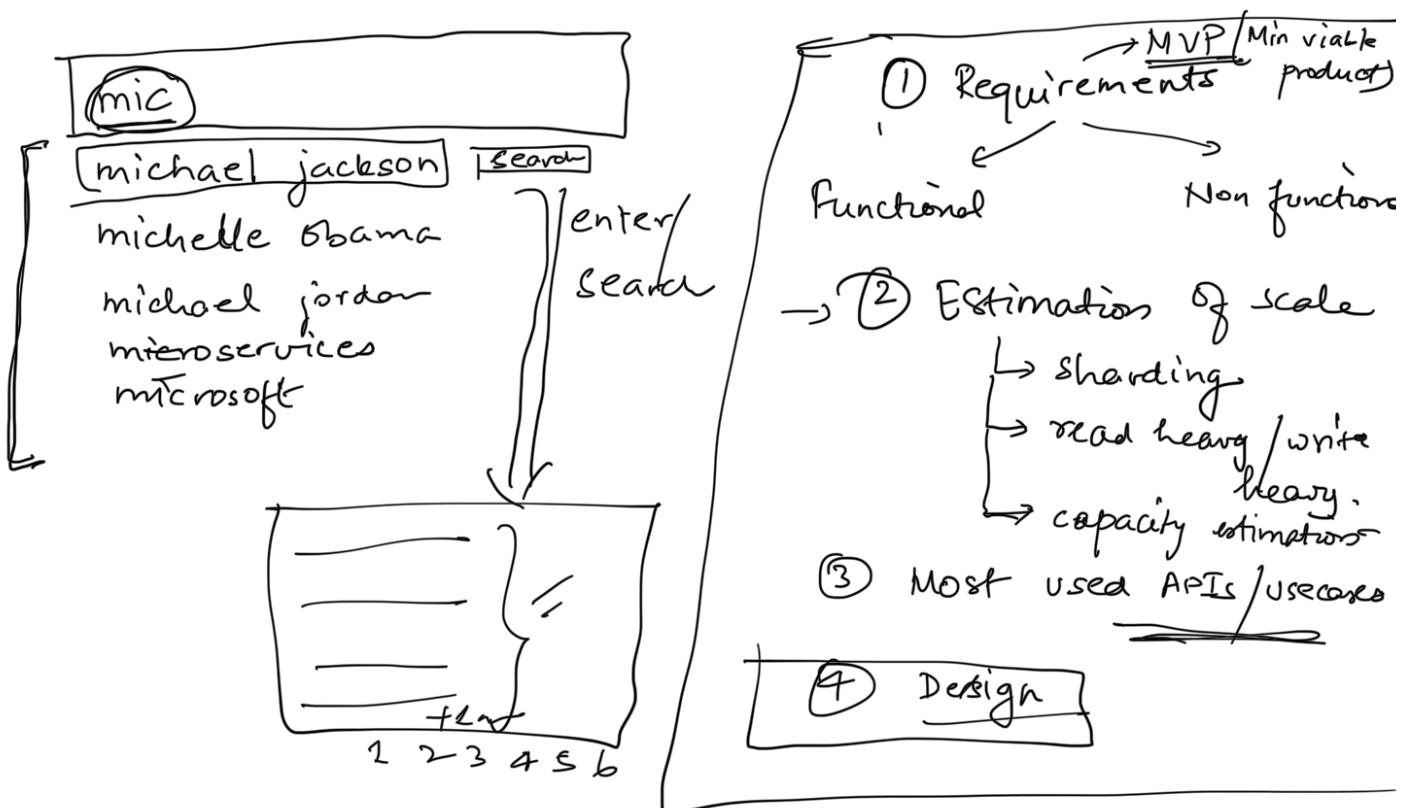
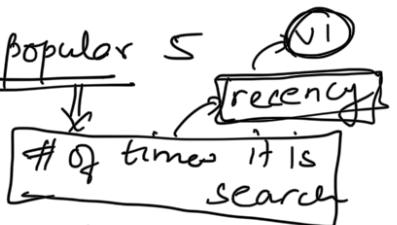


DESIGN GOOGLE TYPEAHEAD



① MVP

- ① how many suggestions (max) → 5 suggestions.
- ② Ordering/ choosing → most popular 5 
- ③ personalise (location, system, person based)
- ⇒ ④ suggestions should have typed text as strict prefix in VD → NOT NEEDED
- ⑤ 3 letters before showing suggestions.
- ⑥ spelling mistakes → NOT NEEDED
- ⑦ case sensitivity → case insensitive FOR VD

Non-functional requirement

Consistency →

- ① Highly available

mic - —

② Latency → very low latency

Estimation of scale

Search term → frequency

Searches → 8 billion / day.

