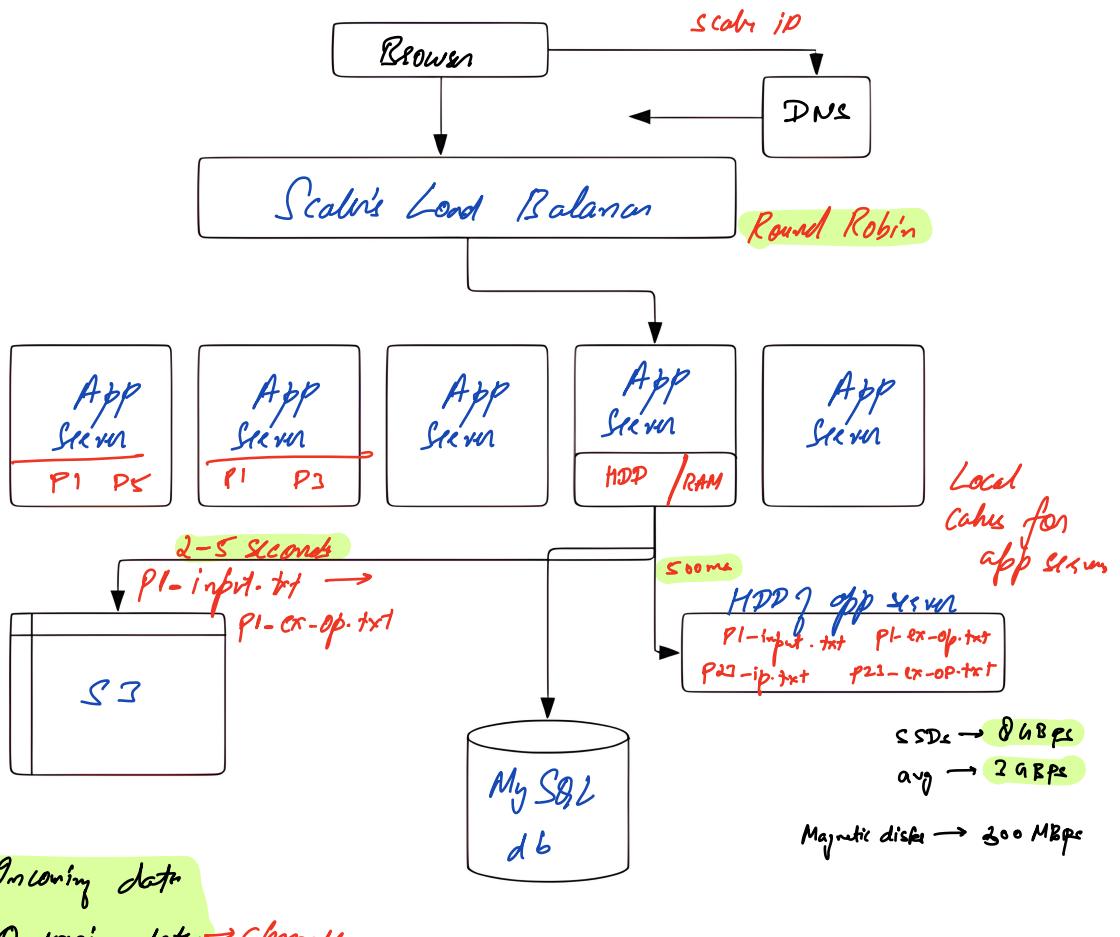
Scala

Online Code Judge

~300,000
participants

3 hour

5 problems



? What does the AppServer need to evaluate a solution?

OSn-id

Problem-1-input.txt

Problem-id

Problem-1-Expected-output.txt

meta data about osn-problem

Scan for each testcase / ML / TL

? How large can this data be?

Test Cases can be super large

500 Mb per file

Input arrg of 10^6 items.

To evaluate 1 submission we need to load 1GB of data

Current Caching Architecture

Source : S2

Each app server has a local cache on HDD
to store frequently used data

Why Not RAM?? \because we have $\sim 150\text{GB}$ of data \rightarrow too long for RAM

Why cache? why not directly load from S3?

\because S3 load takes 2sec
(avg)

HDD load takes
0.5 sec

\because S3 load costs money
(egress)

free

Don't App servers have duplicate cache?

Yes & No.

