

## Agenda

### Approaching HRD Problems in Interviews

#### Google Typeahead



→ ① Step By Step

{ ② Google Typeahead

→ F<sup>m</sup> Reg

→ Non F<sup>m</sup> Reg

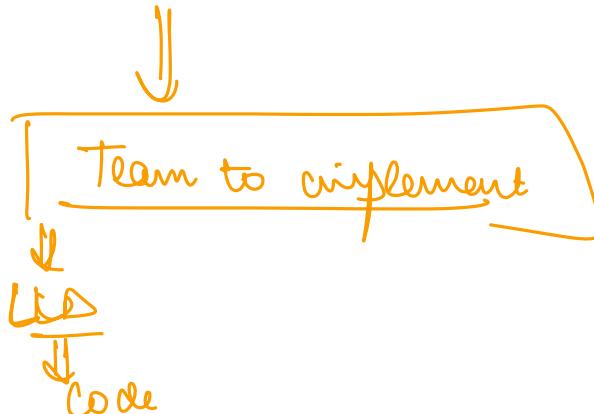
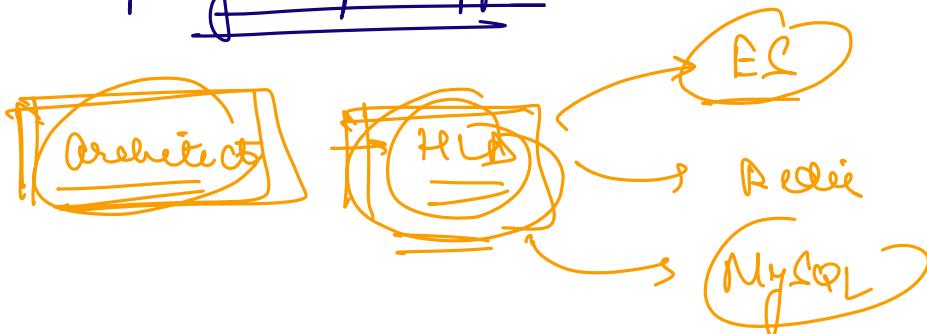
→ Estimation of Cost