8 Billion new search terms / day.
10 ↓

$$\frac{0.8 \text{ B} * 365 * 10 \text{ terms}}{10 \text{ years}} * 100 \text{ bytes} = 86 * 365 * 100$$

$$= \underline{\sim 300 \text{ TB}}$$

8 Billion queries per day

8 Billion / 10 new queries / day

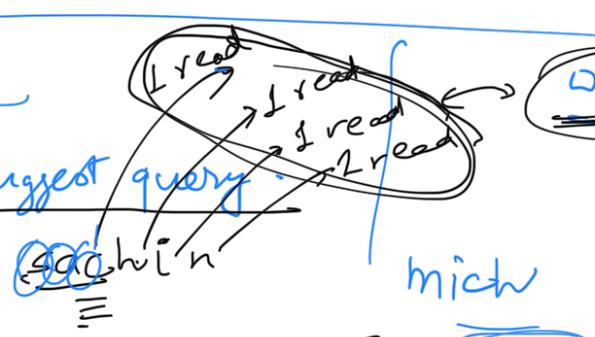
8 Billion / 10 * 365 * 10 new search terms.

$$\Rightarrow \text{bytes} \Rightarrow 8 \text{ B} * 365 * 100 \text{ bytes}$$

$$= 8 \text{ GB} * 365 * 100 \approx \underline{\sim 300 \text{ TB}}$$

Read

every autosuggest query



every search



Sachin tendulkar

[Enter] 



① Traffic → queries per second

↓
9:30 - 11:30

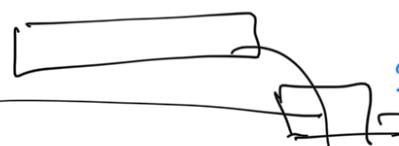
8 Billion searches / day
+ writes.

16 Billion reads / day

$$\underline{24 \text{ Billion}} = \frac{\underline{24 \text{ Billion QPS}}}{24 * 60 * 60}$$

$$= \frac{\cancel{1000 * 1000 * 1000}}{3600}$$

$\approx 250,000 \text{ queries/sec}$



$\approx 250,000 \text{ queries/sec}$

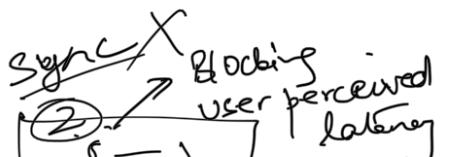
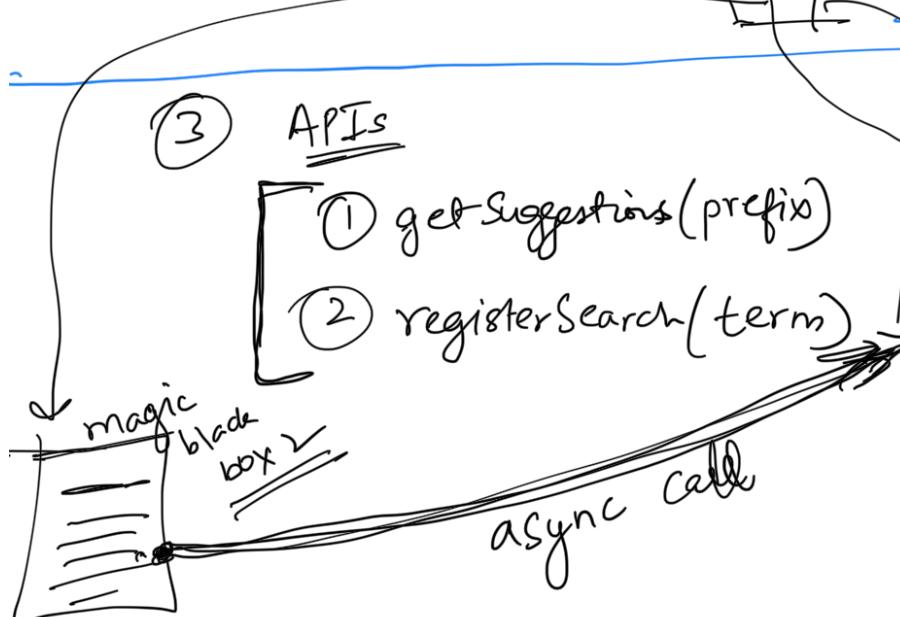
③

APIs

- ① getSuggestions(prefix)
- ② registerSearch(term)

getTopSuggestions(prefix)

