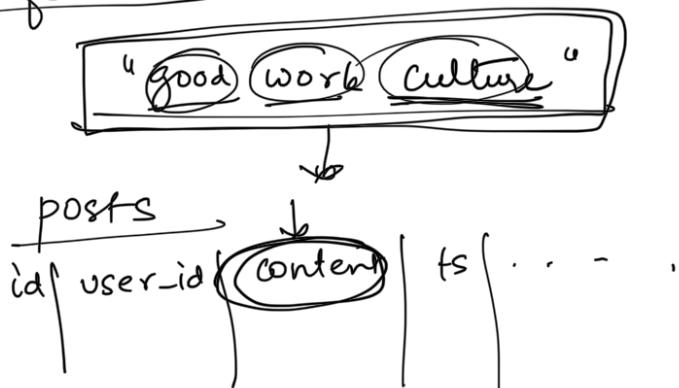


# FULL TEXT SEARCH

- ## Search feature

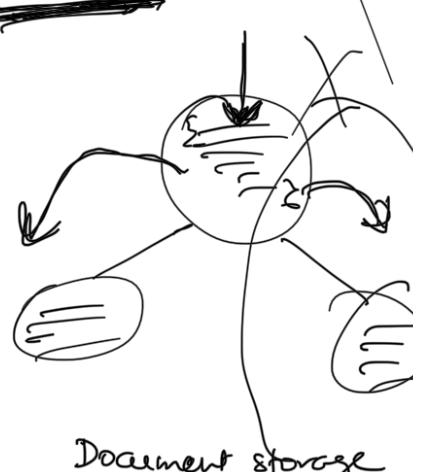


```
SELECT * FROM posts  
WHERE content  
LIMIT 10
```

Slow

~~$O(N)^*$~~   
 $O(\text{substring check})$

LIKE  $^u \%$  / good work  
~~Culture~~)  $^u :$



# NoSQL

key, value

CF

## Document storage

{ postid  $\Rightarrow$  "content" }

· 100million

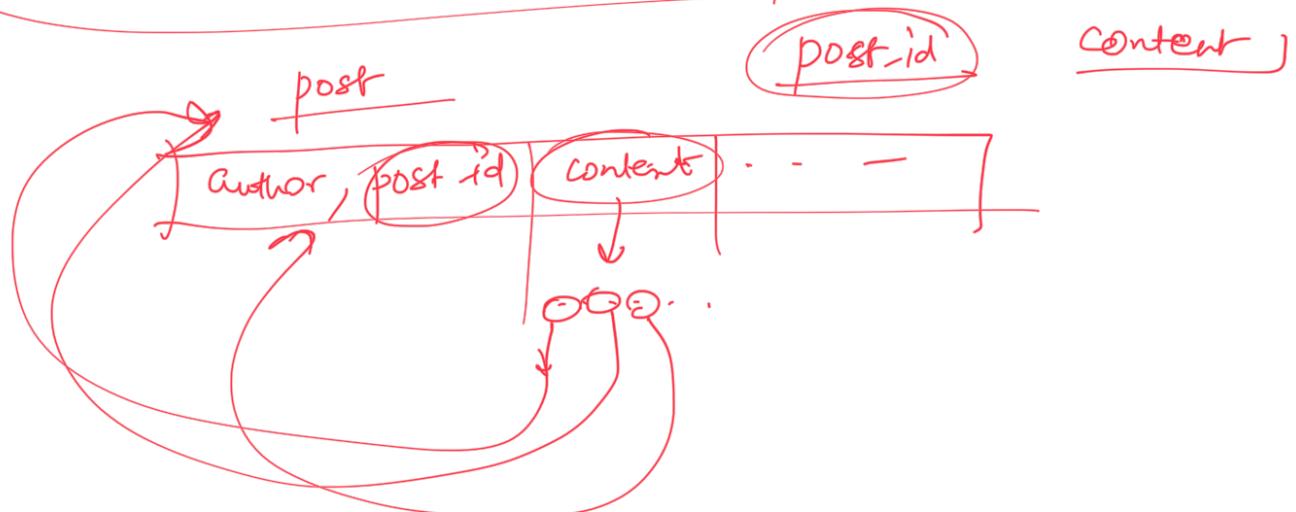
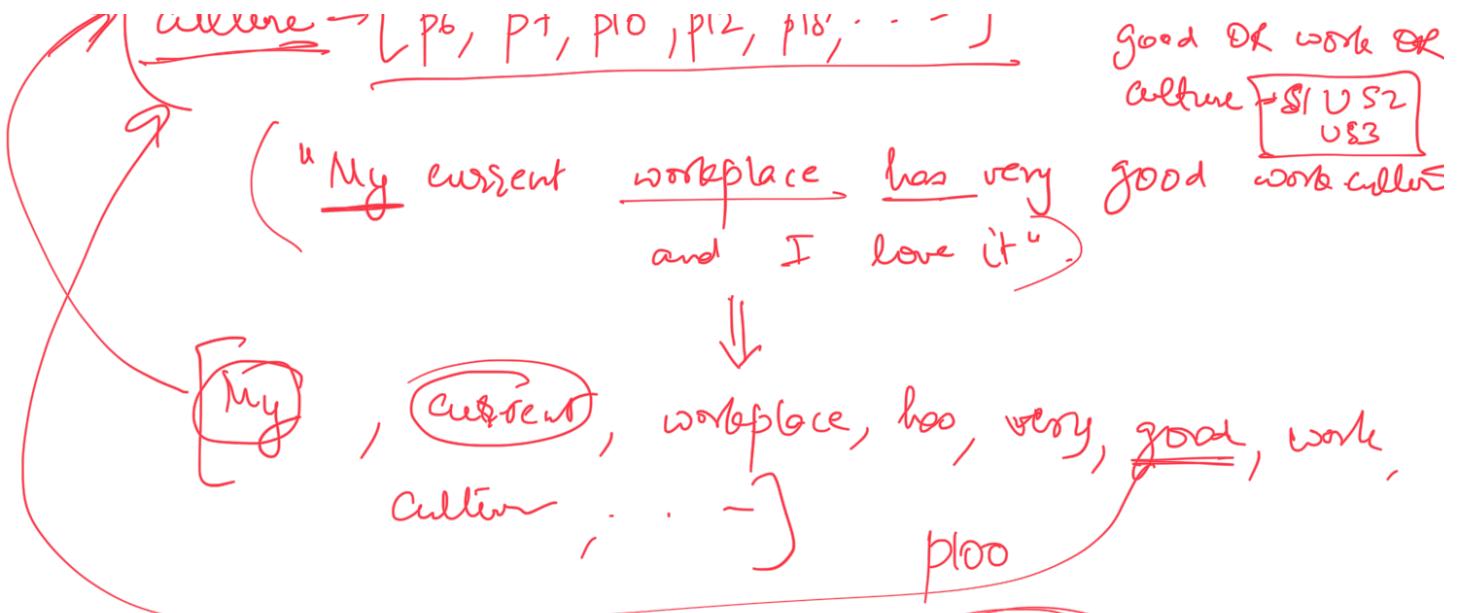
good work allure