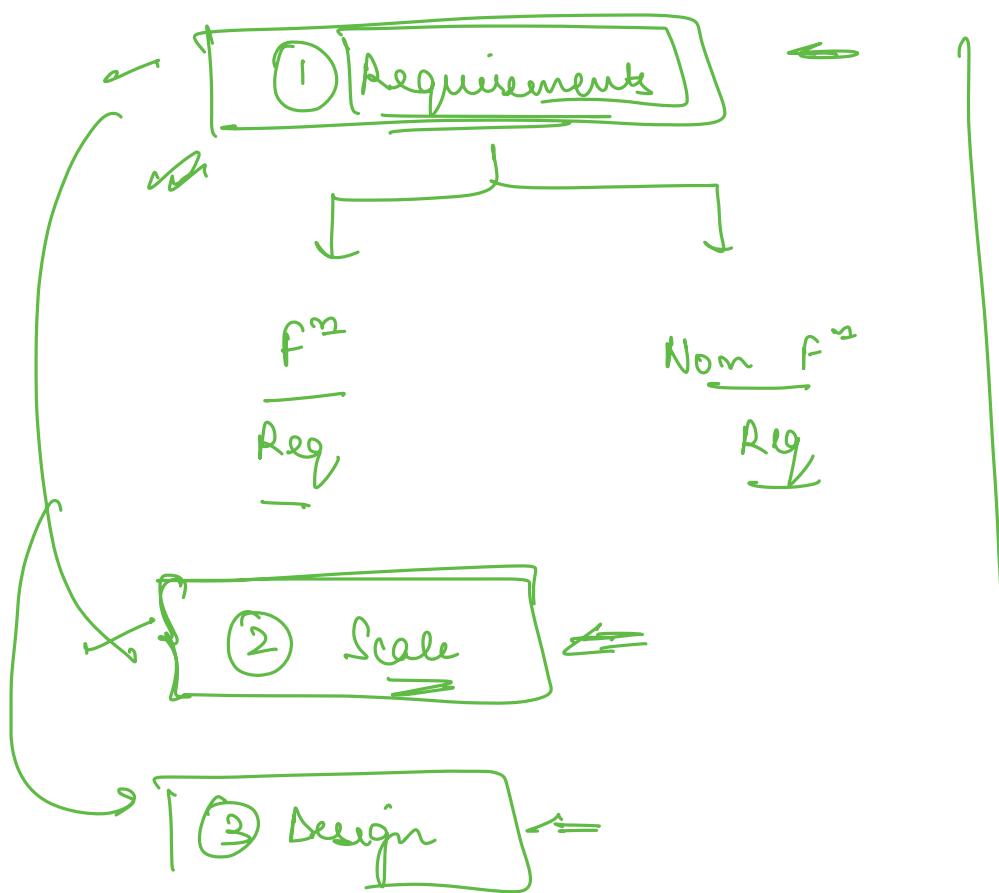
→ Design

③

How to solve recency  
in typeahead suggestion

#### ① Step By Step Approach





In an interview

Ambiguous Problem Statement  
(45 min)

- 1  $F^n$  Req
- 2 Nom  $F^n$  Req
- 3 Estimate Scale

Interviewer will help you with business metrics from which you will have to calc Scale

4 Design

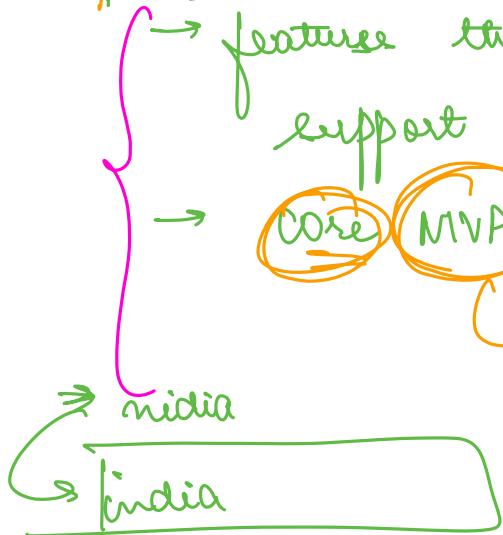
(Alongwith details about working of design)

5 Follow-Ups

→ No code in HLD regardless of type of company.

### Details of each step

#### ① F<sup>n</sup> Req



→ features that you should be able to support in your design

→ Core MVP features of the system

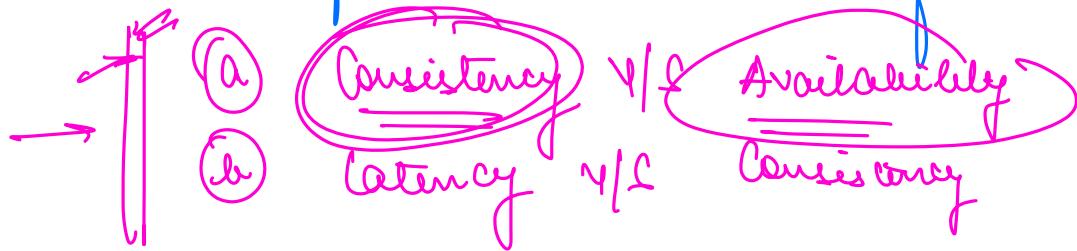
→ Minimum Viable Product

↳ smallest set of req to be able to do a POC for an idea

#### ② Non F<sup>n</sup> Req

→ how those features in F<sup>n</sup> req should work

→ expectation wrt those features



### ③ Estimation of Scale

→ how much load will my system be getting

- (1) Requests / sec      |      Queries / sec → Read
- (2) How much data my system might have to store

Why

- (1) Do I need to do sharding or not.
- (2) Read Heavy or Write Heavy or Both

↳ helps to choose correct DB

Multiple DB

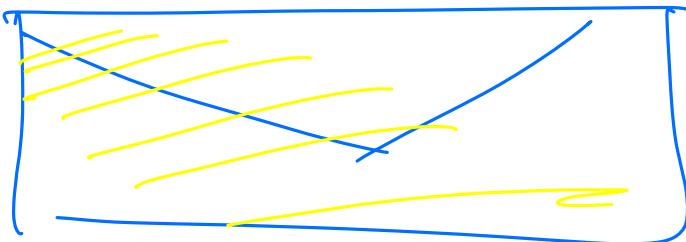
Or reduce

try to tame down  
one type of query

- Guesstimate ↳ Take an initial # and from that, basis multiple assumptions come to load req.
- Interviewer will not hold you accountable for wrong - assumptions