Magic
Black
Box



User



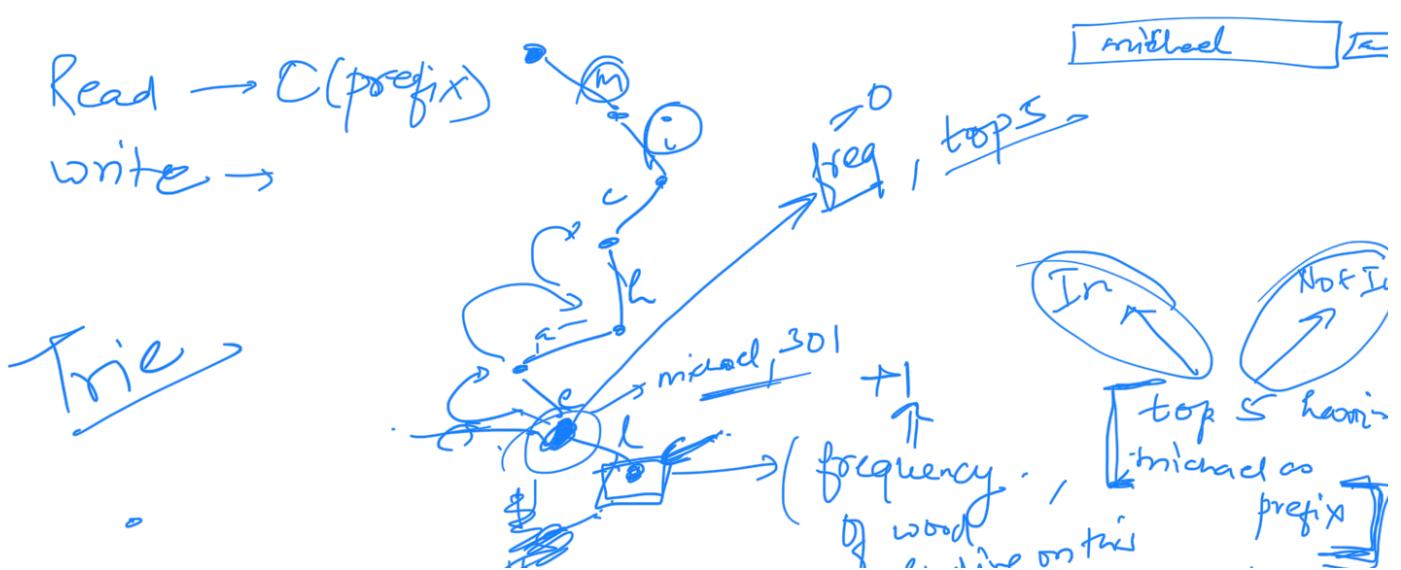
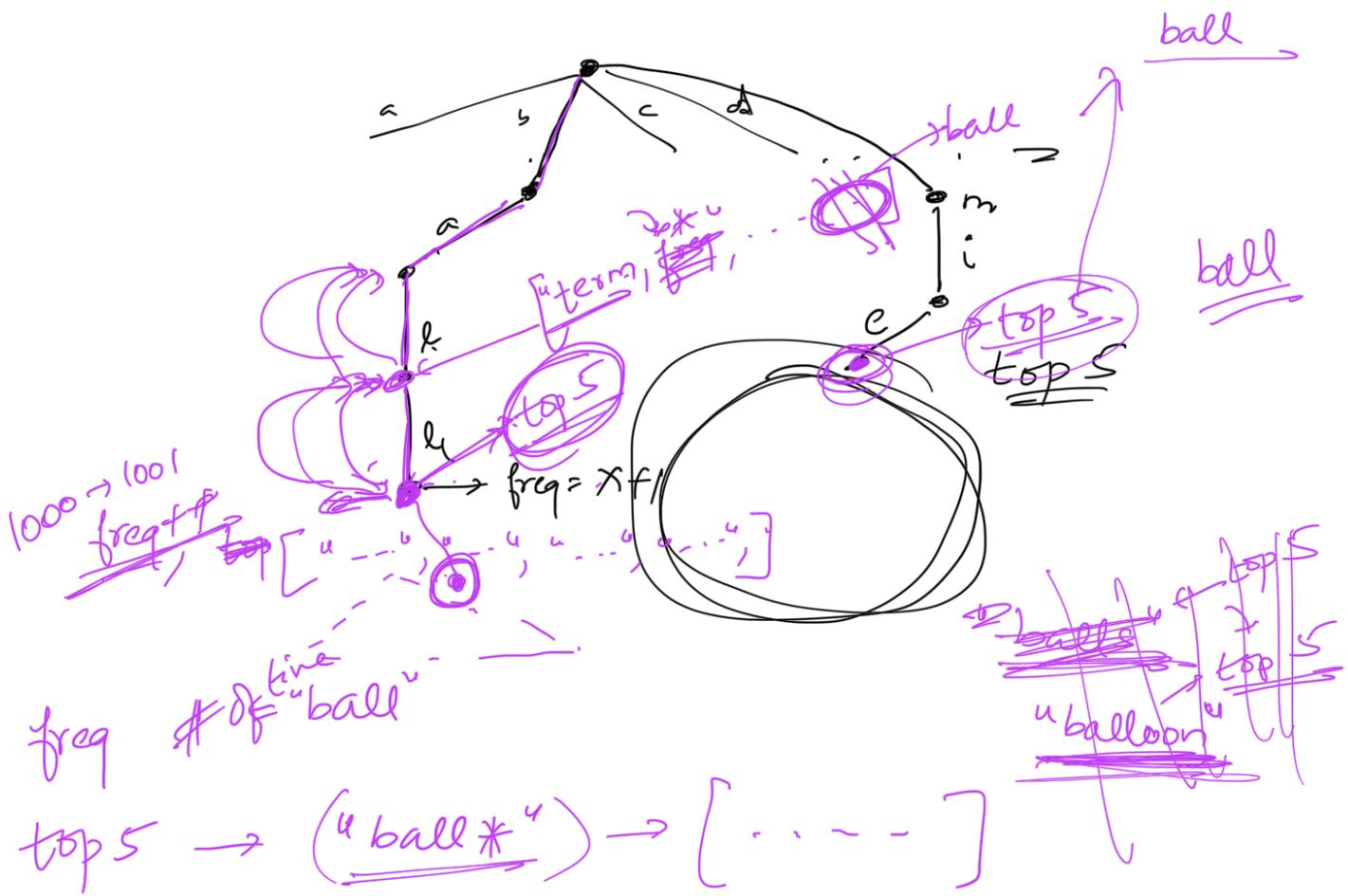
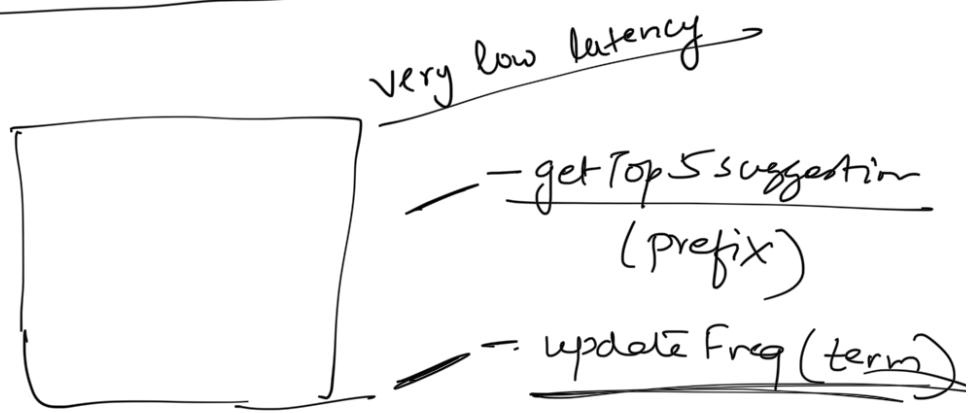
Asynch