## Inverse index :-

The diagram illustrates two sets of points labeled  $p_1$  through  $p_{100}$ . The first set, labeled "good", contains points  $p_1, p_5, p_{10}, p_{15}, \dots, p_{100}$ . The second set, labeled "work", contains points  $p_5, p_6, p_8, p_{10}, p_{12}, p_{16}, \dots$ . The sets are enclosed in brackets above the labels.

S1 n S2 n S3

word → poet,

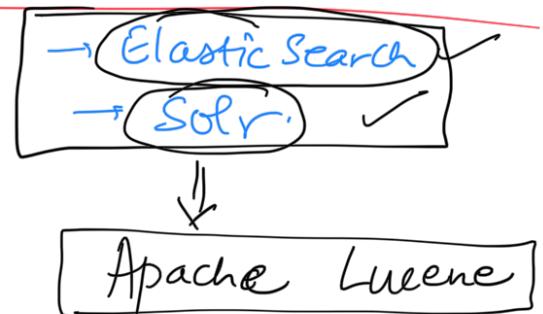


"fun" marathon<sup>4</sup>

I ran a marathon → matches

best ← good

## Fuel text search



### ① Word elimination

→ ["a", "an", "the", "of", "on", "by", ...]

special &

(2) Tokenisation → Break document/post into words

[ "ran", "marathon", "lost", "weekend" ]