Each App server has only 1 ip & one file
for each prob

Actions appear same on duplicates

APP Server 1

P1
P2
P3 } cached

APP Server 2

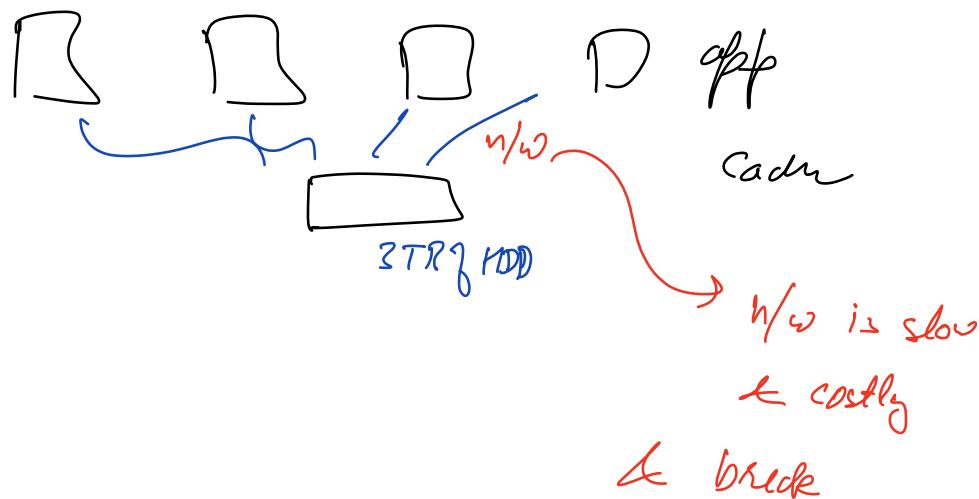
P1
P2
P3 } cached

Can we cache all problem futures?

No. ∵ Scale has 2000+ problems

≈ 3TB of data

If there is duplicate data in cache, why not a global cache?



On any day (even with contexts)

~ 100 - 150 problems are being solved

?

Change to the testcases

Requirement:
→ able to change any problem's testcases anytime
→ low latency for updates (1 min)

Question → Cache Invalidation

Update testcase in S3 → Cache testcase is stale

Write Around

X TTL → add cache for each appsrcm → all stored last fetched time for testcase

P1 - input.txt
P1 - out.txt

P3 - input.txt
P3 - op.txt

P10 - ip.txt

Disk

{
P1 : 8pm
P3 : 5pm
P5 : 3am
P10 : 9:57pm
Hash Map → in Memory}

TTL ⇒ 1 hour

↓
Optimal value??

High TTL → 2 hours

too stale data → content is over.

No Good TTL

Low TTL → 1 min

every cached file is out of data very frequently
too many cache misses!!

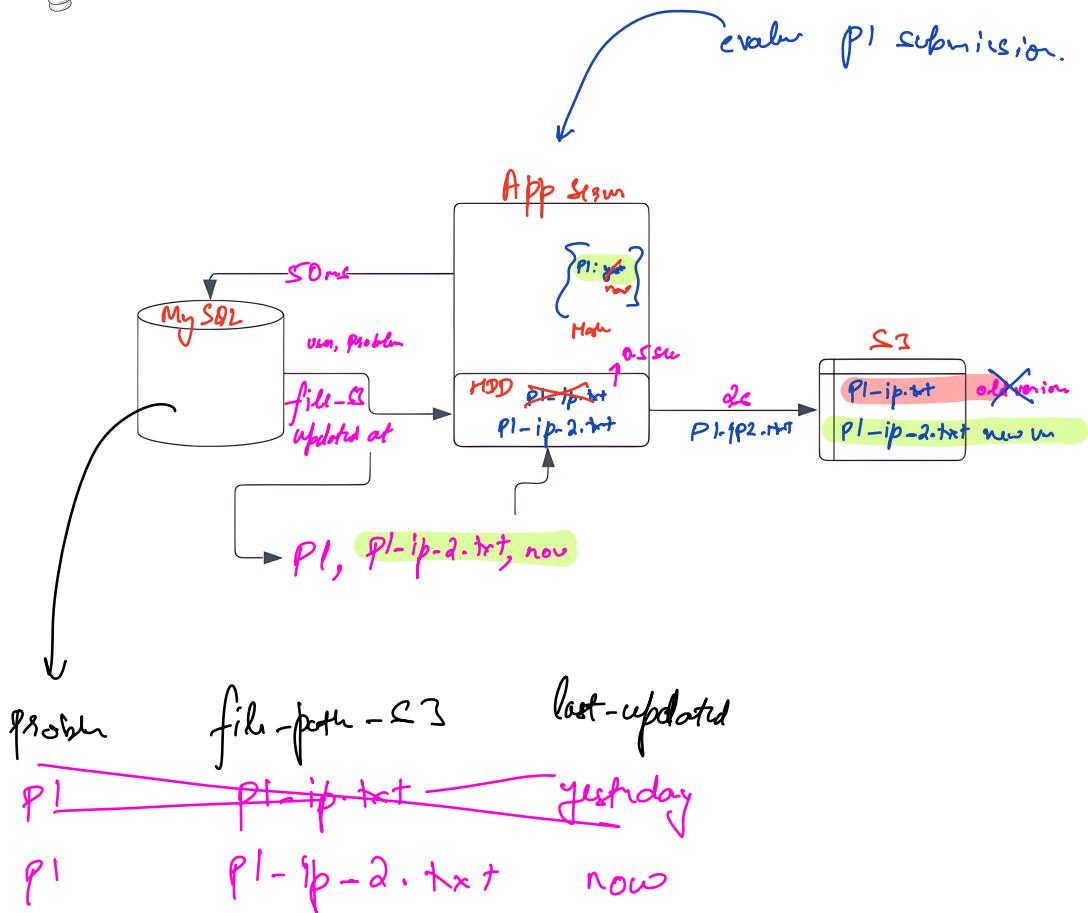
X Bad



Solution

Write-through Cache

1000 app session
due contact



Yes! App Session must hit MySQL DB for
every req. → non-isem
we need to hit DB to
get user & problem info
anyway.
do stuff in 1 db call!