## Back of Envelope Calculation

- (1) no need of lengthy calculation
- (2) can approximate



$$\begin{aligned}
 & \left\{ \begin{array}{l} 50 \times 10^6 \\ 60 \times 24 \times 60 \\ 50 \times 10^6 \\ 50 \times 25 \times 50 \end{array} \right. \\
 & \xrightarrow{\approx} \frac{50 \times 10^6}{50 \times 25 \times 50} \approx \frac{10^6}{10^3} = 1000
 \end{aligned}$$

$\rightarrow$  approx on wrong side

- (2.5) list down APIs you will be building

## 3) Design

- $\Rightarrow$  attend the classes
- $\Rightarrow$  give MOCK Interviews.

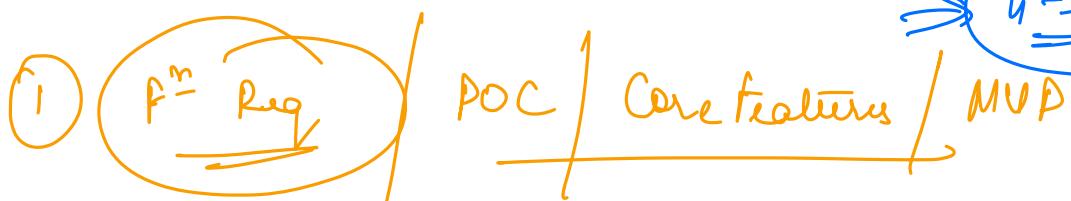
getFeed()  
createPost()

## 4) Followup

## DESIGN

GOOGLE TYPEAHEAD

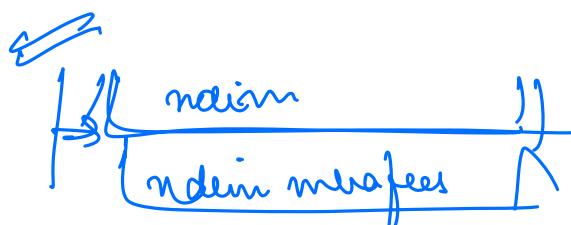
Q) → what is typeahead (RAAE)



- ⇒ a) can be called after any # of chars.
- ⇒ b) 10 Suggestions / request =
- ✓ c) ~~spellcheck~~
- ✓ d) ~~personalization~~
- ✓ e) give suggestions based on popularity.  
→ to MOST frequent matching queries.



suggestions based on matching prefix



$f^n$   
~~Reg~~

- ① queried after every key stroke
- ② 10 suggestions / query
- ③ most frequent queries having current input as the prefix
- ④ Spell check
- ⑤ ~~personalization~~
- ⑥ Search by simple ASCII char.
- ⑦ ~~popularity~~

