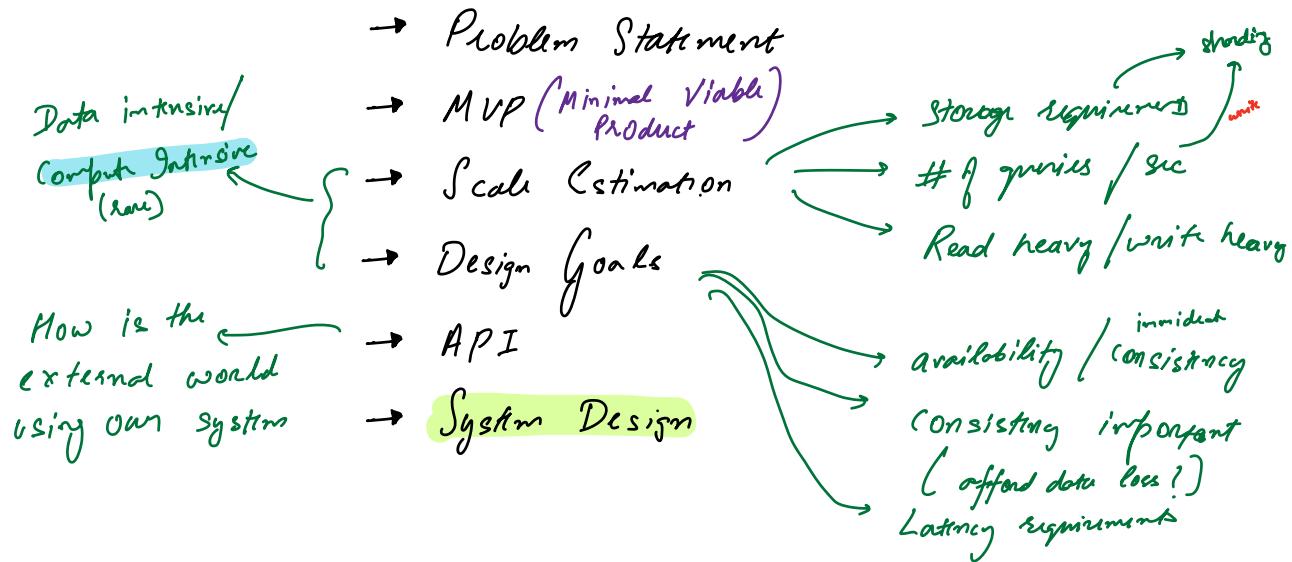


Agenda

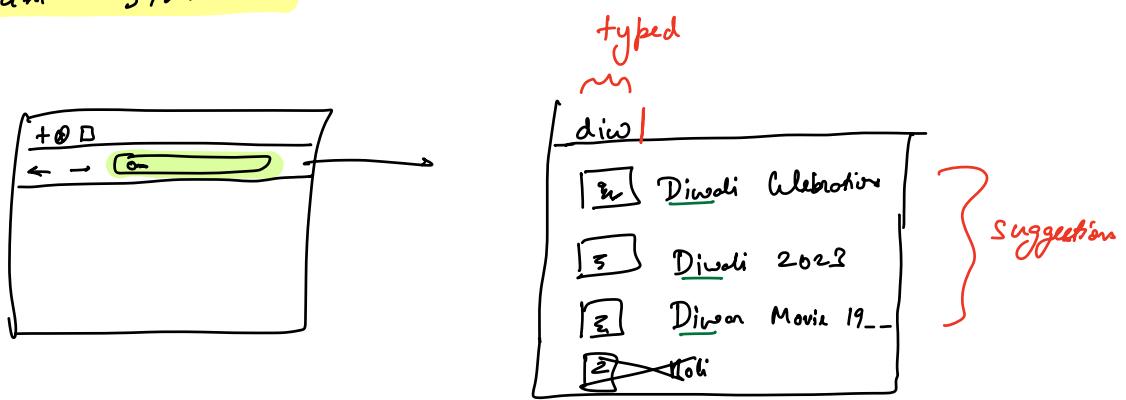


Case Study 1 → del.icio.us

PACELC
P → Availability
A → Consistency
C → Consistency
E → Latency
L → Consistency

Design Interviews → free flowing → discussion based

① Problem Statement



How do you design a faster ahead system

Scale: Google

Engineering Architect @ Google → Search

② Minimal Viable Product

① Maximum # of suggestions

5 is more than enough

② What should we suggest

↳ popular searches

③ Personalization

not important for MVP

④ If no relevant results, should we show some (inrelevant) results anyway?

No → show only relevant results

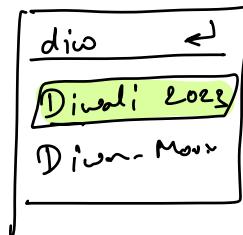
⑤ Fuzzy matches / spelling errors

→ not needed for MVP

→ strict **prefix match**
(begins with)

⑥ Rank results

Popularity of search queries
(diff from typed-in queries)