CASE STUDY



Contest Leaderboard

300,000 participants

5 problems

3 hour duration

? How can we compute the ranklist?

User's submission → every time a submission happens, the true
+ve ↗ ✓
-ve ↗ penalty
Rankings have changed

? How frequently does the ranklist change?

140 changes/sec

Estimate # submissions

avg per sec → 5

total = 1.5M

$$\text{qps} \Rightarrow \frac{1.5M}{3 * 3600} = 140$$

Each submission takes 10sec to eval

$$\rightarrow \# \text{ eval needed} = 1400$$

$$\frac{\# \text{ sec}}{\text{time for eval}} = \# \text{ eval sec}$$

$$\frac{1400}{10} = 140$$

? Is immediate consistency important?

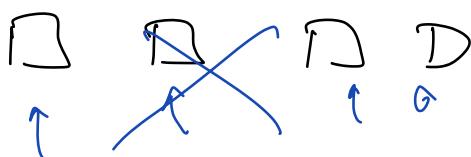
Eventual is fine

↳ max delay that will be okay (30 min @ scalar)

5 - 30 mins



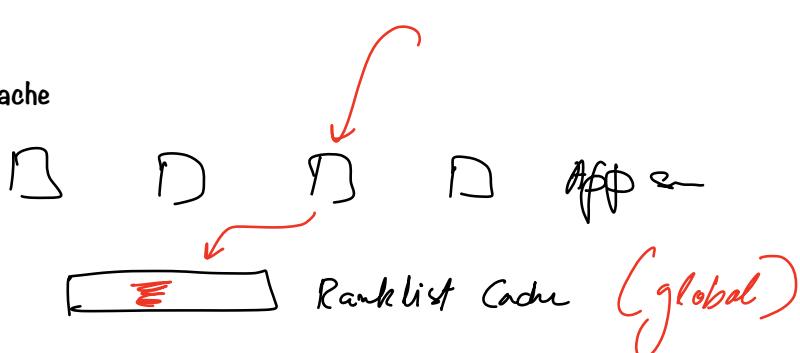
Local cache



TTL / Write Around



Global cache





Redis

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

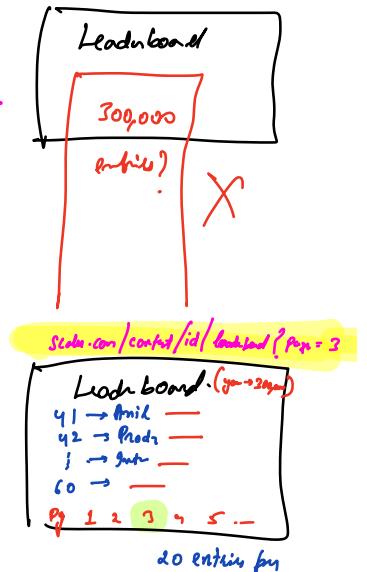
100K req/sec

Key - Value

- ① Very powerful → adv. data structures → sets / arrays / sorted sets / bloom filter
- ② Ultra fast
 - Coded in C
 - Single threaded
 - async io
 - multi-threading
 - Locking.
- ③ in-memory

My Score & Rank

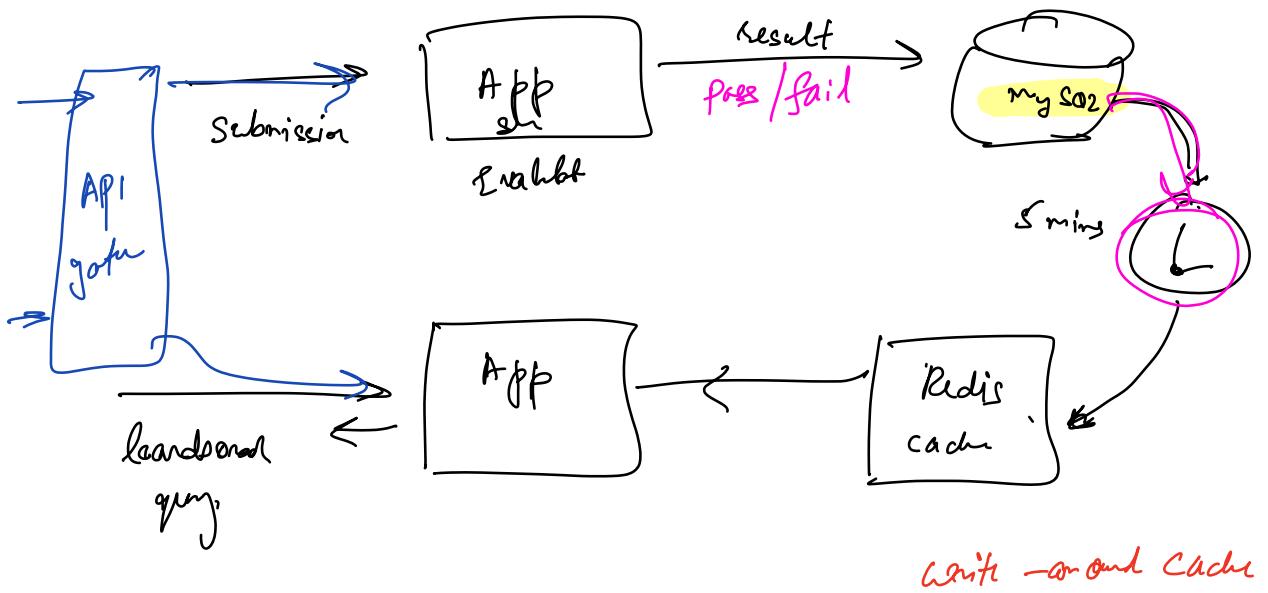
Key	Value
100k+ score	{User-id = 1}
300m score	{User-id = 1}
User-id = 1	{Rank; score; P1 → P2 → P3}
User-id = 2	{Rank; score; P1 → P2 → P3}
User-id = 3	{User-id = ---}
;	{User-id = ---}
Rank = 300,000	---



