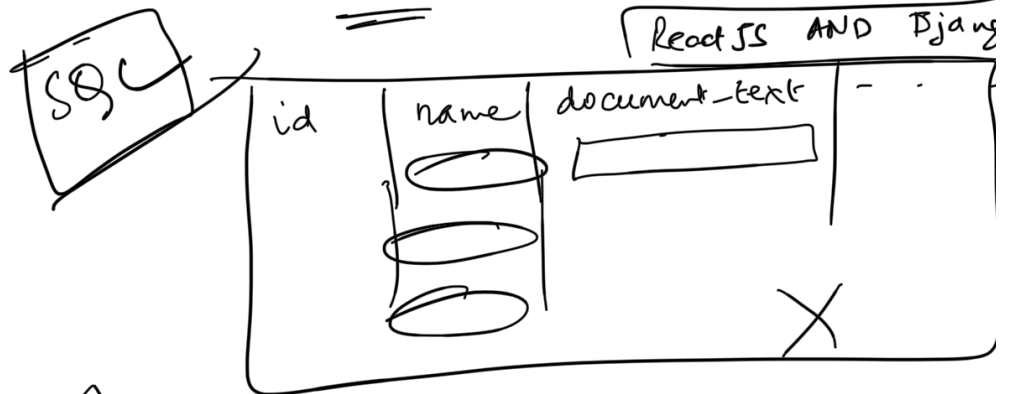
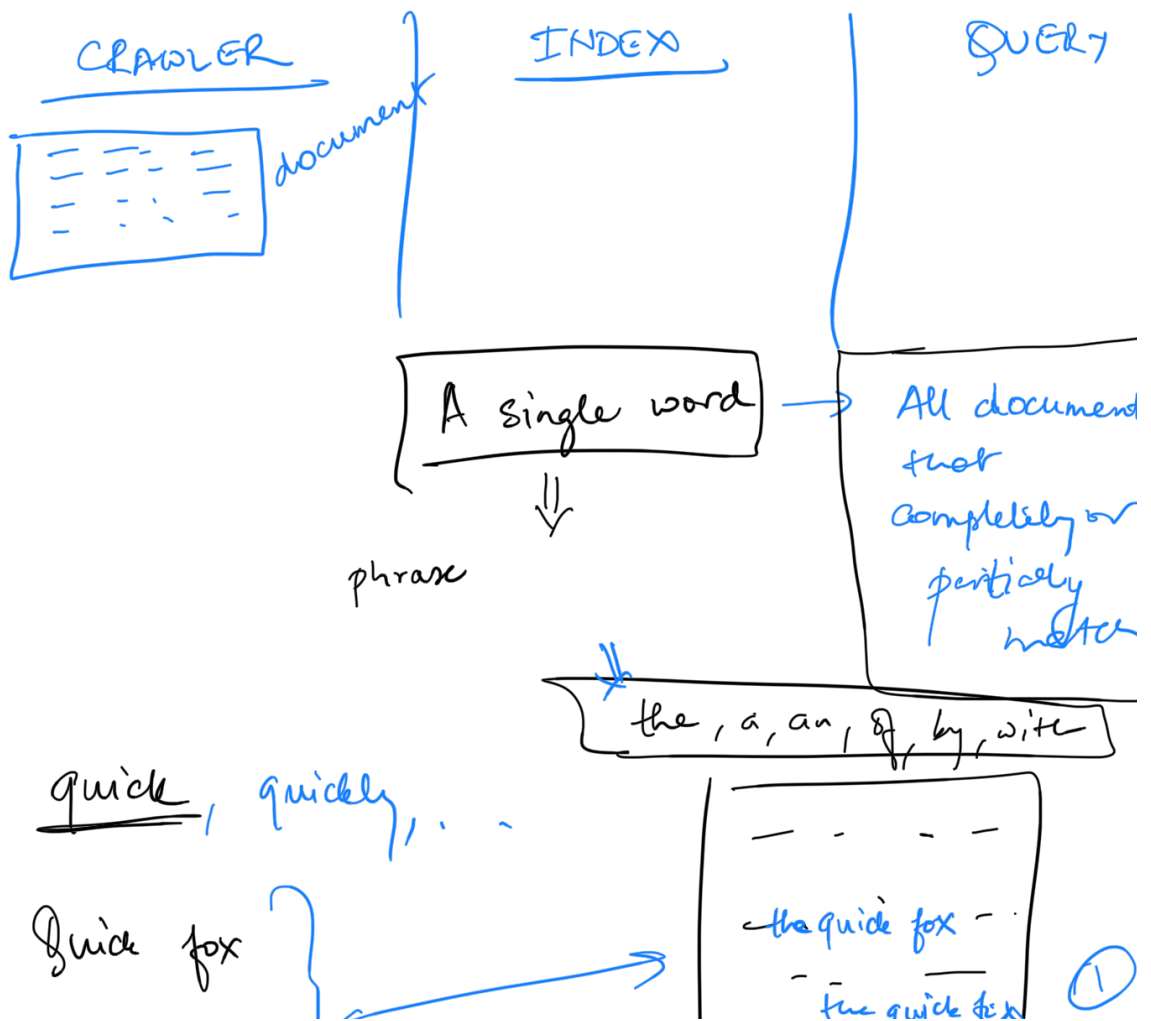