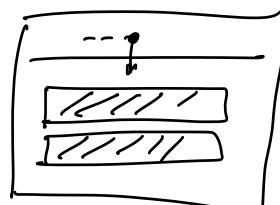
geo-personalized
not for MVP

Show more recent results on top?
not in MVP

⑦ Min number of letters that must be typed before we show typeahead

3

⑧ Show suggestion for every keystroke

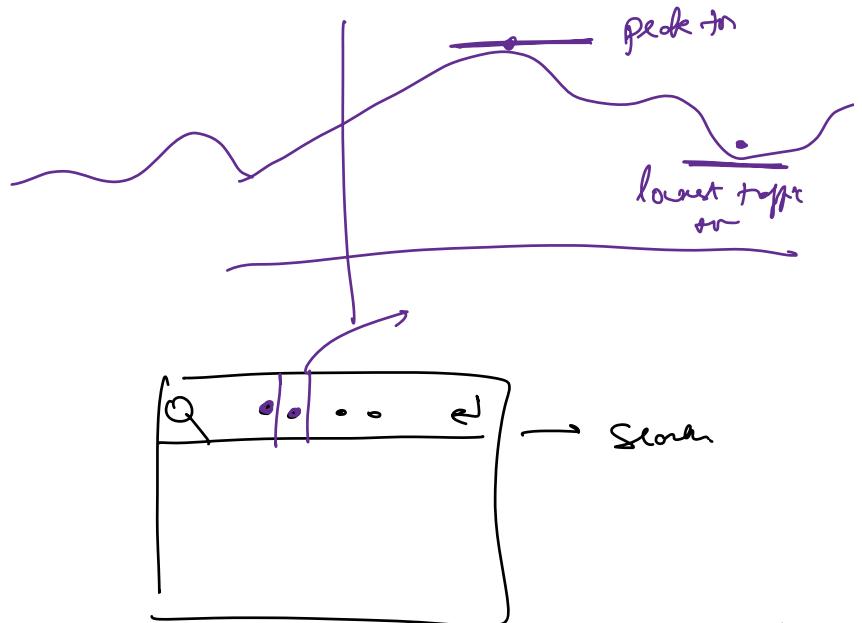


③ Scale Estimation

Volume of queries

→ ≈ 10 Billion Search queries per day

$$\left(\frac{10^{10}}{86,400} \rightarrow \frac{10^{10}}{10^4} \approx 10^5 \text{ queries/sec} \right)$$



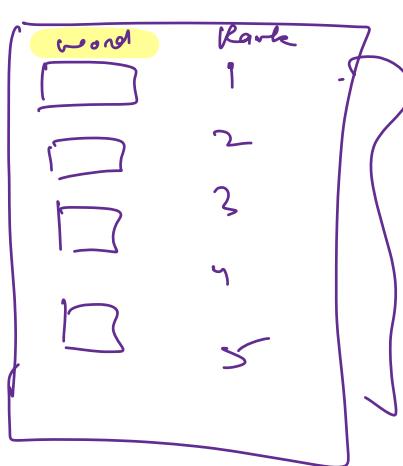
→ 7 typeahead queries
for each search query.

→ 7×10^5 typeahead
queries / sec

→ 1 million typeah / sec



Scalability of data → do we need sharding



10⁸ per sec

32 letters	8B
Search queries	# of hits
Dicordi	2 B/M
Rogn Au	100 M
Holi	800M
Diven Mow	2 M/M

$$32^8 \cdot 2B + 8B \\ \leq 100B \text{ per sec}$$

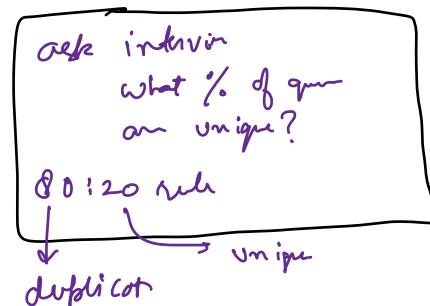
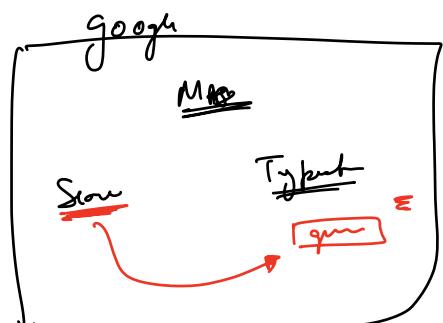
10⁸ searches per day
 ↳ # of unique queries per day
 ↳ never been searched before
 ≈ 15% queries are unique

10 years of Post data

$$\cancel{10^{10}} + 365 \cdot 10^8 \cdot 100B$$

$$1.5 \times 10^9$$

$$1.5 \times 365 \times 10^{12}$$



$$1.5 \times 2.65 \times 10^{14}$$

$$6 \times 10^{14} \rightarrow 10^{15} \text{ Bytes}$$

$$\begin{aligned}\text{Giga} &\rightarrow 10^9 \\ \text{Tera} &\rightarrow 10^{12} \\ \text{Peta} &\rightarrow 10^{15}\end{aligned}$$

$\approx 1 \text{ PB}$

$\approx 1000 \text{ TB}$ of data

is sharding needed?

- ① data is too large
- ② worth partitioning
many

Yes!!

Josh
Linux Tech Tips

100TB server.

Read heavy / write heavy

10B scan per day \rightarrow units $\Rightarrow 10^5 \text{ sec/sec}$

$\Rightarrow 10^5 \text{ units/sec}$

$\Rightarrow 10^6 \text{ reads/sec}$

Both read & write heavy

④

Design Goals

→ Availability vs Consistency



eventual consistency is okay



→ do we need consistency at all? **No!**
(can we deal with data loss)

We can afford small
data loss

actual → 800M
second → 799 M hits

→ Consistency vs Latency

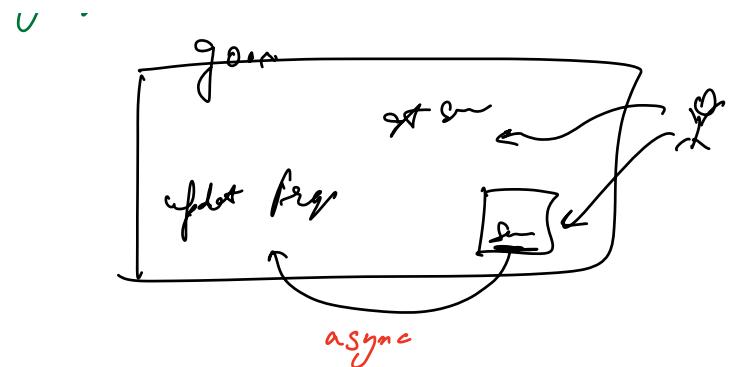
Latency should be very low

⑤ API

→ how will outside world use our system
anything except our service → user
at google → other services