query (all keys b/w rank - 000 001 to rank - 000 020)
Pg 1

Pg 2 → rank - 000 021 to rank - 000 040

$$600,000 * 110B = 66 \text{ Mbytes}$$



300,000 users,

1 hour contest

On avg → each user might look

at ranklist $\Rightarrow 5$ times.

$$\frac{300,000 \times 5}{1 \text{ hour}} = \boxed{140 \text{ qps}}$$

10:48 → 10:55

CASE STUDY

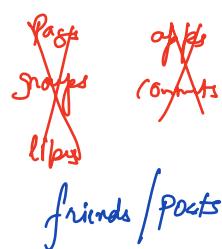


Facebook News Feed

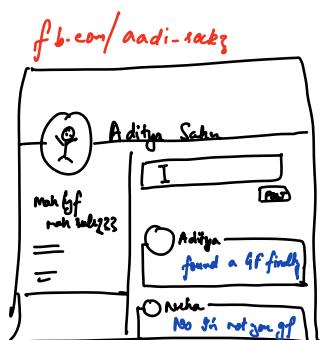
? How can we compute the News Feed?

DB schema (Relational)

Users		friends		User-posts		
id	Name	user-id	friend-id	post-id	profile-id	Content
1	Arijit	1	3	3	3	I have a gf
2	Neha	3	1	2	3	3f
3	Aditya	2	3	1	1	No gf

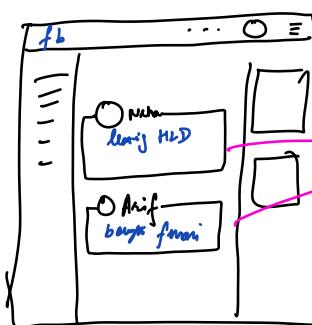


Profile Page



any post made by
Aditya / his friends
on Aditya's profile

News Feed



recent & popular
posts made by
friends of Aditya
on their own
profiles on other
people's profiles

Select * from user-posts
where profile-id = 3 → Aditya

Order by created-at desc

limit 20

- ① get friends list of Aditya
- ② join on user-post table to find any post on any friend's profile
- ③ Score
 - Recency
 - ↳ Popularity
 - ↳ Demography
 - ↳ Sort by user-affinity / perform

?

Does all the user data fit on a single machine?
 ↗ posts / friends lists
 $\text{fb} \rightarrow 2B^+ \text{ users}$

No!!

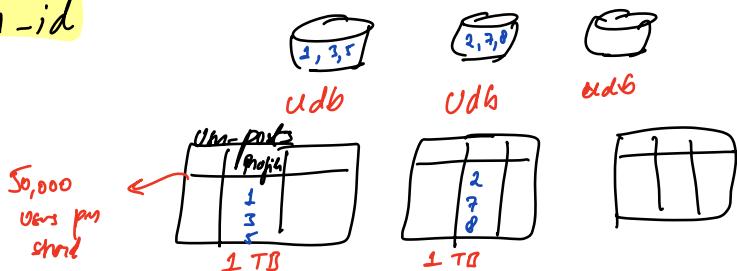
∴ shard the data

?

What is my sharding key?

↳ ① all the fragment queries must ideally hit 1 shard ✓
 2 shards ✓
 n shards ⚡ bad

User-id



Profile → exactly

1 shard

∴ all posts made on
 a given user profile
 ↳ inside that user's
 shard.

?

How do I retrieve data for the news feed?

We have to hit $\frac{\text{all}}{n}$ shards!!

worst case → if Aditya has n friends

we might have to hit n shards

fan-out

limited by slowest shard



$1 \leq n$ } fanout

?

How can you optimize the construction of the News Feed?