TEXT SEARCH SYSTEMS



Google Search



A quick for

think thought

quick for

word → document, freq, positions

ice

① → stemming → running, runs, run, ran
run

Apache Lucene

② get rid of insignificant words

Ram 15 running

Ram ran

Ram run

Index

word → [(docid, position),
]

Ram ate maggi
eat

run → [(1, 10),
(1, 15),
(2, 3)] documents
large

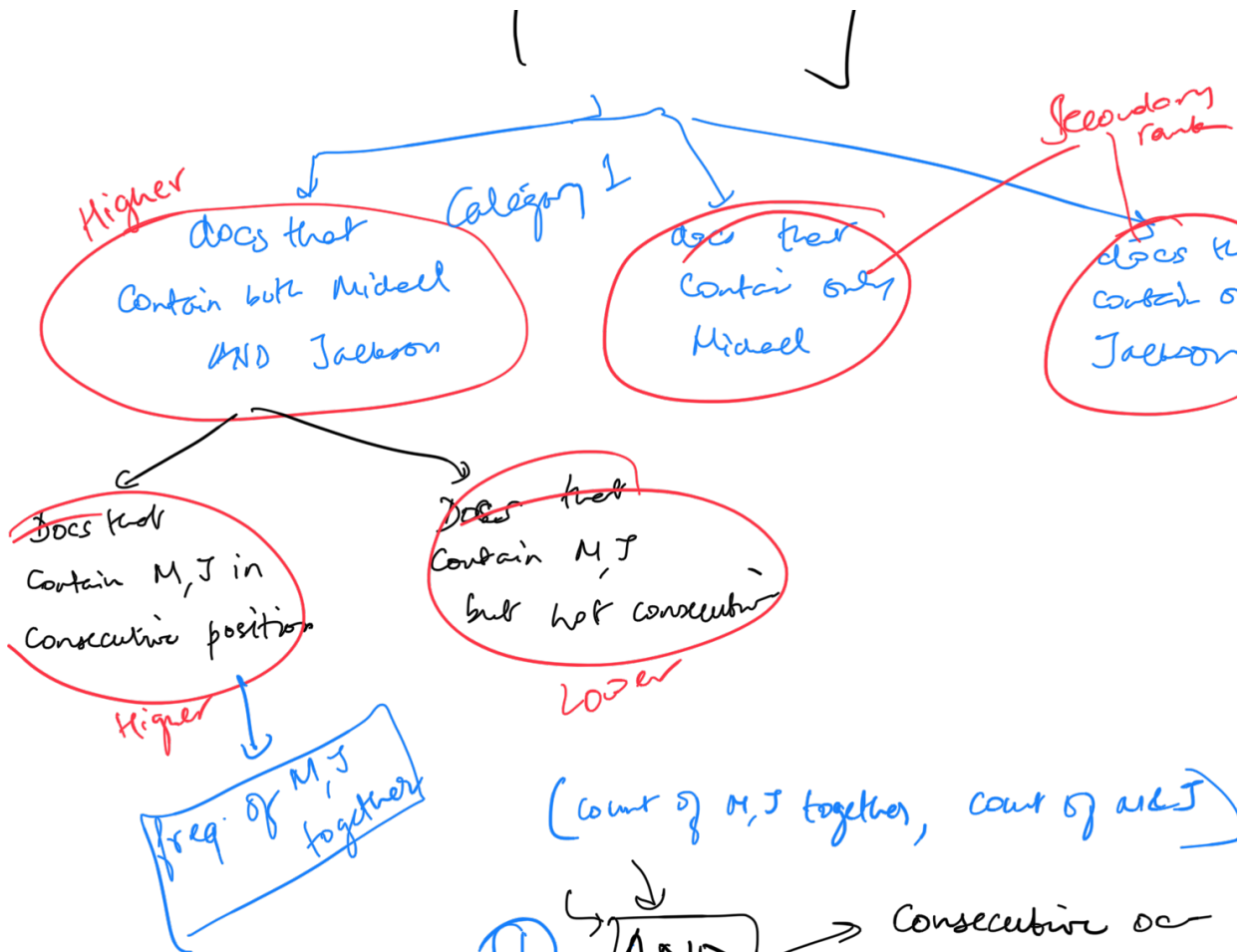
word
↓
documents
+
frequency per
+
position of occurrence