① get Typeahead Suggestion (partial-query)
user | upto 5 results

② update proxy (search-query)
google's instant search

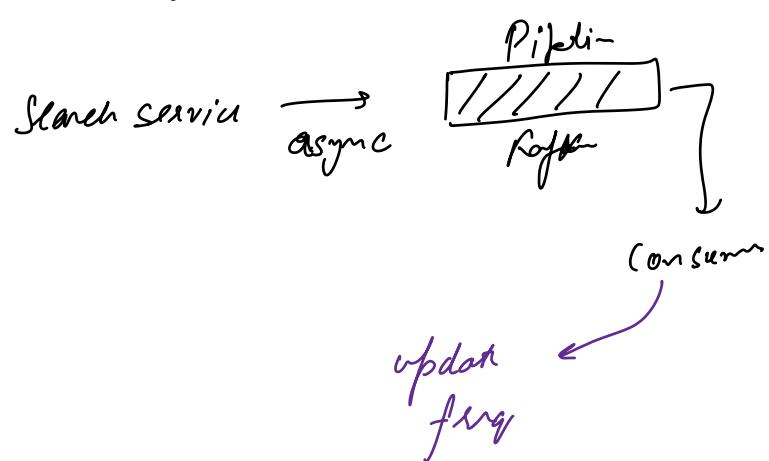


10:42 → 10:50

Designing Typohead

Update Freq (short-term)