↳ VO: most frequent asked query across

complete history

↳ VI: trends; give more weightage to recently popular query

Sa

Caria nizya

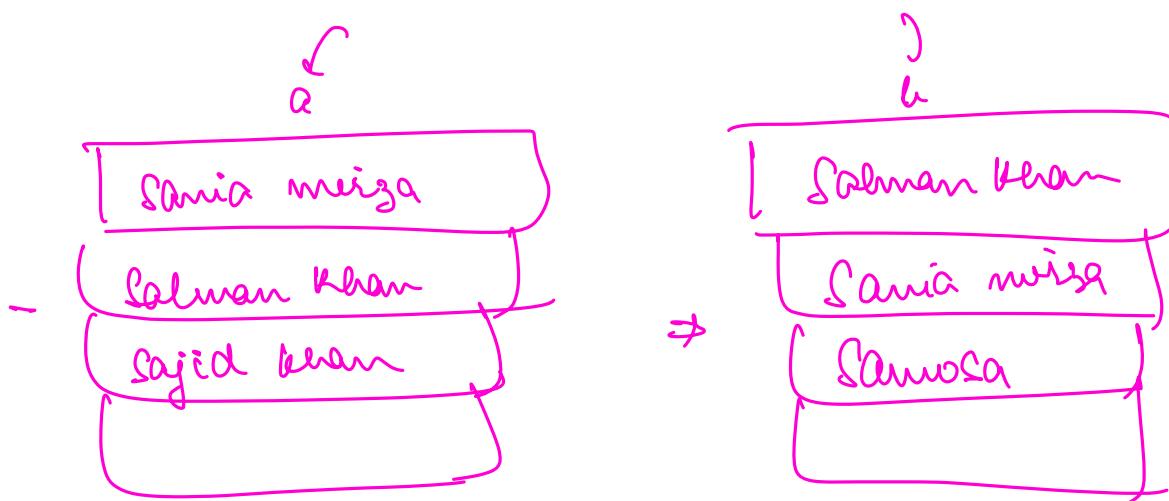
Salman Khan

take it up at end if we have time

⇒ follow up

② Non F<sup>m</sup> Reg

① Avail V/S Consistency  $\Rightarrow$  Available



ii) Latency V/S Consistency

$\Rightarrow$  V.V.V.V low latency

$\Rightarrow$  Competing against typing speed  
of people

### ③ Estimation of Scale

- ⇒ ① Storage
- ② QPS

nd  
nda  
ndan  
—

Storage  
—

$26 \times 26$   
26 15

→ Store freq. of every query in the post

(S)

Shahrukh  
Salman  
Shahrukh Khan

India  $\Rightarrow 10M$

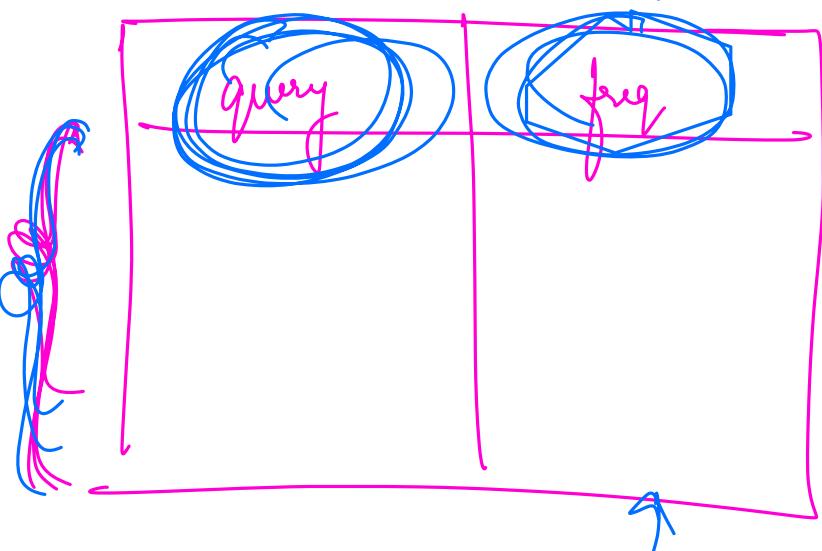
Salman  $\Rightarrow 5M$

Shahrukh  $\Rightarrow ?M$

Shahrukh Khan  $\Rightarrow 2M$

160TB

long  $\Rightarrow 8B$



32TB

Users with Google / day  $\Rightarrow$  500 M users

# Search/user/day  $\Rightarrow$  20

$\Rightarrow$  total searches / day  $\Rightarrow$  10 B searches / day

% of unique queries  $\Rightarrow$  20%

$\Rightarrow$  # of unique queries / day  $\Rightarrow$  2 B =  
length of each query  $\Rightarrow$  40 char

$\Rightarrow$  Size of queries / day  $\Rightarrow$  2 B  $\times$  40 B

= 80 GB

Let's provision our system for 5 years?