Potential Solutions

→ diff sharding key

→ pre-calculate News feed for each user & cache it

↳ ① estimate amount of data

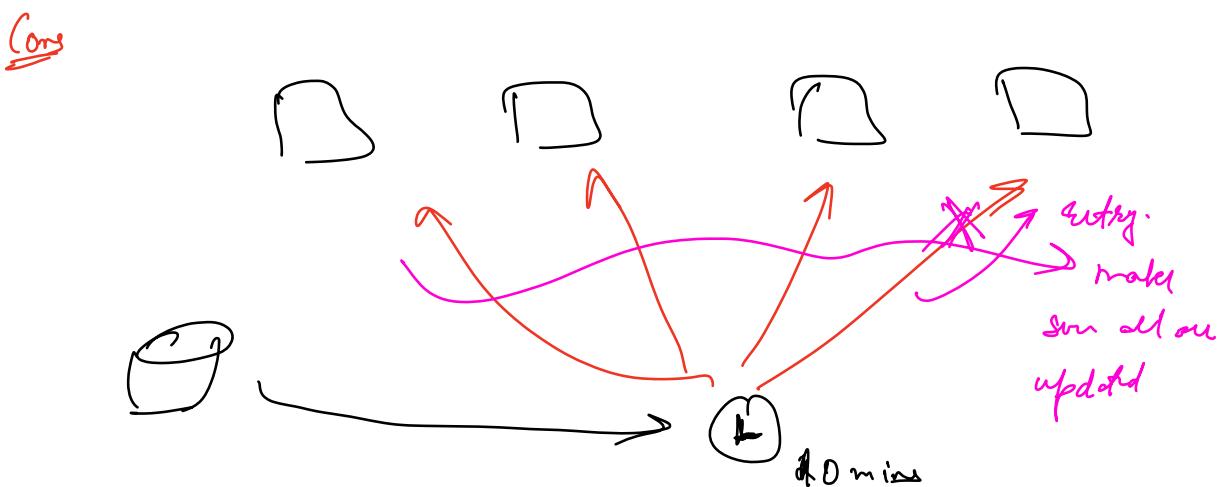
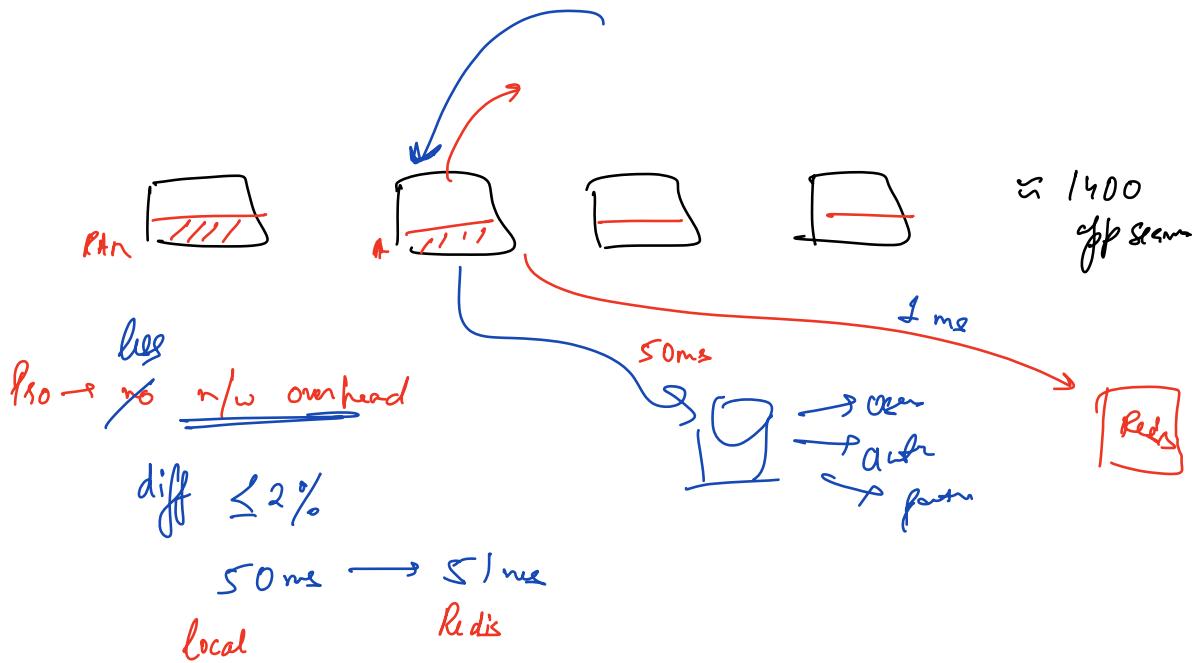
↳ ② how will you update the newsfeed

→ on each post → push to each friend's news feed algorithm.

→ person ^{most} all friends are in same shard → impossible

↳ clustering → hard problem

↳ fan-out unites



Calculation

assume 1 server can handle 1 req at any given time

N servers

Each server takes 10 seconds

10 secs \rightarrow N requests get served.