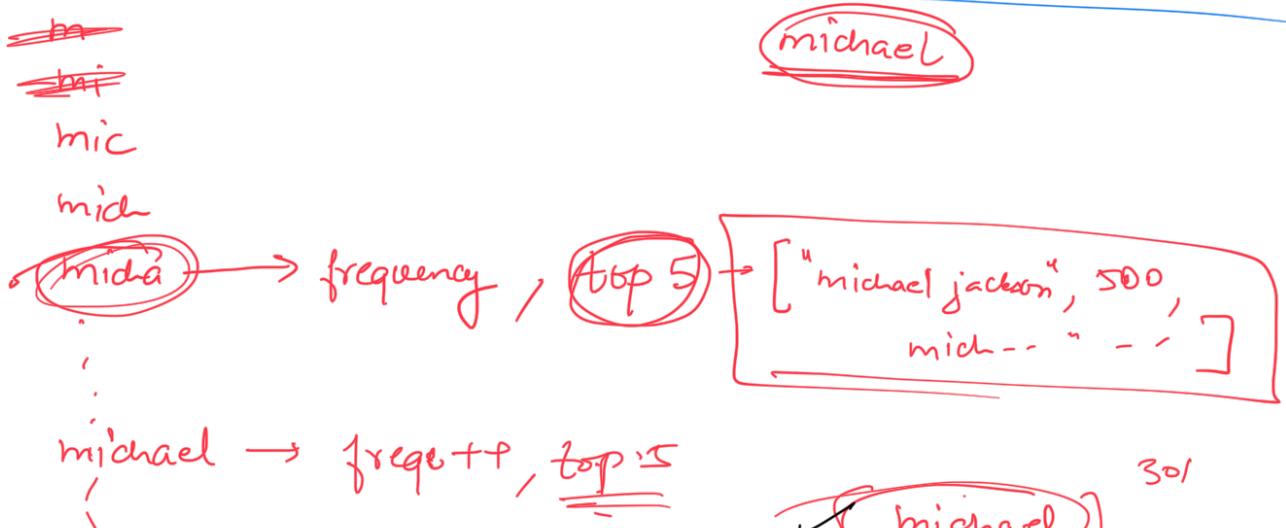
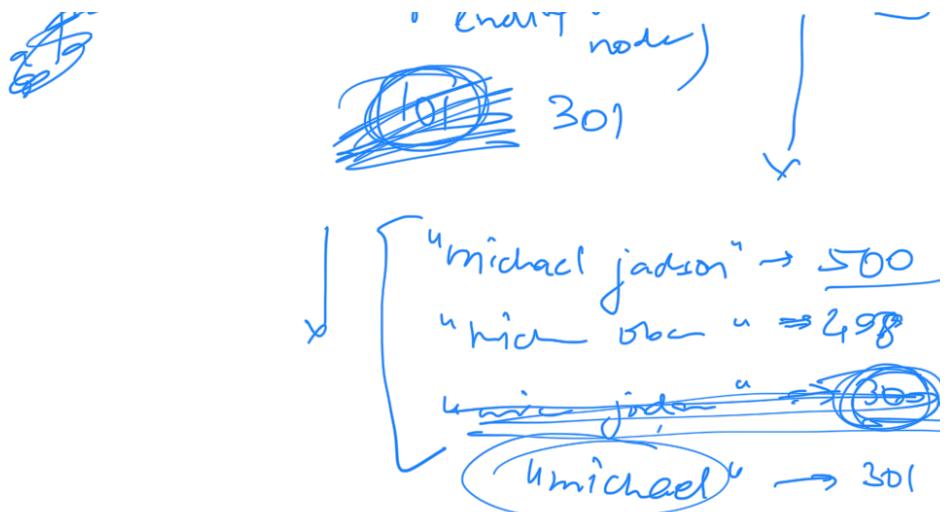
- Non Blocking,