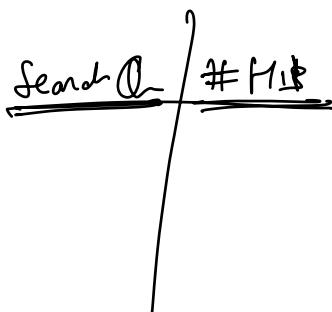
X



Key - Value

center → Maintain

read →



X won't work
for read queries

get Typoahead Suggestion (Prefix) → top 5 results
avg 10 letters ~ 10

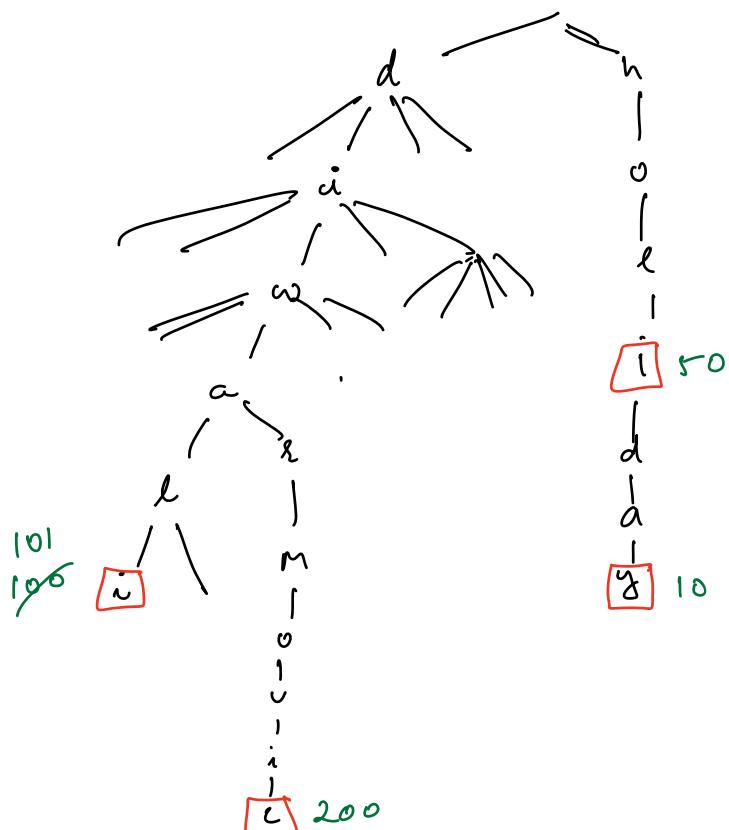
$$(26+10)^{10}$$

36¹⁰

$$\begin{array}{r} 36 \\ \times 36 \\ \hline 10 \end{array}$$

True Approach

Construct a tree of slender species



node = children [36]
is terminal
hits
for 5 results
with counts

get TBB suggest (Partial)
diagram

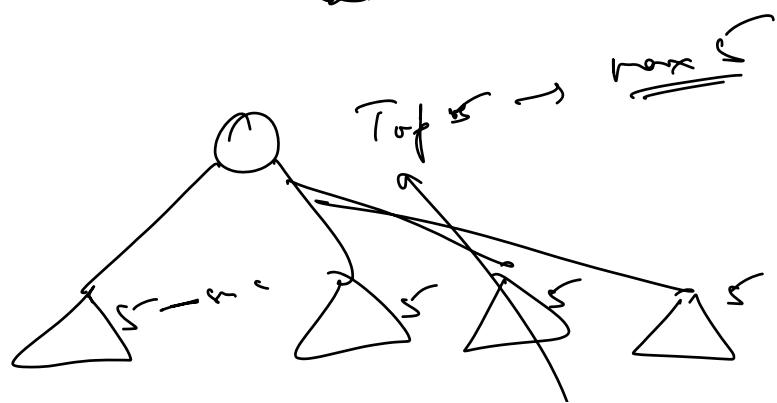
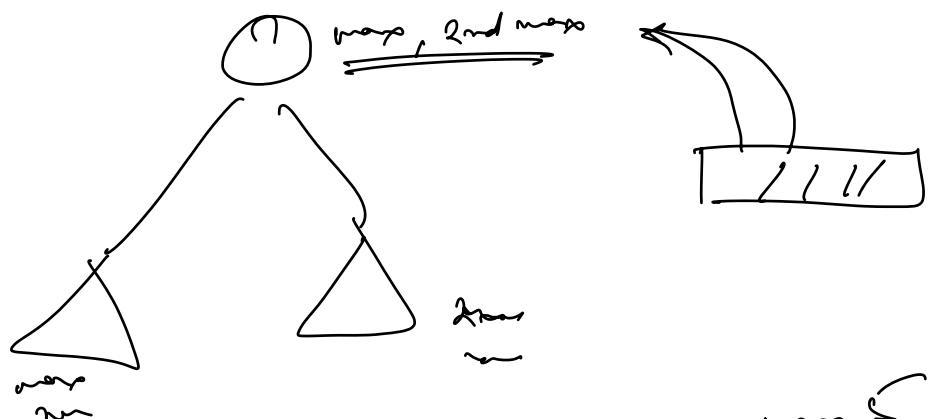
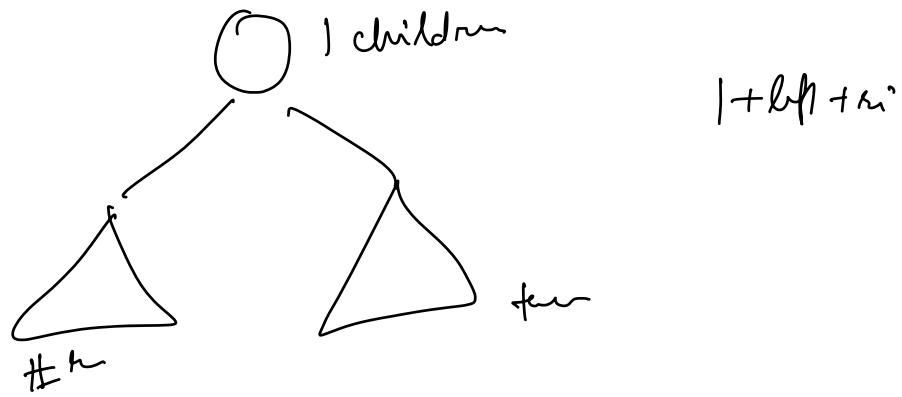
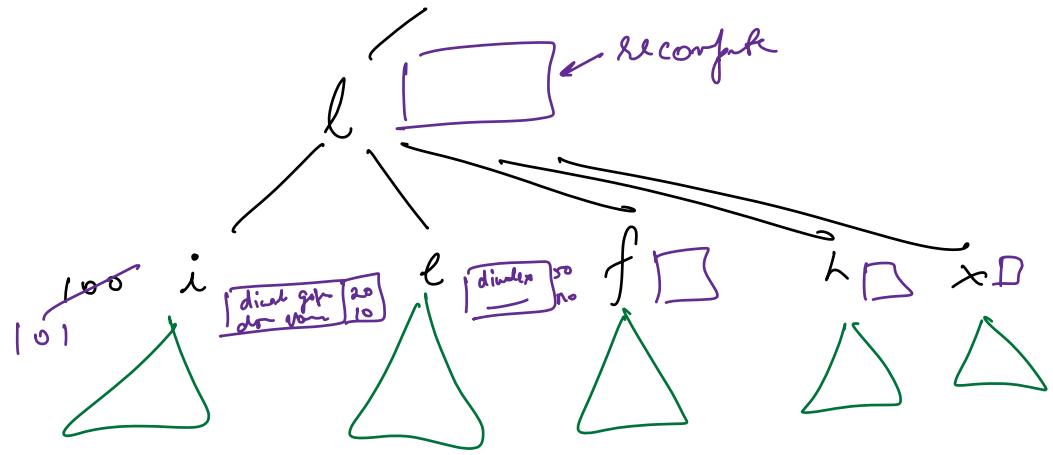
Wfdati frgr (Stora-far)
diseasi