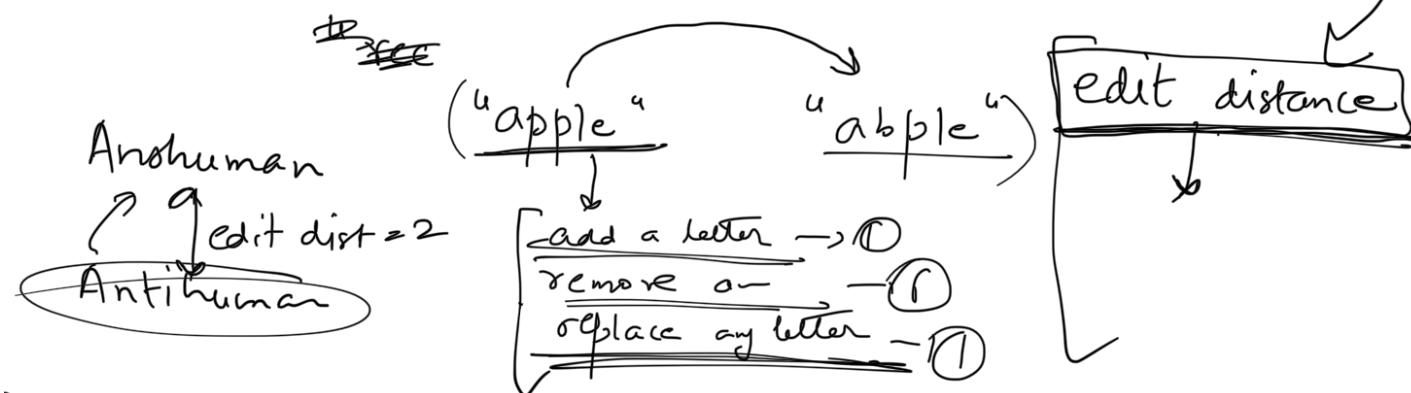
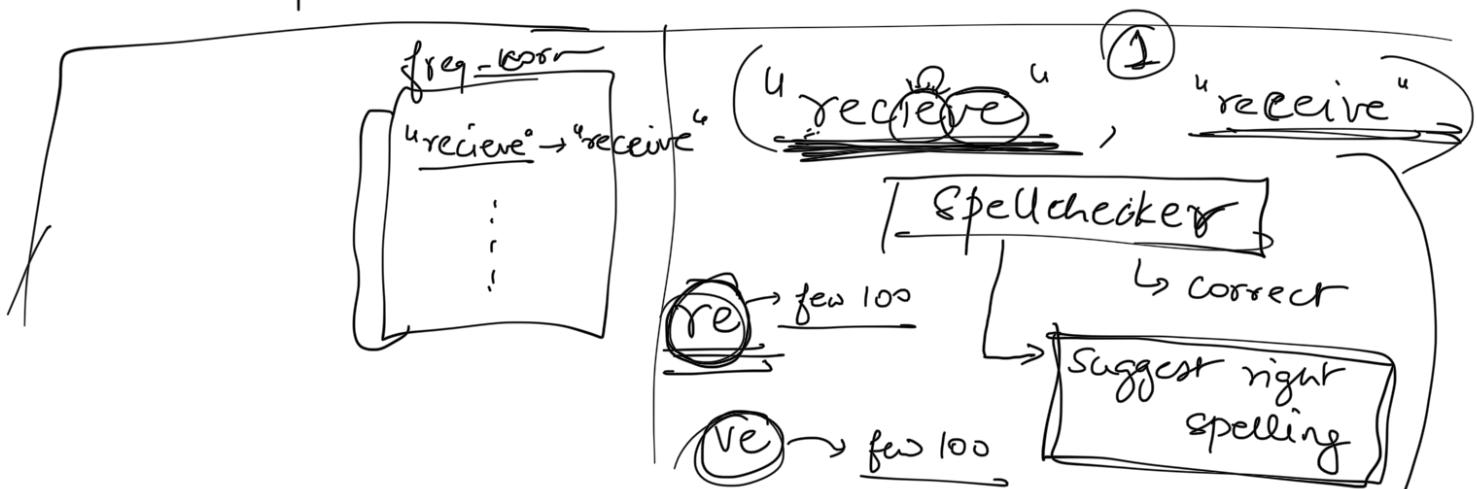
(3) Stemming / Lemmatisation.

ran → run

ran  
running → run  
run

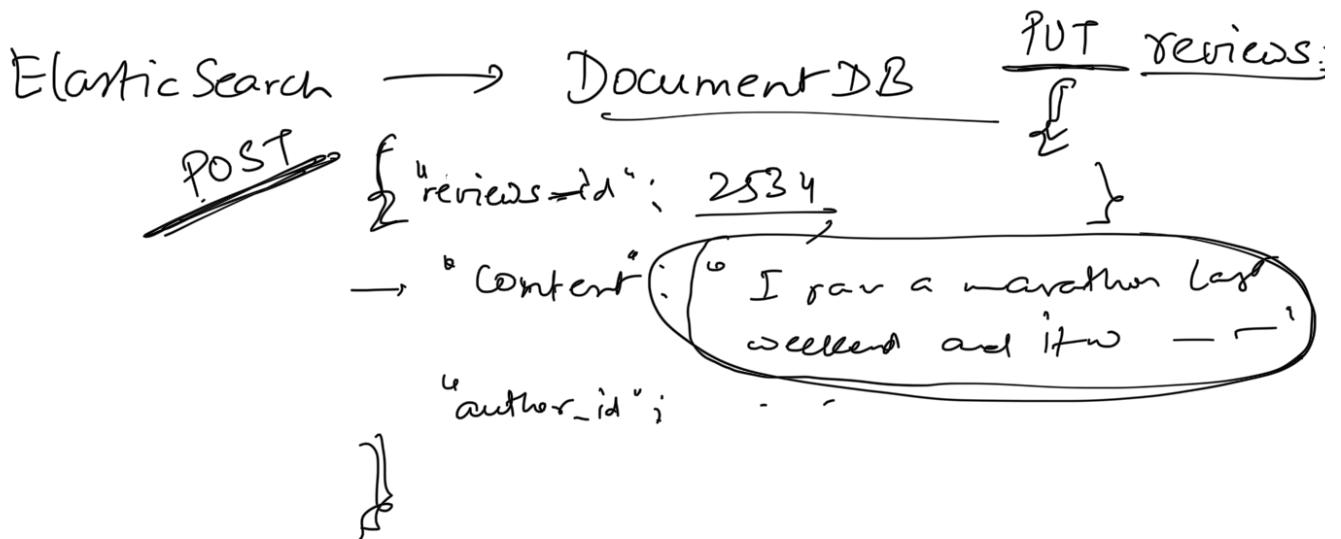
[ "run", "marathon", "lost", "weekend" ] → good  
best  
better  
good

