# DESIGN GOOGLE SEARCH TYPEAHEAD

mi (mci)

michael jackson
michelle obama
michael jOrdan

Queries per second

1. Functional requirement
   MVP (Min. Viable Pro...
2. Estimation of scale
   → 1 Sharding
   → 2 Read vs write
   X [ 3 QPS → 1 machine &
3. Design goals/ non-fu
   rcqi.
4. API design
   → system design

---

## 1 MVP

1. Given **3+** characters, give me **top 5**
   suggestions of searc

   Top = most freq. search
         ↓
         most trending.

mic

[ ✓ — Personalisation → my searches, search from my region,
  ✓ — Spelling mistakes → X          search from
       → all suggestions will have            brows
         typed text as strict prefix

---

## Est. of scale

1. Shard → YES

   phrases → count

   → new search term

michael jackson  Enter
michael jackson

6-8 billion
searc
day

10 billion se
per day

30% searches are
new

3 billion new
phrases/day

3 B * 50 byte

* QPS * 20

"365" ⟶ ⟵                    * 365 ⟶

⟶ 1 PB

= 3 * 50 * 365 * 20

= 1000 TB ≃ 1 PB

② Read heavy vs write heavy

Both read        ⎡ ① Read → "mic" → get Top 5 s~
& write heavy.   ⎣ ② Write → any search → update
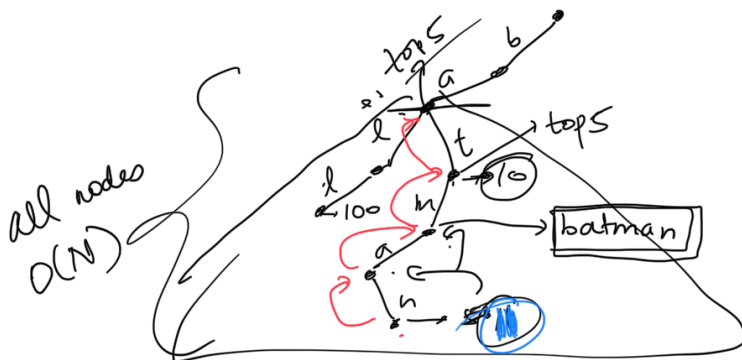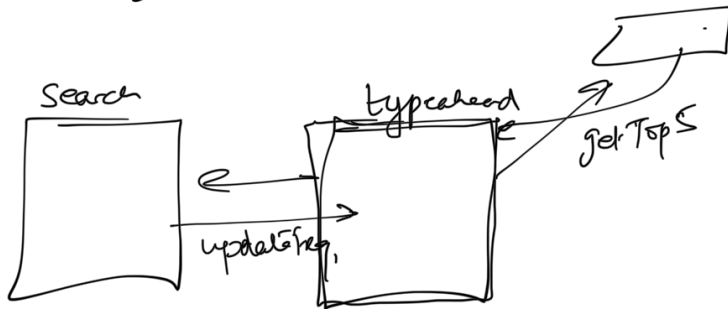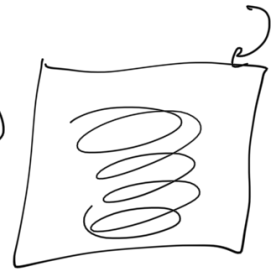                                            count/fre