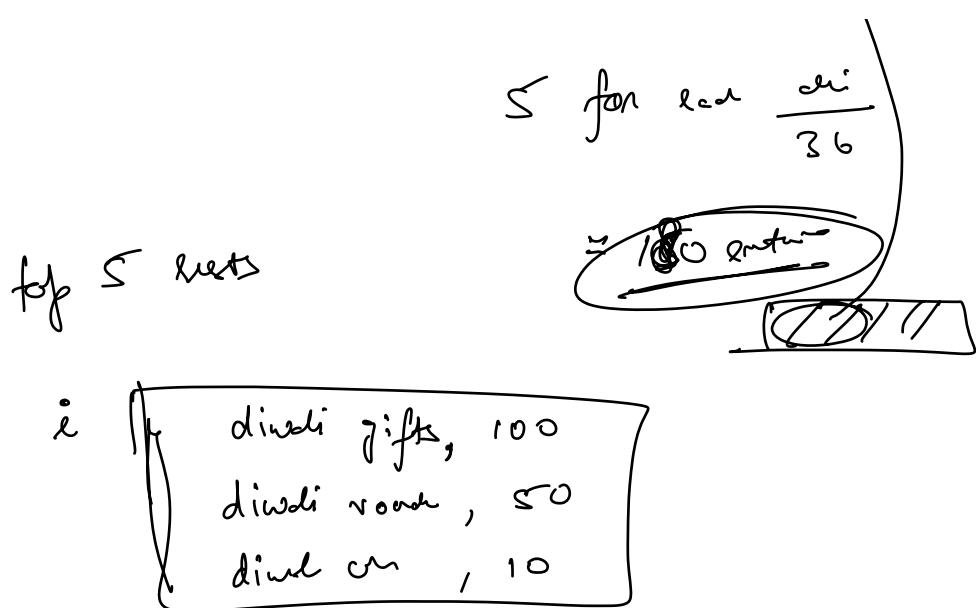
32 def

32 days

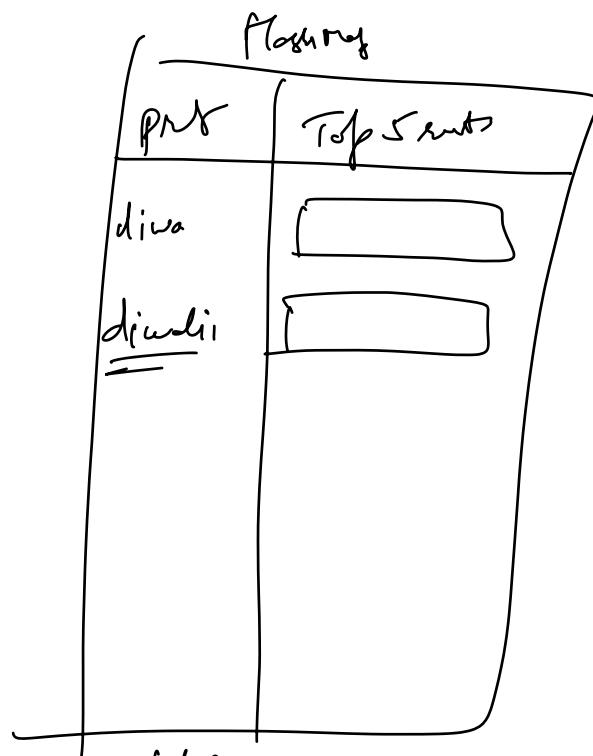
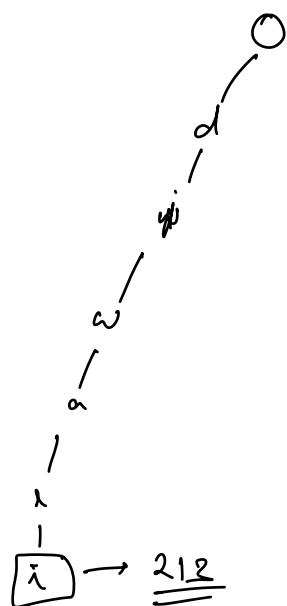
32 days

a ω i $\dot{\theta}$ diwali





Tree + Hashmap



fun

213	# hits
-----	--------

HashMap Approach

H1

Partial query	top 5 hits
d	[diu —, ——]
di	[— — —]
diu	[— —]
diva	[— —]
diwali	[———]
Ran	[———]

H2

Search query	# hits
diwali	100/101
diwali ch	50
Ranbir	20
:	

get Typer for Sys (qpx) → loop in H1

update freq (diwali) → H2 → update count
H1 → update top5

d
di
diw
diva
diwal
diwali

?

freq & m/s change

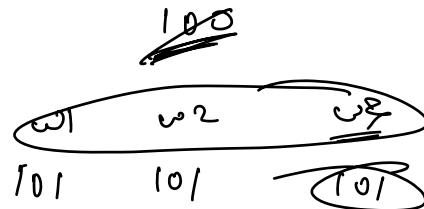
Trie → Prefix match
space optimization

Key - Value Store is automatically sharded based on
key

Optimizing Writes

10x more reads than writes
↳ updating multiple things
cont
each prefix → 30 prefixes

Write → Prevents reads → prevent stall reads
not important for this system
→ prevent writes → always needed



① Sampling

do we care about exact counts? \times

We will write only for every 100th search

1000

count = 0

updatefreq (start at) {

if (count == 100) {

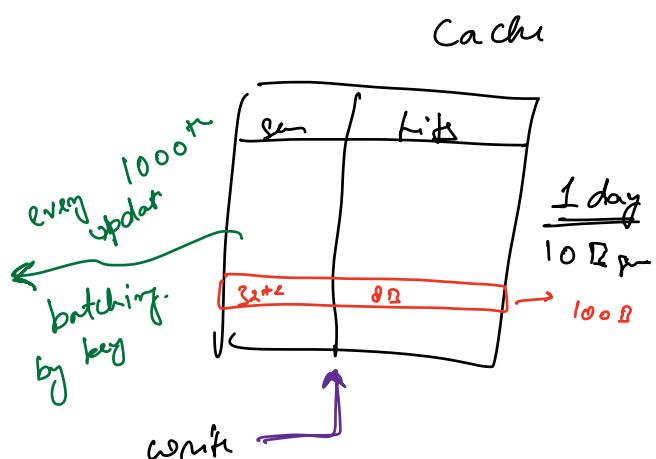
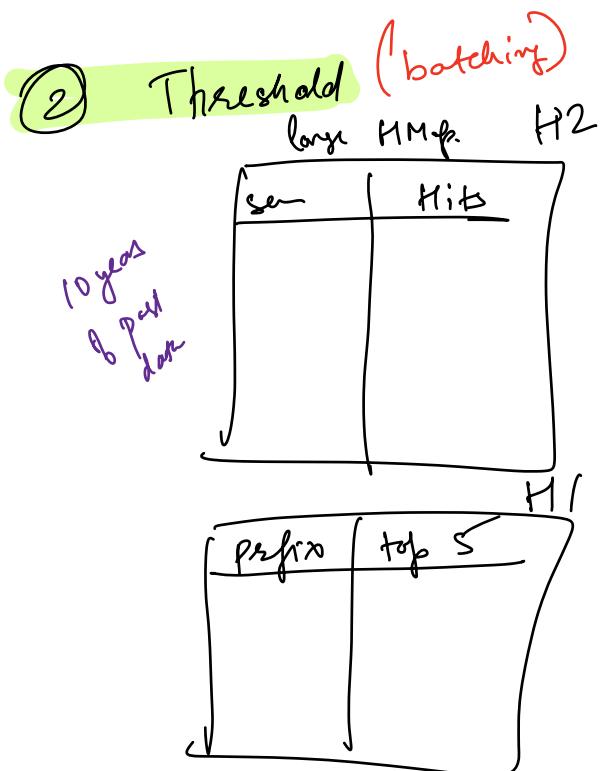
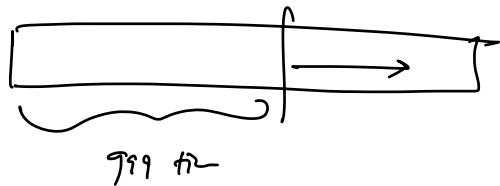
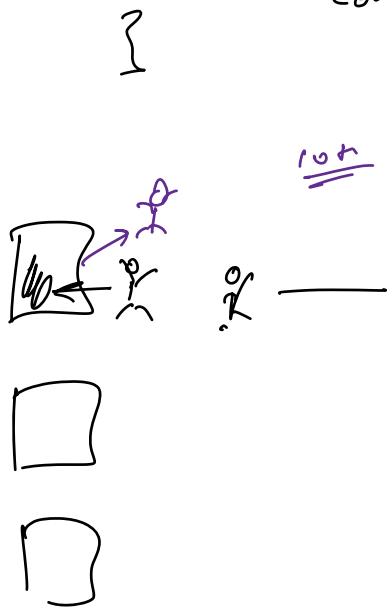
actually update
Count = 0