(4) Stop words (abusive words, controversial)



Inverse Index

"word" → [ (doc\_id: [pos1, pos2 ...], doc\_id2: [...]) ]



GET

}

Syntax

"Content": "run marathon"

AND, OR  
exact

- UseCases
- ElasticSearch → {≡}
- Syntax → YT videos

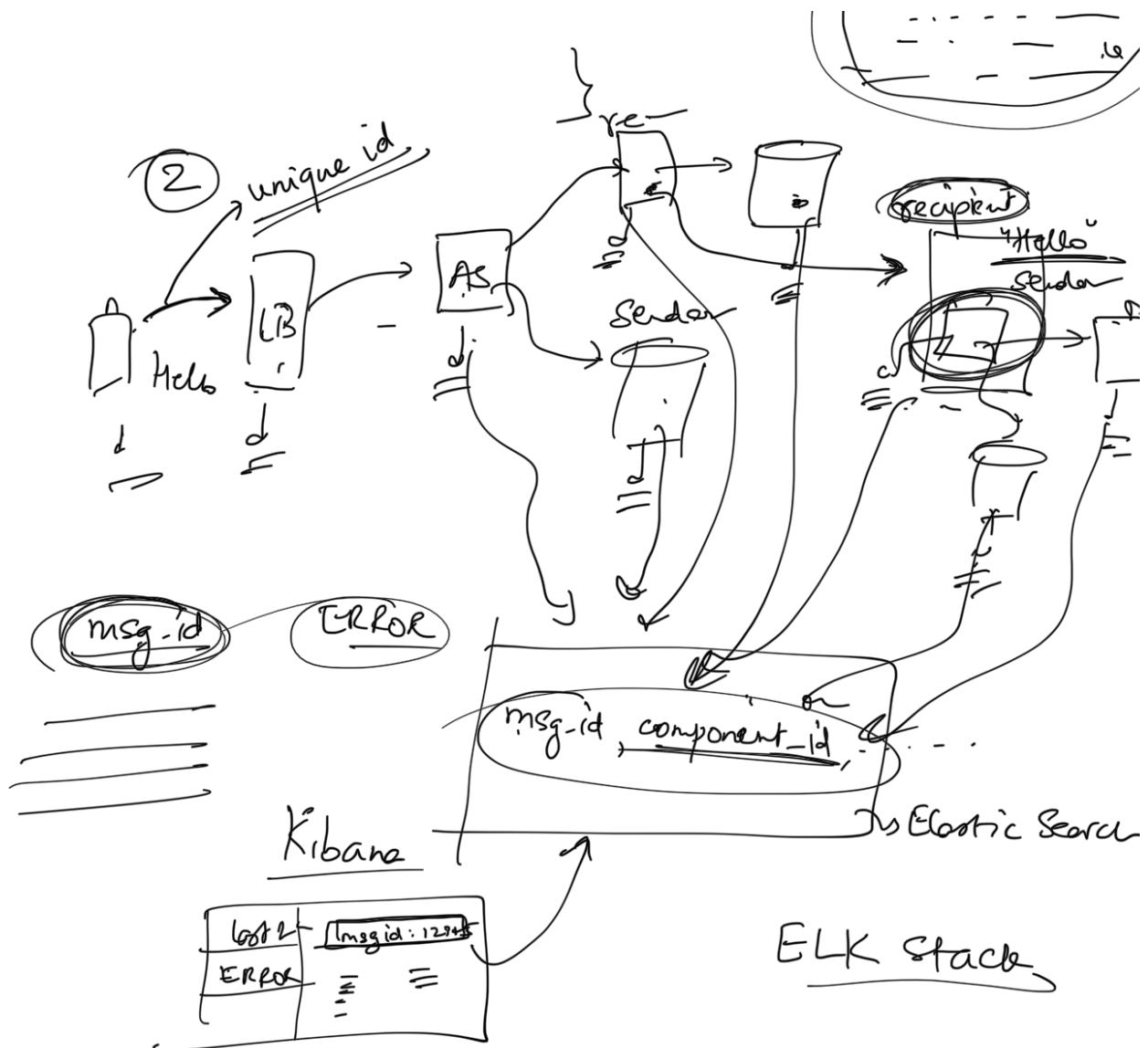
### UseCases

#### ① ATS / Scaler

"java spring boot"

"keywords" → matching resume\_ids.