$$\Rightarrow 80 \text{ GB} \wedge 5 \wedge 365$$

$$\Rightarrow \approx 80 \times \frac{5 \times 365}{400} \approx \\ \approx 80 \times 2000$$

$$\approx 160 \times 10^3 \Rightarrow \text{GB}$$

160 TB

+

19.2 TB

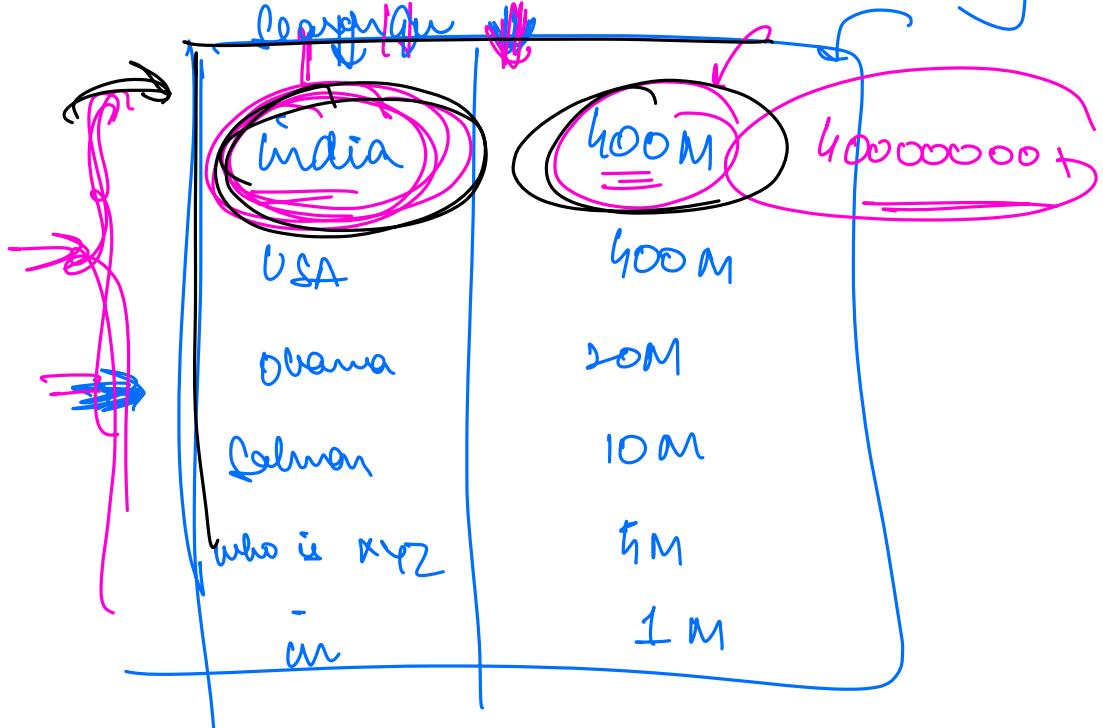
in 5 years

$$2 B \times 365 \times 5 \times 8 B$$

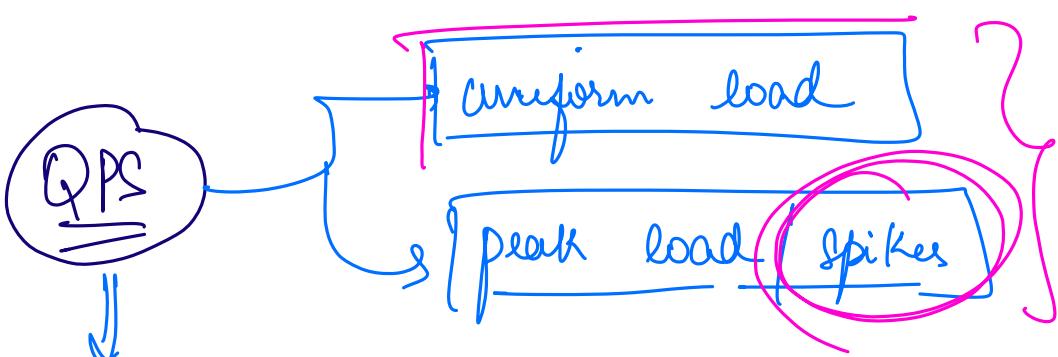
$$\Rightarrow 2 B \wedge 8 B \wedge 400 \wedge 5$$