}

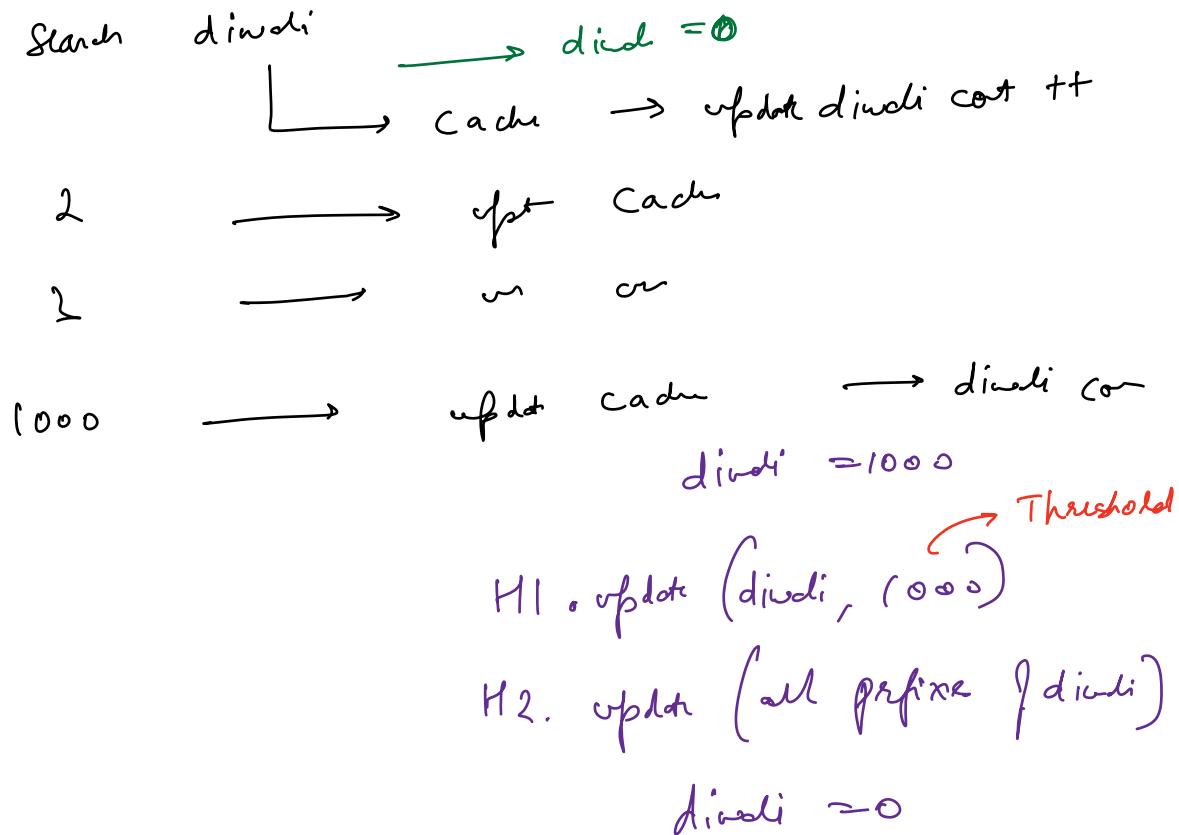
100 M

(M)

Count ++



$$100 \text{ B} \times 10^9 = 1 \text{ TB of data}$$



Recenty Factor
 Rodger Sontu → Top result
 Rodger Fed → wisdom → Rodger Fan.

Time Decay
 Hits
 } decay at \sim fixed rate
 decrease by 10%, 000 off each day

diwali → 10^n
 diwali vr → 1^n

cerebral fluid → 1000
~~cerebral fluid~~ → 50,000

after each day decrease each count by 1%
or 10%

and also by $\frac{100}{L}$ → reduce for date

diwali \rightarrow ~~10M~~ $8,999,900$

diwali celebration \rightarrow ~~1M~~ $900,000$

at day 1st
① decrease each value by 10% .