③ Search on your own website

{ page: " - - - ",

A hand-drawn diagram illustrating a concept. At the top left, the word "Course" is crossed out with a horizontal line. To its right, the word "pause" is written vertically, with a small arrow pointing from the crossed-out "Course" towards it. A large bracket is drawn underneath "pause", spanning most of its height. To the right of this bracket, the word "Content" is written vertically, followed by a series of short horizontal dashes. Another large bracket is drawn underneath "Content", also spanning most of its height. Below the first bracket, there is a curved brace that wraps around the bottom of the "Content" bracket and extends downwards, ending in a small arrow pointing towards the bottom right.

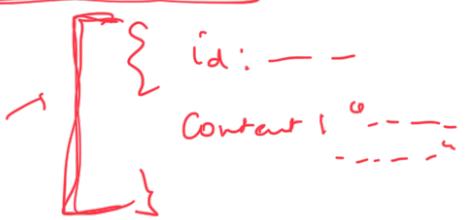
## Non-functional requirements

A hand-drawn diagram illustrating a spot welding process. It shows a vertical metal plate being held by two hands. A horizontal beam, representing an electrode, is positioned above the plate. A red arrow points from the left towards the electrode, indicating the direction of the welding current. Another red arrow points from the electrode towards the text 'Spot' on the right, indicating the point of contact.

- L
- ✓ ① Fault Tolerant
  - ✓ ② HA
  - ✓ ③ Read heavy sys'

### Terminologies

#### - Document

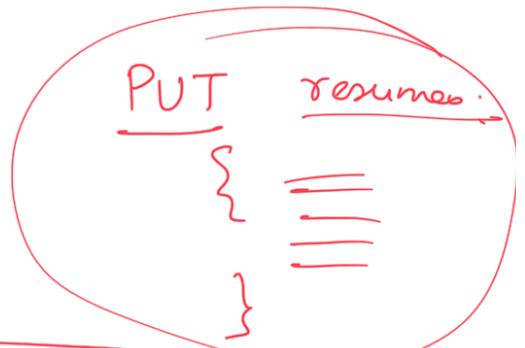


#### - Index

#### - Node

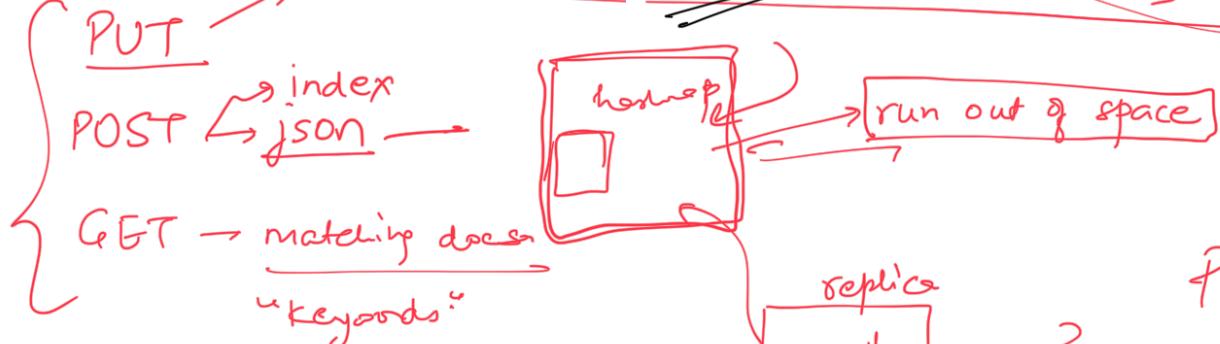
Resumes → LinkedIn Posts.