1 sec \rightarrow $\frac{N}{10}$ req

$$\text{1 sec} \rightarrow \frac{N}{10} \text{ req}$$

$\frac{140 * 10}{N}$

req

Unitary Method

10th class

$$\frac{140 * 10}{N} = 1 \Rightarrow N = 1400 \text{ servers.}$$

Evaluating a submission

I/O Bound
2sec / 0.5sec

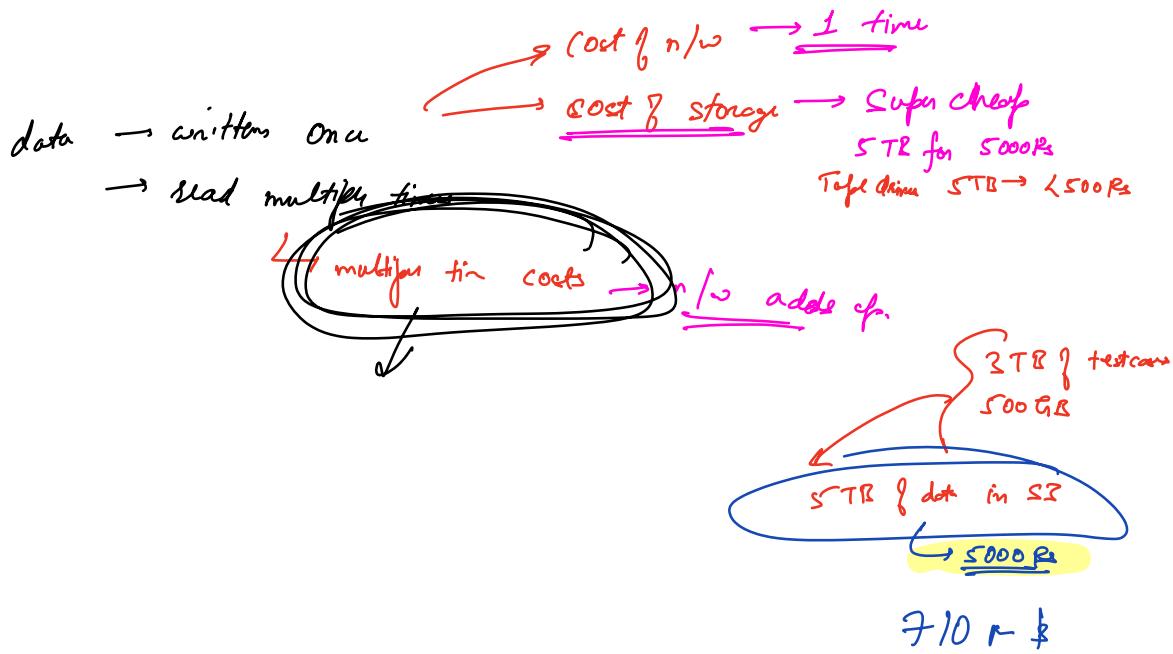
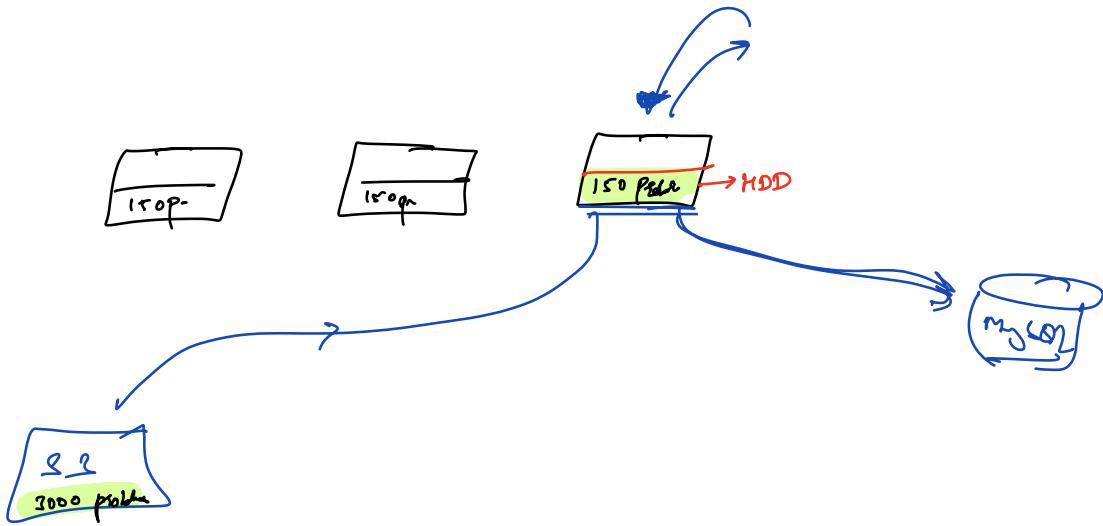
CPU bound