Design goals / Non Functional

(HA) ✓

Latency of (reads) → ┃Super low┃

## API

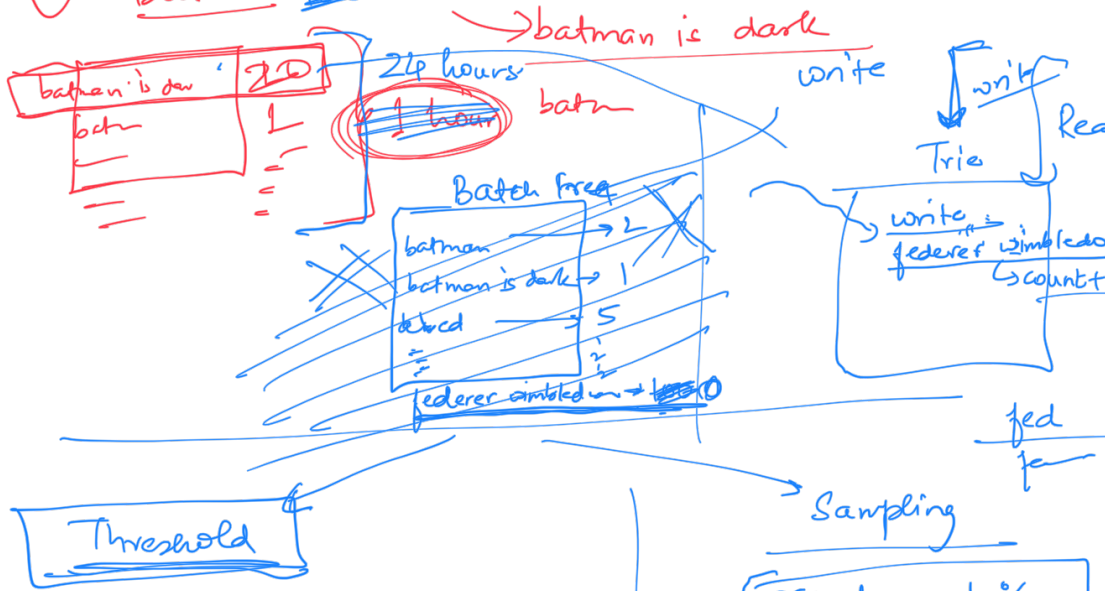⎡ ① get Top 5 suggestions (typed-str)
⎣ ② update Freq (search-term)



Search            typeahead                 get Top 5

          ⟵ update Freq

                                              ba

top 5                    b
         a
   s              t    top 5
  l                  /0
  l
all nodes   100   m
O(N)                        batman

          a        n

batman is dark

node {
  int freq
  list< pair<string, i

⎡ batm·:- 1    to 20
⎢ batm ··· 2    30
⎣ batm --- 3   10

|   |   |   |
|---|---|---|
| batm ...y | 25 |
| batm i—5 | 35 |
| batman | 11 |

① Batch 

[100]
[100]

→ batman is dark

| batman's dav | 20 |
| bch | 1 |
| | |

24 hours
1 hour   batm

write

write
Trie
Read

write
federer wimbledo
count

**Batch Freq**

| batman | → 2 |
| batman's dark | → 1 |
| detoced | → 5 |
| | |
| federer wimbledon | → 10 |

fed
fe

Sampling

random 1%
Searches

0 - 99

1% random write →  if (rand()%100 == 0) {
                        process
                        the wr
                    } else {
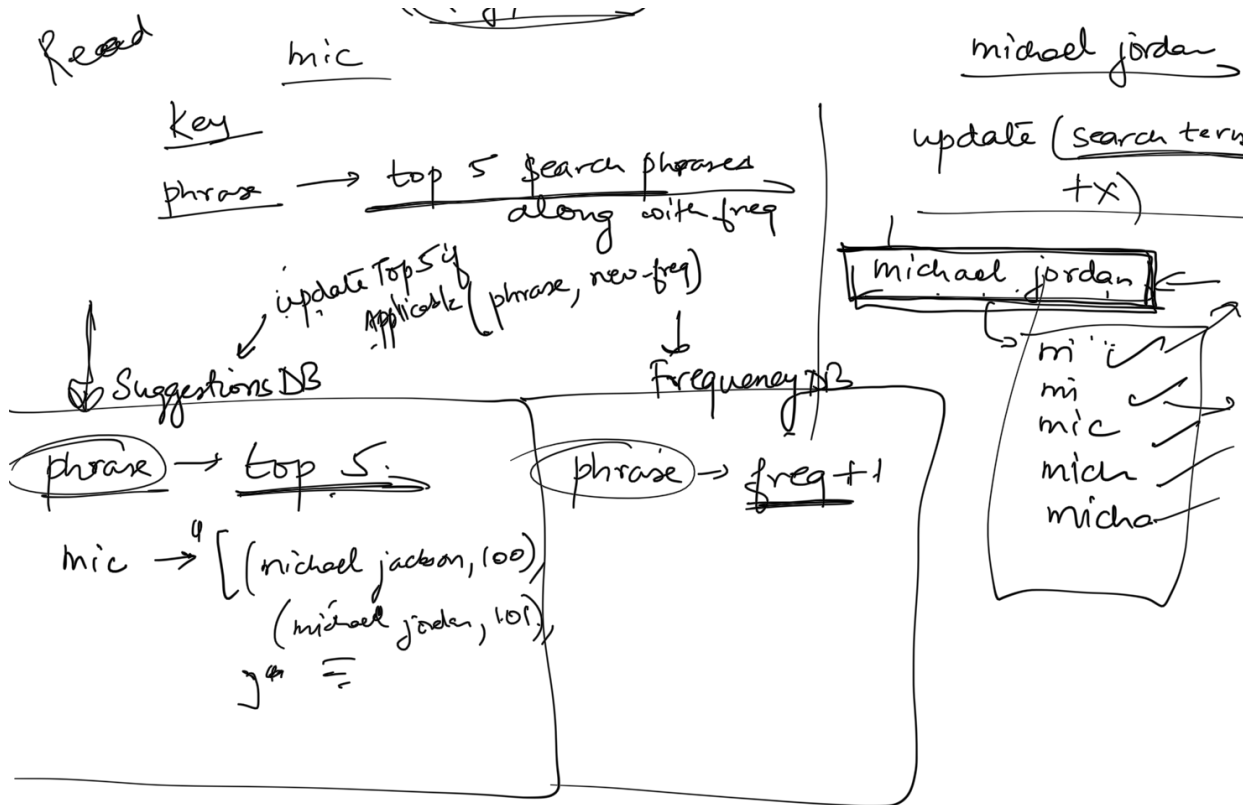                        // ignore
                    }