DB ↔ Index



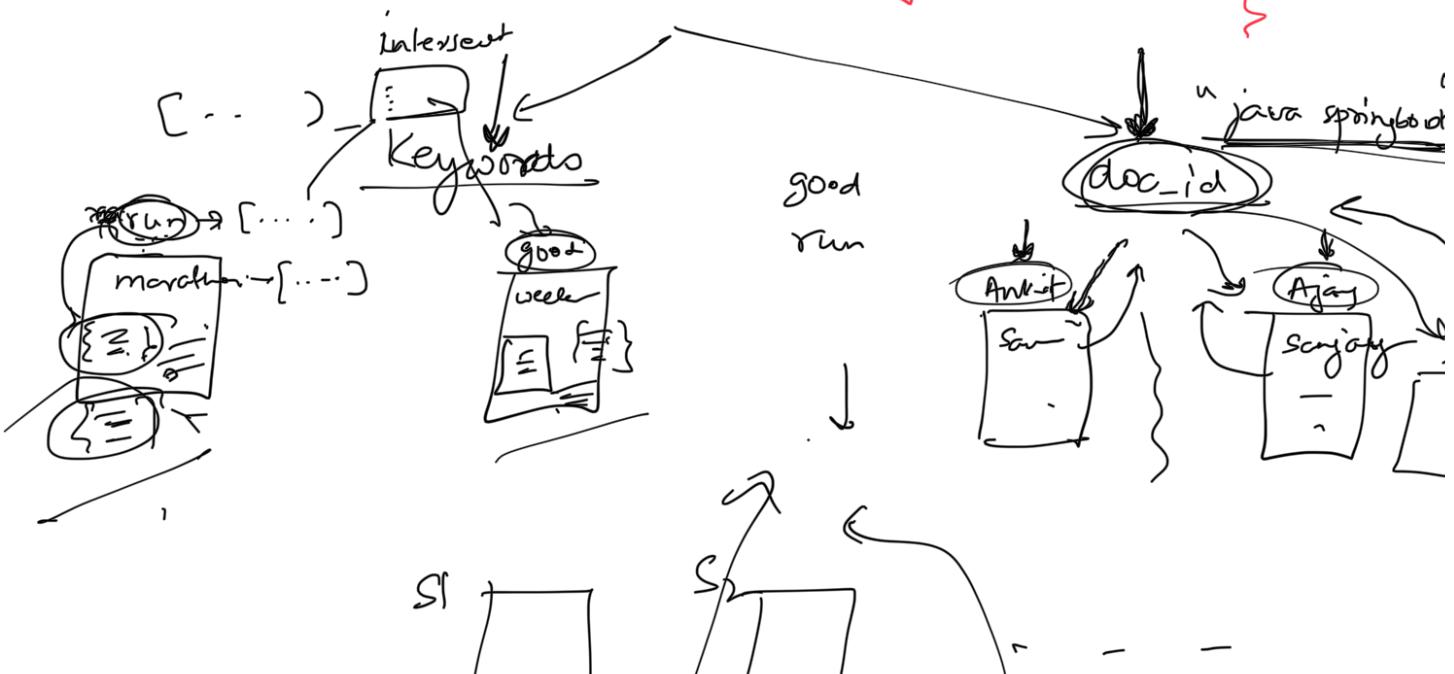
create new DB/index

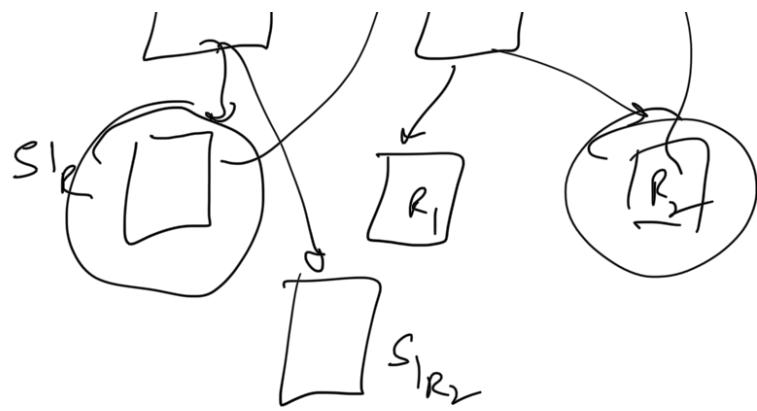
index



POST, res

{  