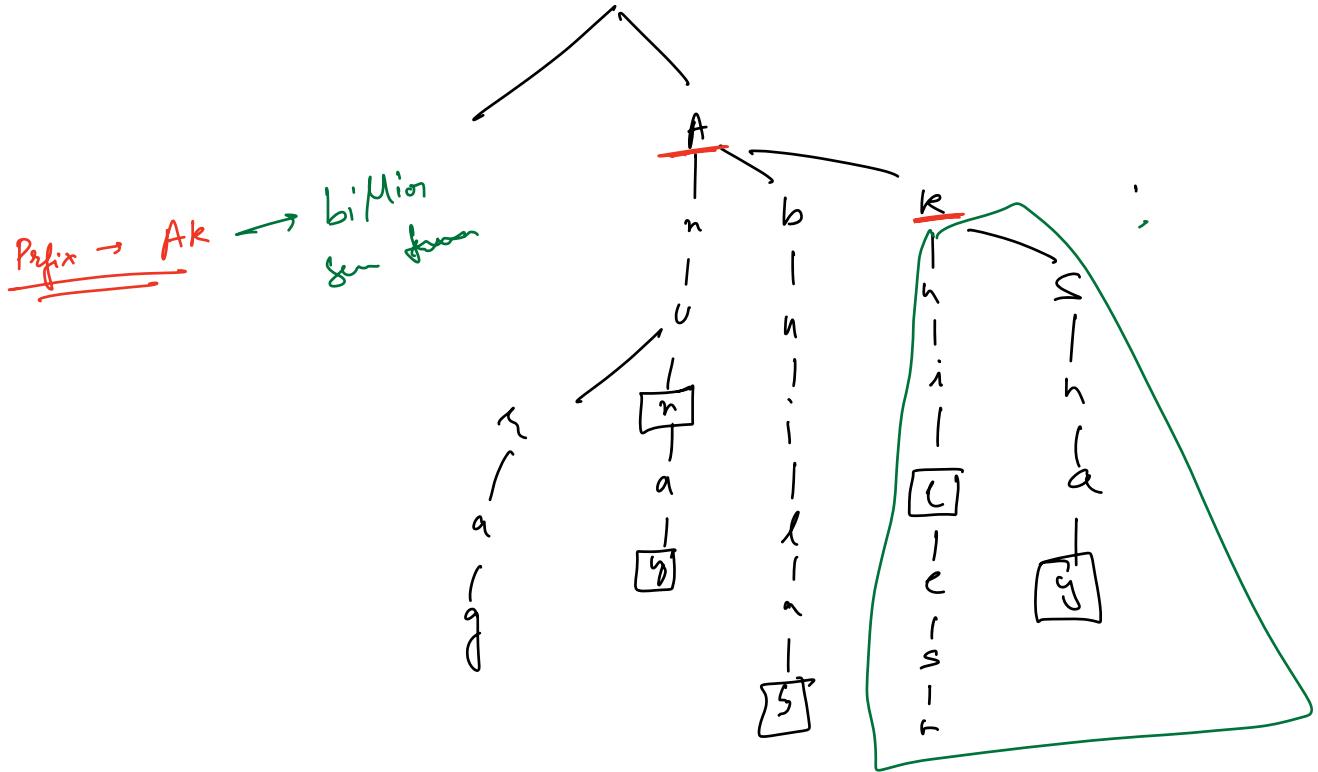
② dev each value by 100

~~44,900~~
~~45,000~~
cerebral fluid \rightarrow ~~50,000~~
cerebral palsy \rightarrow ~~10,000~~
~~X6~~ \rightarrow ~~50~~ $9,900$

1000 reads per sec
10 worth per sec

1M reads
10,000 worth/sec

initially \rightarrow procu \approx 1000 machines



Sampling

SAMPLE-SIZE = 1000

Count = 0

```
update freq (search-term) {
    if (Count == SAMPLE-SIZE)
        send update Reg (search-term)
    Count = 0
}
```

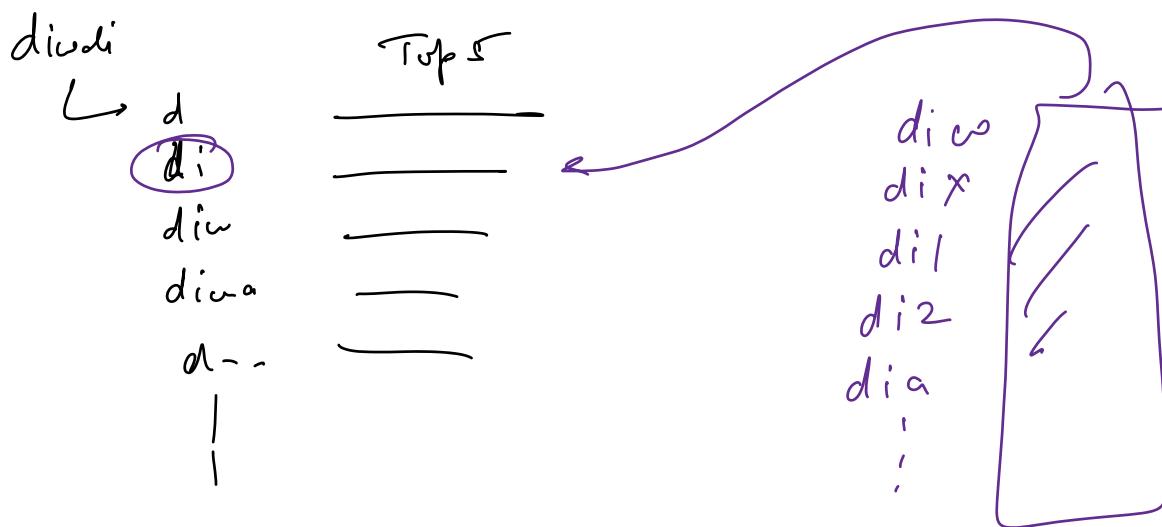
}

Count ++

Batching

DATA-SIZE

```
update f (sum +) {
    freq[s+t] ++
    if (freq[s+t] == 1000)
        send update Reg (sum +, 1000)
    freq[s+t] = 0
}
```