- ~~Doesn't add to~~ user later

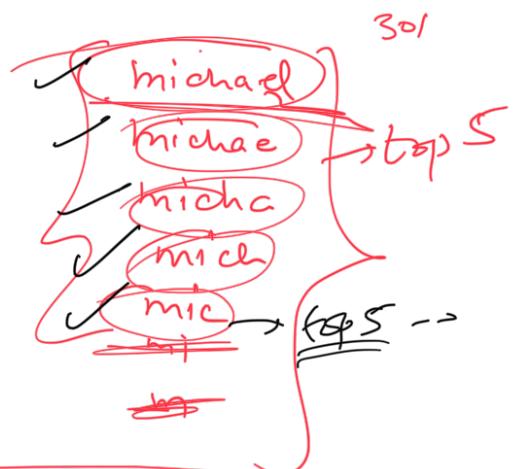
- Background





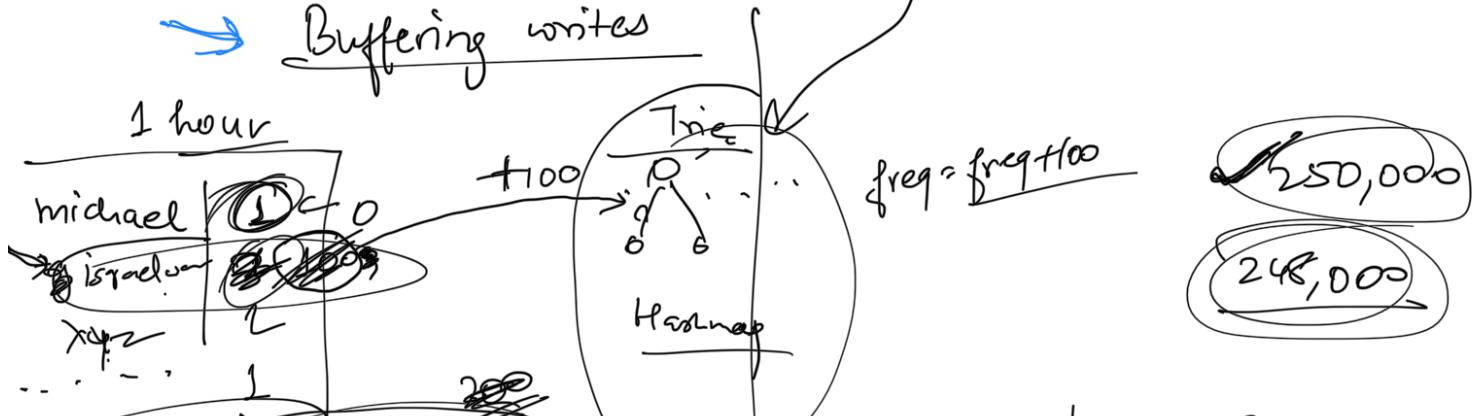