---  
---  
---}



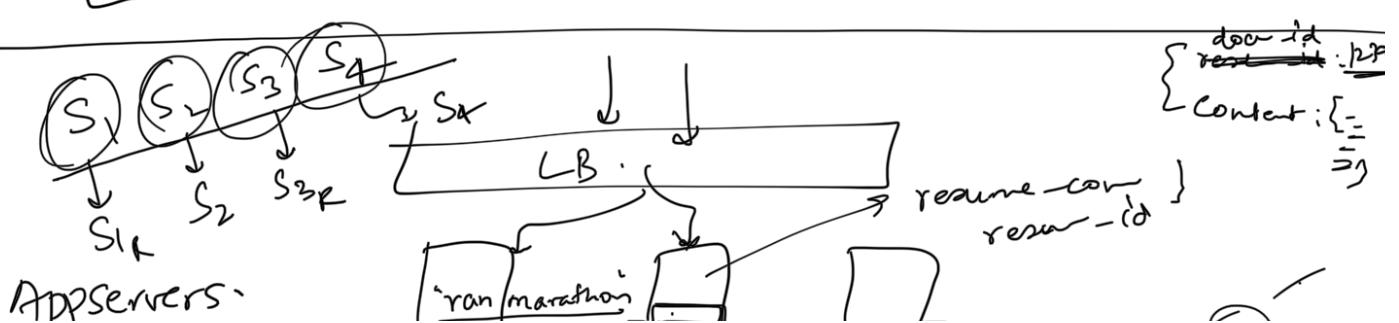
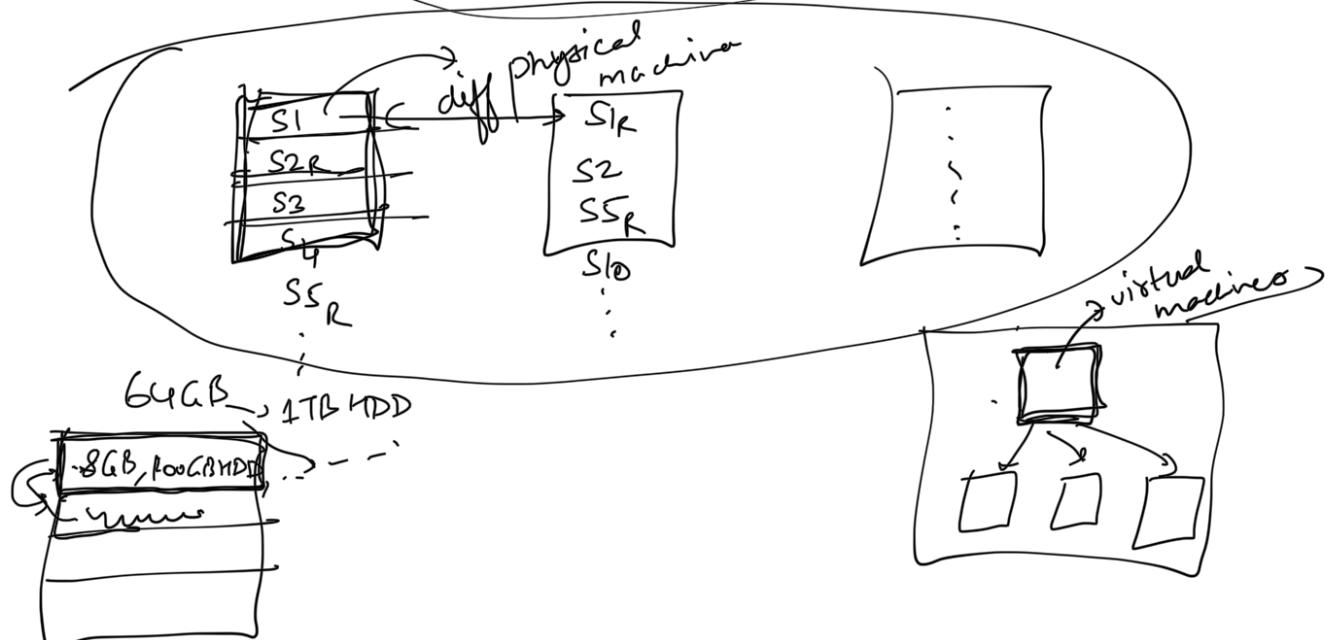
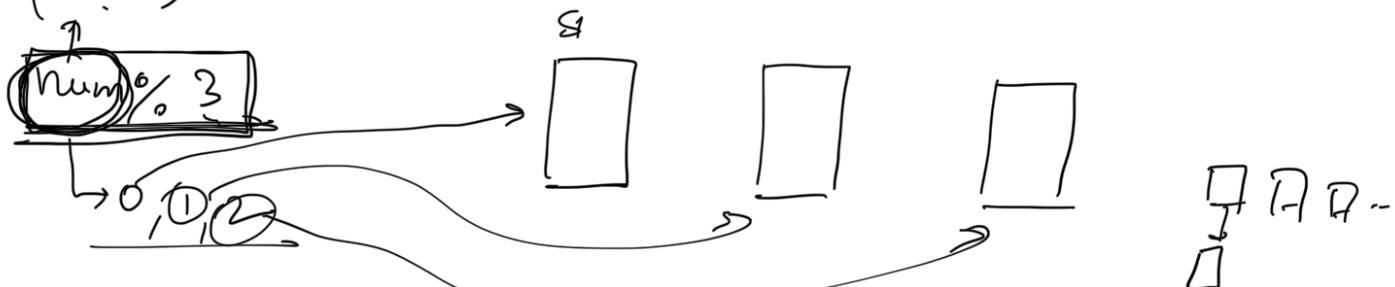