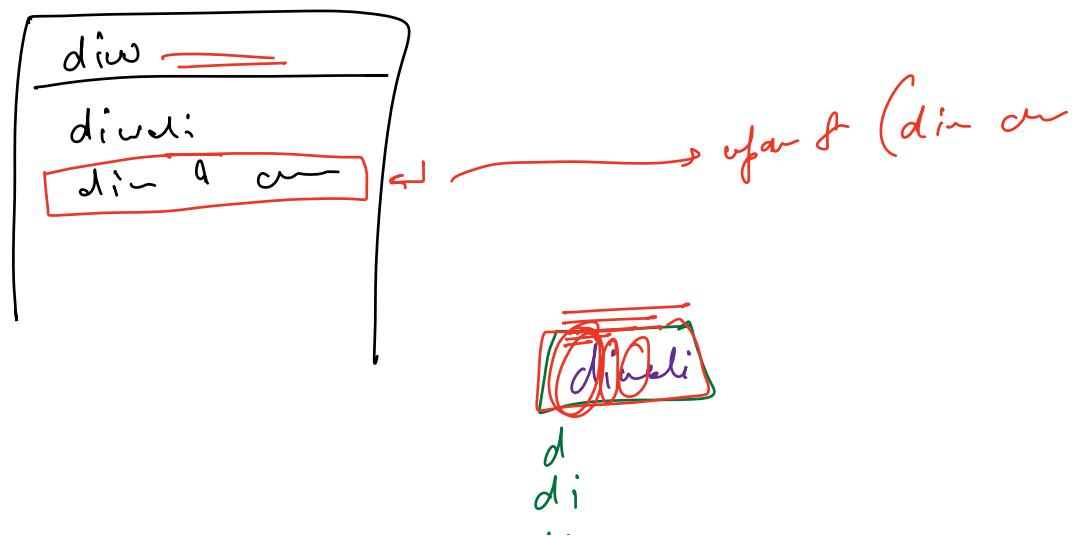


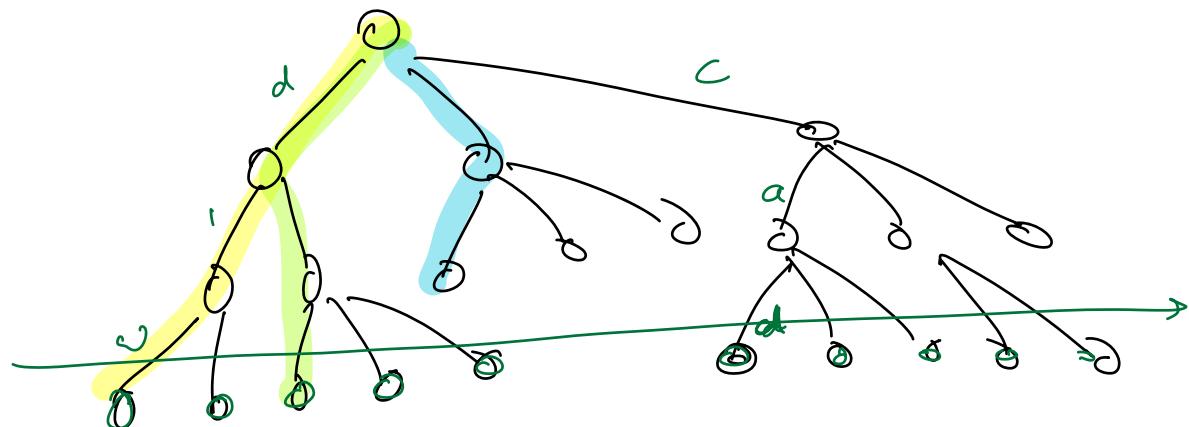
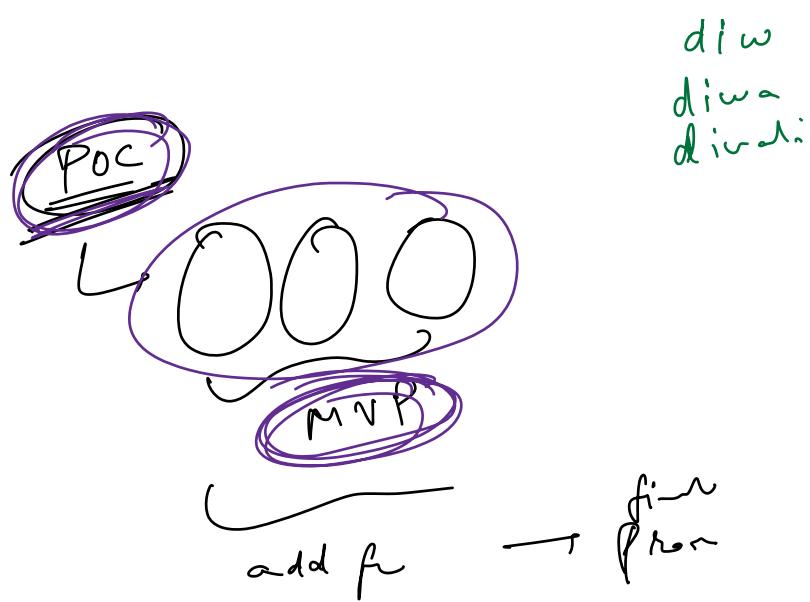
~~DB1~~

Prefix	Top 5
d	diwali
di	
<u>diw</u>	
R ₀	

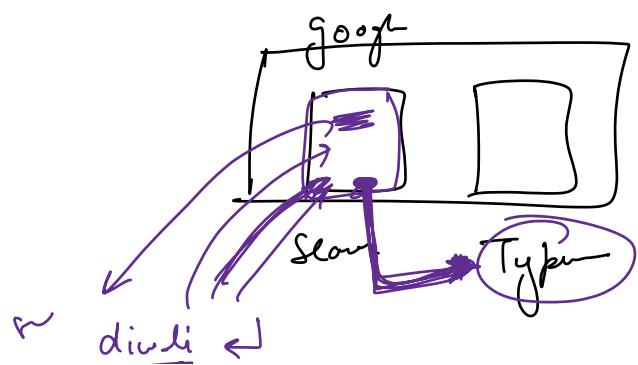
~~DB2~~

Searant	# hits
diwali	1M
dim ar	100,000
run	1M





Shand fry



Dian

Dianli → A
R ~ L_{O2} ~

~~11~~