10sec \rightarrow pure CPU time

you should never multi-thread a CPU bound

process (OS Lecture)
Threading

leaderboard server \rightarrow I/O bound



Modern Hardware

Typical Server

200 GB ↔ 2TB HDD

2GB ↔ 64 GB RAM

2 ↔ 16 core

100 Mbps ↔ 10 Gbps n/w

Most Expensive

Server

5 PB of HDD

12TB of RAM

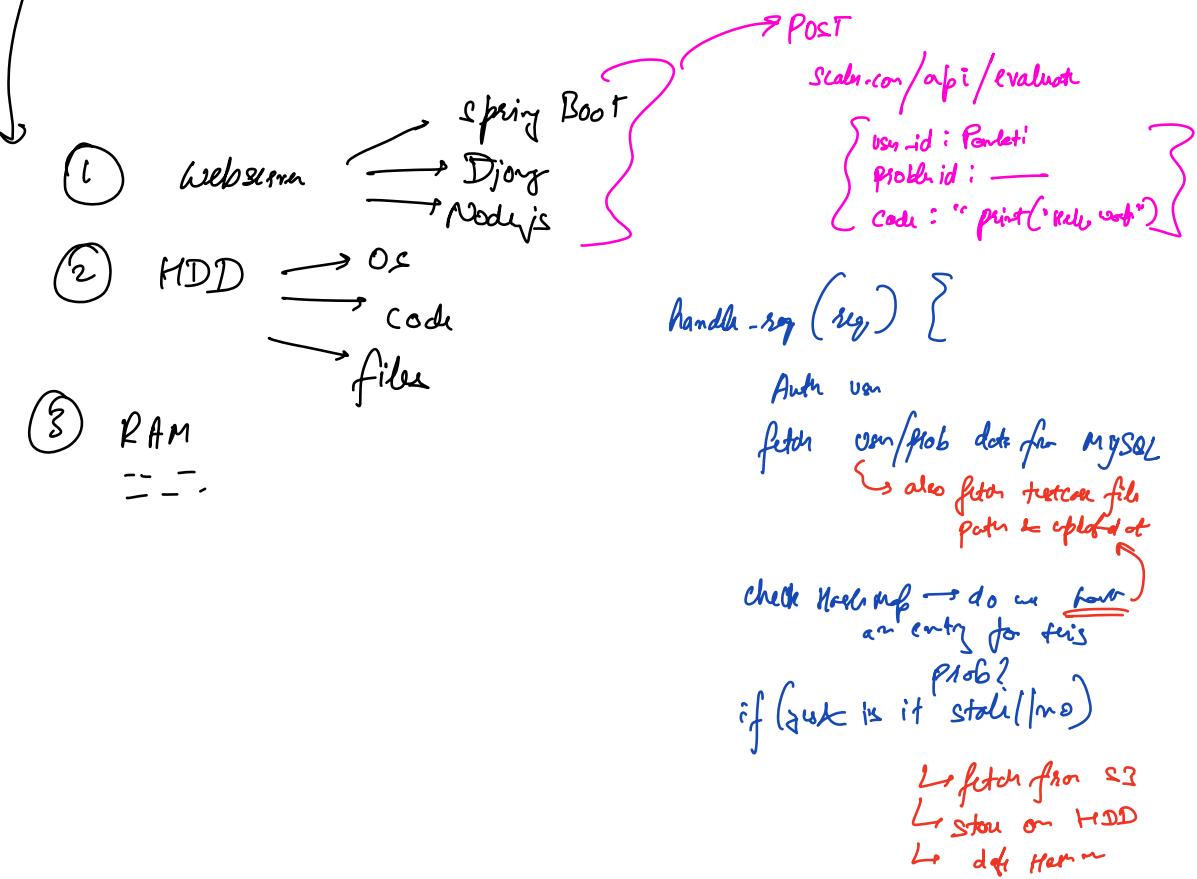
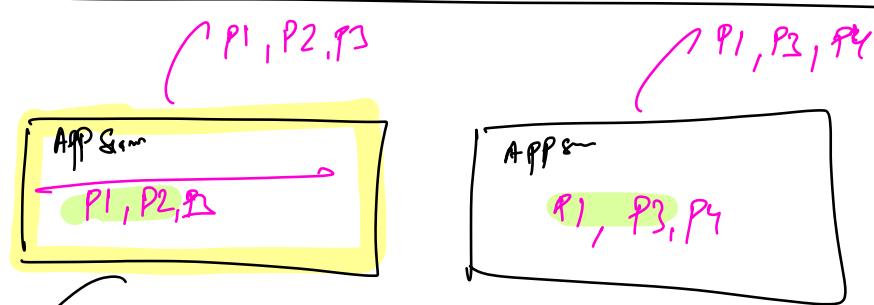
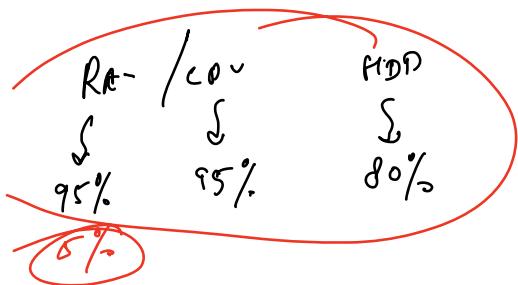
600 core