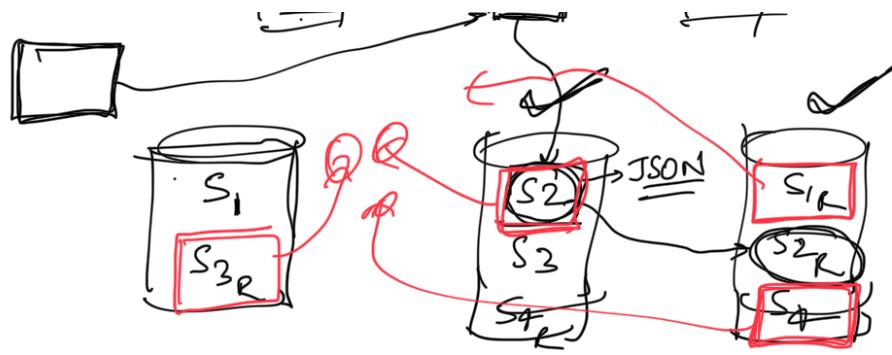
① Sharding - Key : document-id

Shards : ~~100~~ [50]  $\xrightarrow{\text{hash(doc-id) \% num\_shards}}$

Replicas : (3)  $\xrightarrow{\text{150 machines}}$

$\text{hash(doc-id)}$





(4) shard  
↓  
2 replicates

S1 → S1_K
S2 → S2_K
S3 → S3_K
S4 → S4_K

① Request to all shards

↓  
every shard  
↓  
random pick

inverse hashing  
doc

② Wait for response

HA → low consistency.

ES:

① Indexing takes time

② Good for reads, (HA)

↓ slow for writes.

③ Read

~~michael~~

mikael jackson

Ranking

michael jackson

Regular exp:

Score: --

1.00

# can

michael

0.8

michael jackson great dance

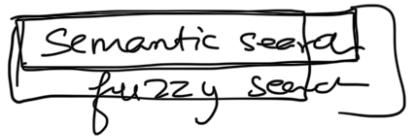
0.5

"Keywords"

→ matching pages

Ranking

→ Page score  
→ Page ranking



- Large narrowing score
  - ↳ Click rate
  - ↳ time spent



if words · contains ("semantic") ·

<sup>e</sup> antihuman

anshumany

antihuman'

add a letter —①

remove a letter - Ø

replace a letter -①

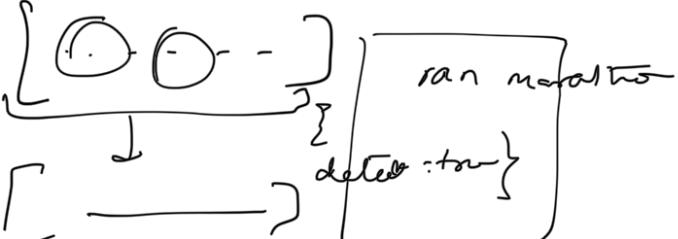
GET \_resumes

text: "van marathu  
OR: { "van",  
"mar"

## Soft delete

permanent delete

$\gamma_1 = -$



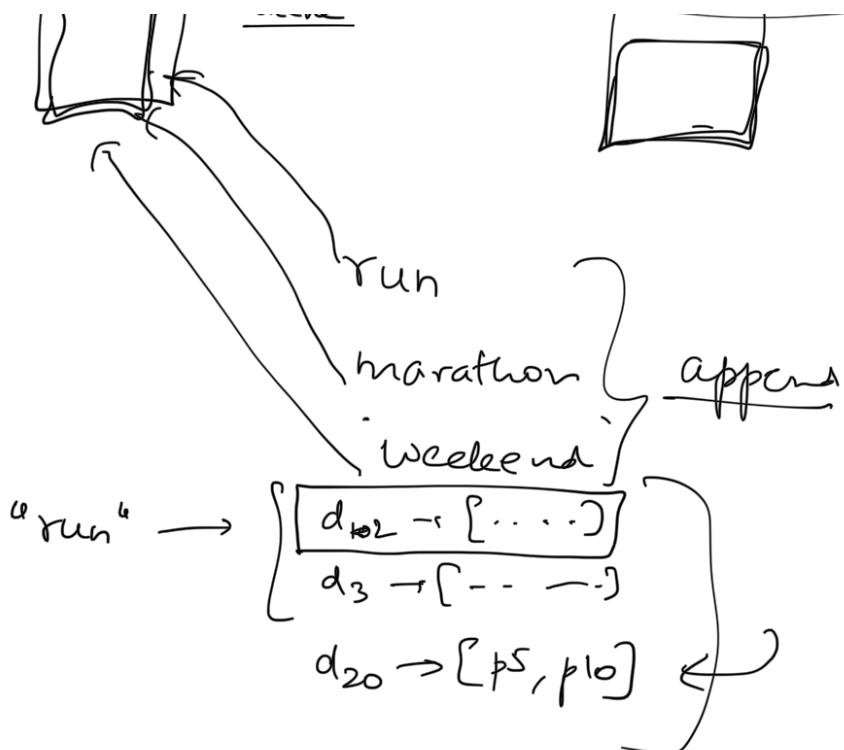
$d_{11} x$   
 $d_{12} x$   
 $d_{10} x$   
 $d_{22} x$

S1

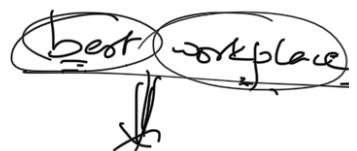
doc, docs . - doc<sub>20</sub>

Lucene

$S_2$  —  
 $\boxed{d_{10}, d_{12}} \boxed{d_{15}, a_2}$



best → good



Good workplace