$$\Rightarrow 16 \text{ GB} \wedge 2000$$

$$\rightarrow 32000 \text{ GB} \Rightarrow 32 \text{ TB}$$



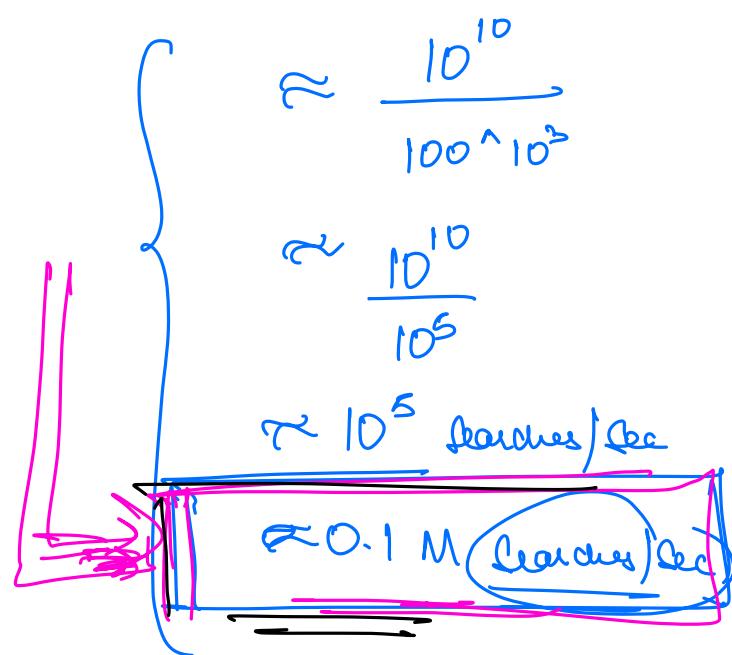
= {  
 → 192 TB of data just to store frequency  
 of every search over 5 years in google }  
 {  
 → can't fit in 1 machine  
 → Sharding } } } } }



→ maximum (peak) qps

→ 10 Billion Searches / day

$$\Rightarrow \frac{10 \times 10^9}{24 \times 60 \times 60} \approx \frac{10^9}{25 \times 4000}$$



Assume a user writes 10 chars before searching

on average

$$\Rightarrow 0.1 \text{ M } * 10 \text{ / sec}$$

$$\Rightarrow 1 \text{ M suggestion requests / sec}$$

$$\Rightarrow \# \text{ Reads} \Rightarrow 1 \text{ M reads / sec}$$

$$\Rightarrow \# \text{ Writes} \Rightarrow 0.1 \text{ M writes / sec}$$

Read Heavy

Write Heavy

Both

$\times < 1\%$  of prev one

$\Rightarrow$  Other thing doesn't even matter

$\Rightarrow$  Multiple DB for diff loads

$\Rightarrow$  or somehow reduce one

$\Rightarrow$  other queries are statistically insignificant

→ type of query

Summary

(i) F<sup>m</sup> ✓

(ii) Non F<sup>m</sup> ✓

(i) Availability

(ii) Low latency

(iii) Scale

| q     | Count |
|-------|-------|
| india | 50    |
| dot   | 20    |
| long  | 40    |