Threshold

② Read v's
write heav
× 100

updateFreq ( search_term ) {

    if ( rand()% 100 != 0) return ; // 99% q
                                    the C

    // logic of ~~write~~ updateFreq

SAMPLING

}

(Key, value)          .. mic

Read          mic                          michael jordan

Key                                         update (search term
phrase → top 5 search phrases                      +x)
         along with freq

     update Top5 if                        ┌─────────────────┐
     Applicable (phrase, new-freq)         │ michael jordan  │◄──
                                           └─────────────────┘
 ↓ Suggestions DB      Frequency DB              → mi  i
                                                   mi
(phrase) → top 5.   (phrase) → freq +1             mic
                                                   mich
 mic → [(michael jackson, 100),                    micha
        (michael jorden, 101),
       ] =


                              for(i=0 ; i< str·length; i++)
get Top5 Suggestions (phrase_prefix) {    { str[0:i]
                                          }
  return  Suggestions DB · find (phrase_prefix)
}

  update ( search_term, count) {

    cur_freq = Frequency DB· find (search_term)
    new_freq = cur_freq + count
    Frequency DB · update ( search_term, new_freq)
    for (i=0; i< search_term. length() ; i++) {
        prefix = search_term [0: i]
consistent
hashing
sharding key    top5 = Suggestions DB · find (prefix)
                  if ( search_term in top5) {
                                  // update search term
(Key) → value                        in top5
                3 else if ( new_freq > top5's least fr

Suggestions DB-update (prefix, top 5)

---

Michael jackson → 100,000,000

## Time Decay factor

TDF = 2          >1   <2

✓ Michall jackson                    1.7  1.0

| Day 1 | Day 2 | Day 3 | Day 4 | Day 5 |
|-------|-------|-------|-------|-------|
| (100) | 100   | 100   | 100   | 100   |
|       | 50    | 75  175 | 187  | 190   |
|       | 100+100 / 2 | 100/2 +100 | | |

1.9  2

✓ Shania Twain

| D1 | D2 | D3 | D4 | D5 | D6 |
|----|----|----|----|----|----|
| 100 | 10 | 1 | 1 | 1 | 100 |
| 50  60  30 | | 31  15  16 | 8  9 | | 109 |

$$\frac{1}{2^4} \quad \frac{1}{2^3} \quad \frac{1}{2^2}$$

Count = count / TDF

---

Hashmap            vs        Trie

Serialise

HOD

---

TDF = 2

Anshuman

| Day 1 | Day 2 | Day 3 |
|---|---|---|
| 240 | 240 | 300 |

120 + 60 = 180

current + prev / 2    11:59    $f = freq/2$

12:05

Anil

| Day 1 | Day 2 | Day 3 |
|---|---|---|
| 240 | 240 | 300 |

180

An
↓
Anil

---

K, V
- → supports sharding
- → persistence
- → HA
- → high throughput
- → high read heavy syste?

| HC | read | KV |
|---|---|---|
| HA | write heavy | doc cf |

(key), value

Hash
(key)
b

---

mic        26 words.

mic

m →

ã b c d e f ...

ni
Cf

ops

c
d
e
f
j

rowid $\longrightarrow$ {

CF
____
____

____
____
____
____
____
____
____

$\longrightarrow$

10

| row $\longrightarrow$ [ | --- · · ·