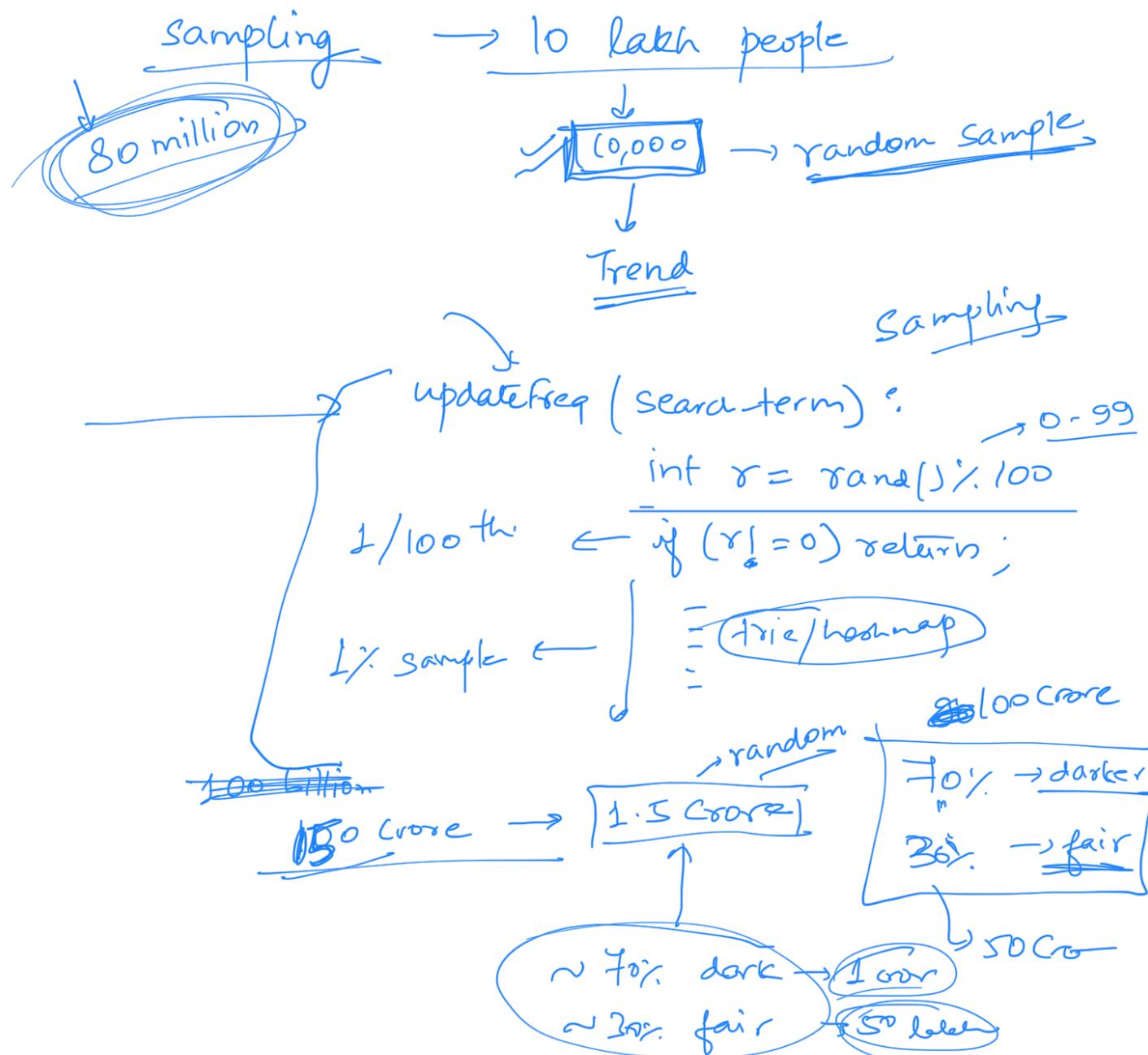
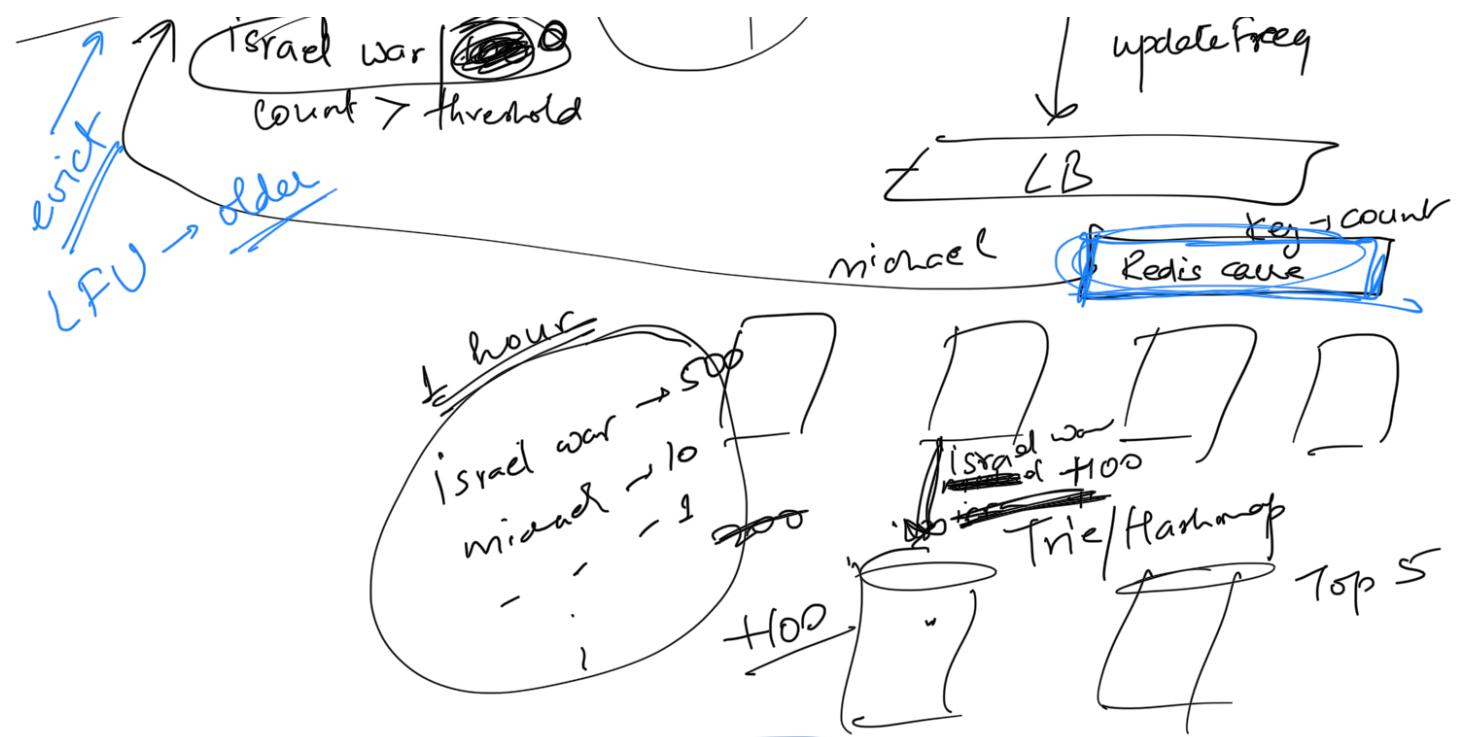
~~getTop5 (key):~~