(i) 192 TB storage over 5 years  
↳ frequency of every search

(ii) # Reads → 1M / Sec

(iii) # Writes → 0.1M / Sec

Whenever a  
Search hits

Whenever a  
Suggestion  
Req

(iv) APIs

inde

→ i

→ in

→ vid

③

## ④ APIs

⇒ `List<String> get Suggestions (query)` // read  
    |M/Sec

⇒ `void update freq (Search Terms)` // write  
    → 0.1 M/Sec

## ⑤ Design

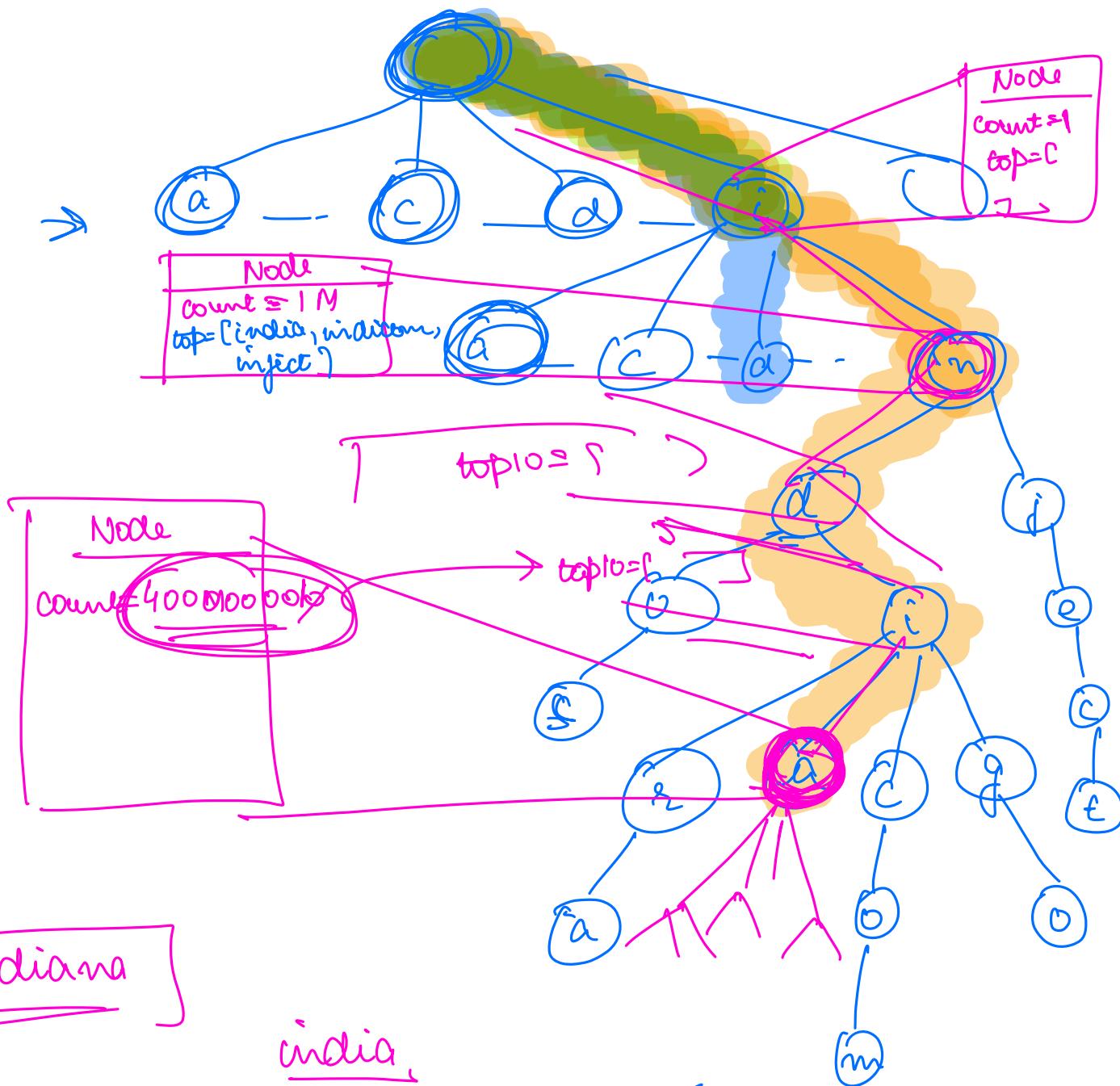
forget Storage was high

`List<String> get Suggestion (String prefix)` {  
    // find Strings that have prefix.  
    // as their prefix.

}

HW:

learn about   ries



Indiana

Node f  
india  
 i → ✓  
 in → ✓  
 ind → ✓  
 indi → ✓  
 india → ✓  
indian → ✗

Node f

int countOfOccurrences Till That String;

List < String > top 10 queries With This As Prefix;

$O(|\text{Size of Prefix}|)$

① get suggestion (prefix) {

// go to node for prefix in tree  
// return node. top 10's

⇒ v. fast  $\Rightarrow$  v. v. low latency

Compared to going through the complete table.

② update freq (Search Term) {

// go to the node associated to SearchTerm  
// update count

// update top 10 of all nodes in path  
all that need.

Samia

"S" →  
"Sa" →

"Sam" →

"Samia" →

"Samia" → ✓

"Samiam" = X  
→

[Sea,

Smoke, Salmon, Salmon, ~~Salmon~~)

Samia

H|W

① Read Trees

② Read typed notes of this class