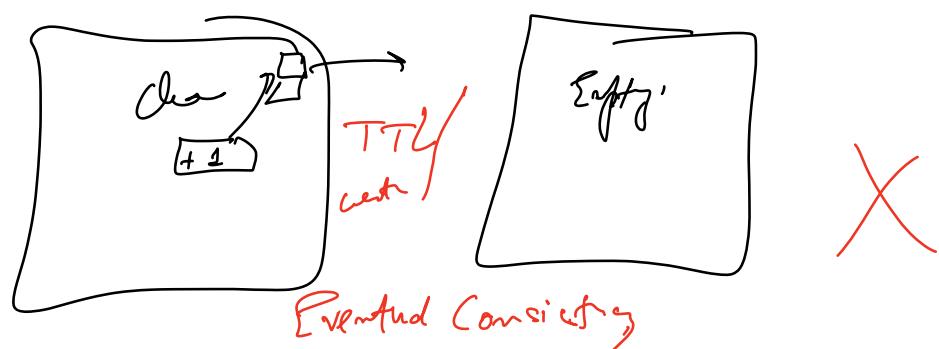
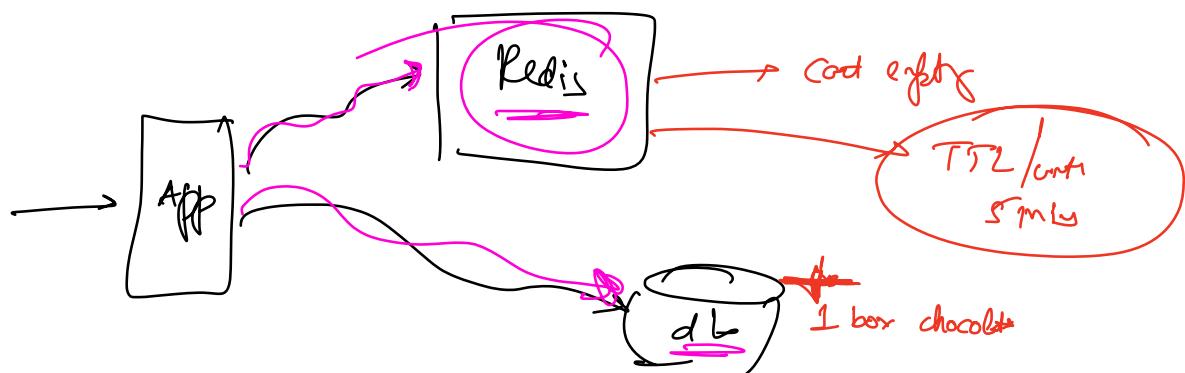
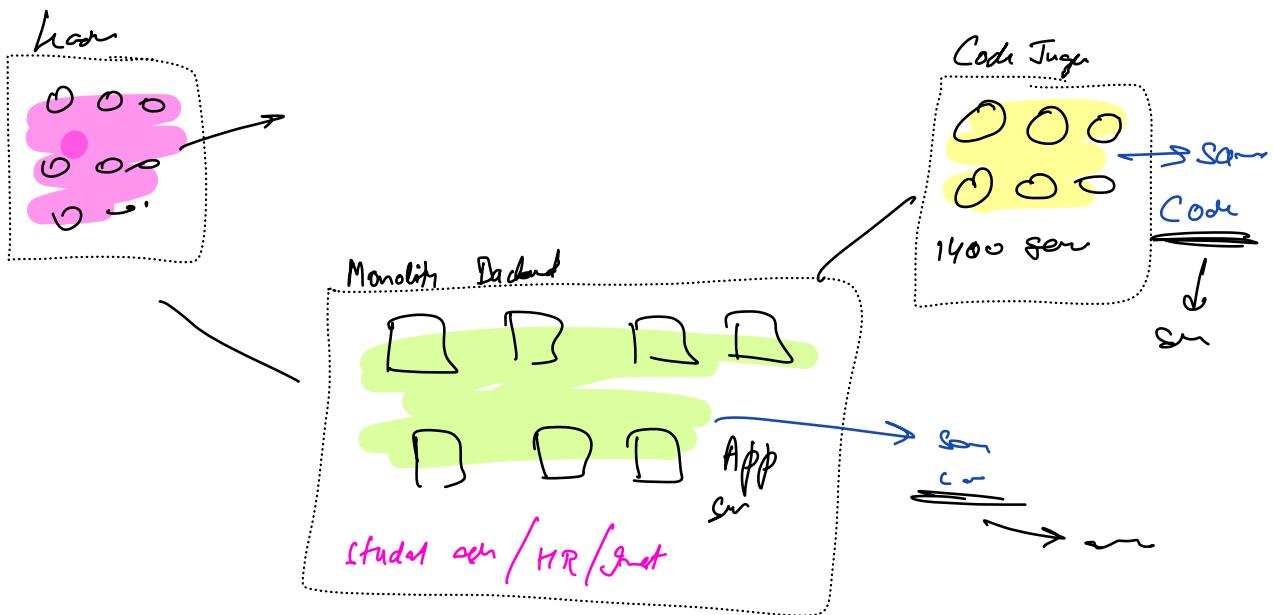
100 Gbps n/w

App server → 100 GB q in ~ Hashmap Cache

Redis server → 60 GB → cache

q in R → OS / - -





Immediat Cons → write though

11. 11.

4



Cache user -> news feed

3

ever
refine

Next
close

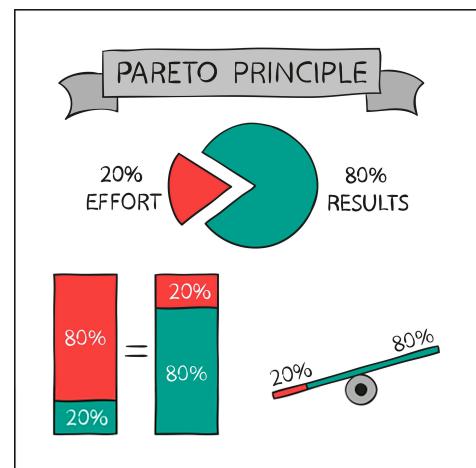
1.

2.

3.



Back of the envelope estimates



? Daily post data

? Monthly post data