① read & write heavy system

→ Buffering writes





① Read/write ratio

② Trie/Hashmap → read fast

③ Shard

④ Recency

Hashmap

mic
read

Key

→ freq, top 5

mic

mic
freq
top 5

me
LT

mic
LT

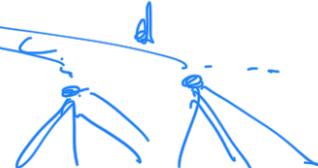
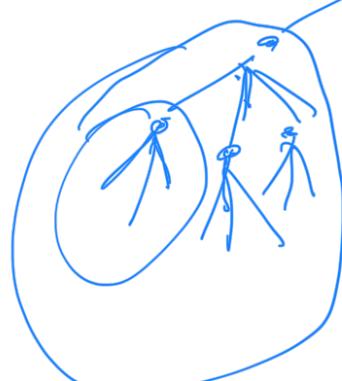
300 TB

→ 300-500 shards

Trie

26

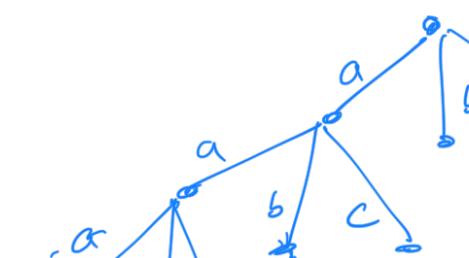
(X)
fix



LT

LT

mic



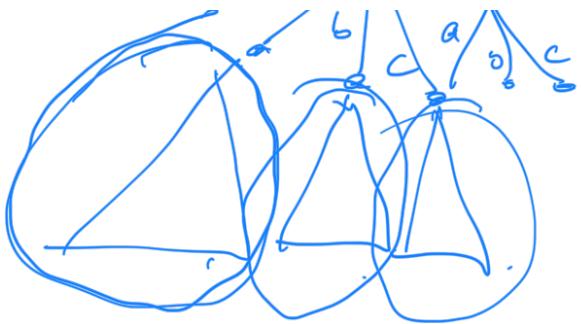
aaa

XXX

mic

shc

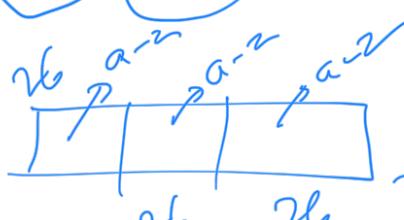
26 * 26 * 28



' - - - - -

$\Rightarrow 17,521$

Consistent here



$26 \times 26 \times 26$

$\approx 8K$

~~200~~

~~350~~

~~175~~

"aaa"

"aab"
"aac"

"nad" ---

israel war $\rightarrow 100$ recency

freq + recency

ishani word $\rightarrow 500$

~~100~~ 2

Time decay

≥ 1

~~tdf = 1.1~~

1.1

$12.5 \dots$

$day 100$

$israel \rightarrow 100$

$ion \rightarrow 2$

day 1

michael $\rightarrow 100$

ishani word $\rightarrow 10$

$\rightarrow 50$

day 2

min $\rightarrow 2.5$

in $\rightarrow 1.25$

$\rightarrow 2.5$

day 3

~~50~~

~~2.5~~

$\rightarrow 1.25$

day 4

~~1.25~~

~~0.625~~

:

~~tdf = 1.1~~

10000

$\frac{1}{1.1} \approx 10$

day 1

michael $\rightarrow 10,000$

day 2

$\rightarrow 10,000$

$\rightarrow 1/2 \rightarrow 1/2 \rightarrow \dots$

4000

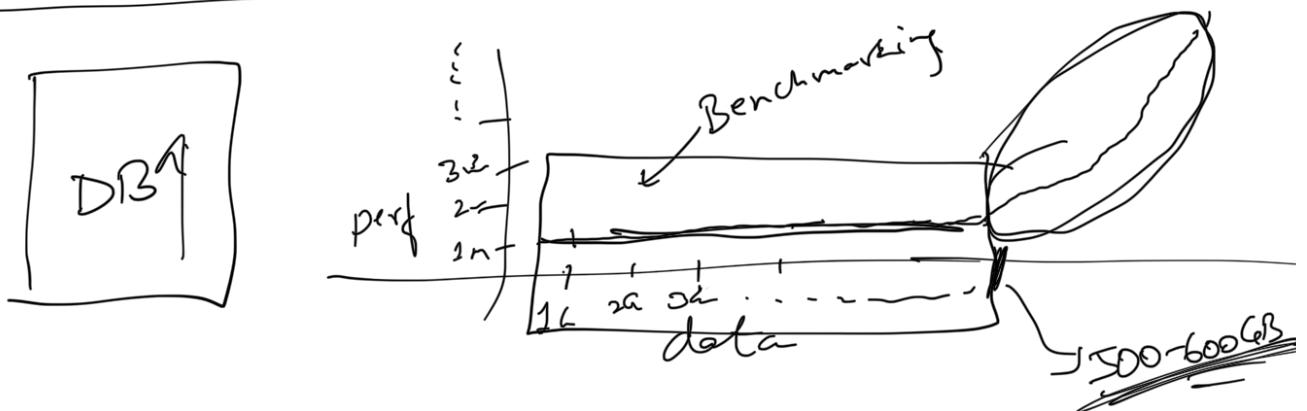
$michael \rightarrow 4000$

$10 \rightarrow \dots$

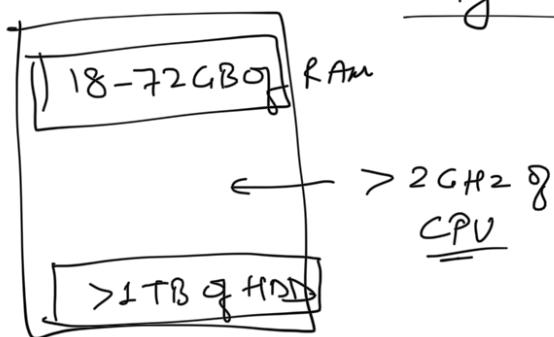
day

$$C_{\text{cur-day}} + \frac{C_{\text{prev-day}}}{td_f^{10}} = \frac{10,000}{\sim 100} \approx 100$$

$$+ \frac{C_{\text{prev-prev}}}{td_f^2} + \frac{C_{d-3}}{td_f^3} + \frac{C_{d-4}}{td_f^4} + \dots$$



~~Redis~~
~~Kafka~~ → microservices



BigTable | DynamoDB
↓
consistent hashing

AWS → EC2 machine



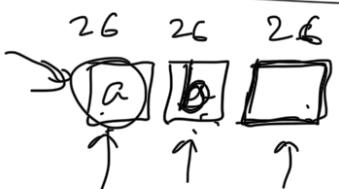
Relational DBs → Tries

Oracle DB

aaa, aab, aac, aad, -----



↓



a-2 a-2 a-2

abd

18000 user



CH

$$\begin{aligned}
 & 26 + 26 + 26 \\
 & + \dots \\
 & 26 * 26 \\
 & + 26 * 26 \\
 & + \dots \\
 & = 26 * 26 * 26 \\
 & \boxed{\approx 18000}
 \end{aligned}$$



(00 share
....)



abc



aba

> 1 subtree
180 Subtrees.

18000
100
= 180