Ram eat maggi

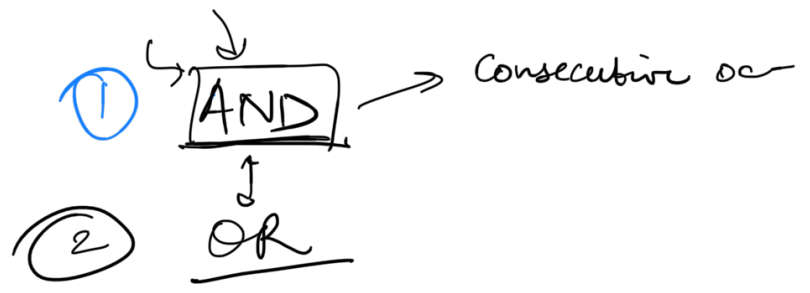
Ram
eat

* Michael jackson

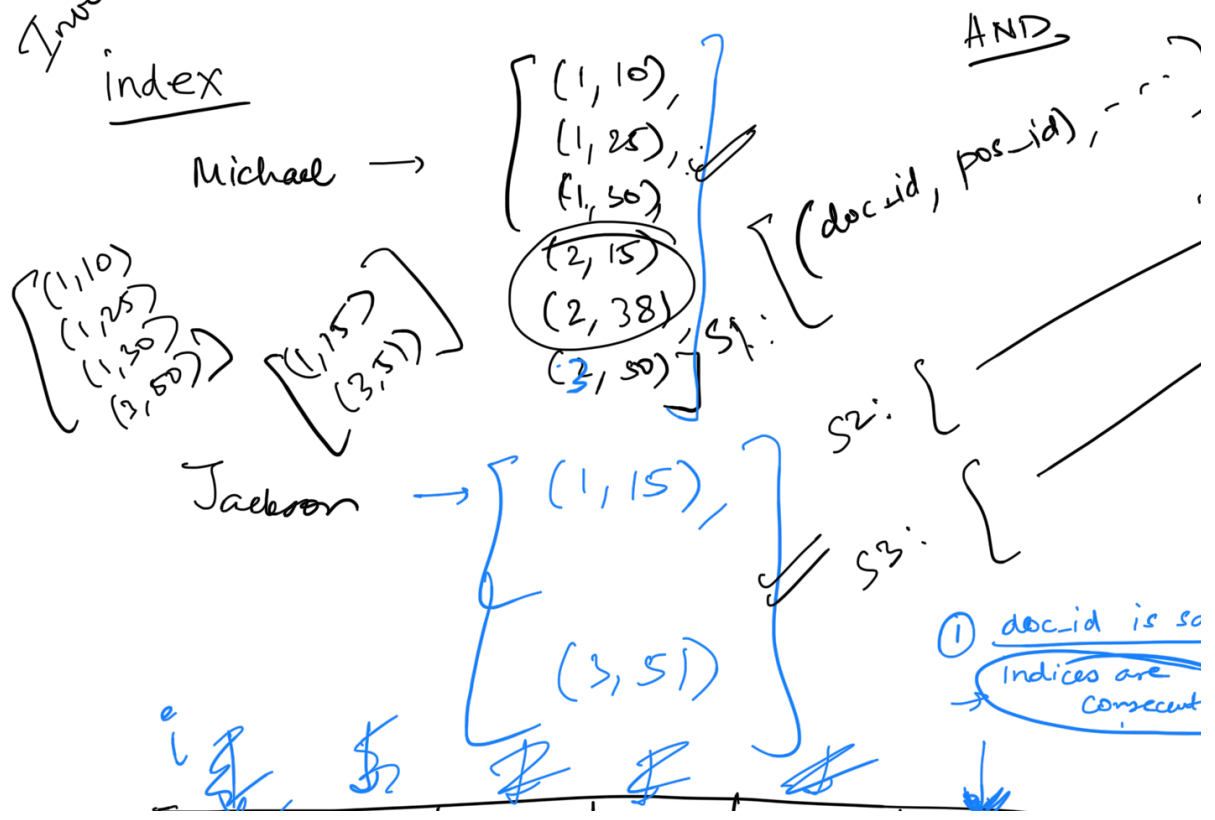
Index
word → [(docid, pos-id)]



(count of M, J together, count of and J)



Inverted index
index



A1:

(1,10)	(1,25)	(1,30)	(2,15)	(2,38)	(3,50)
--------	--------	--------	--------	--------	--------

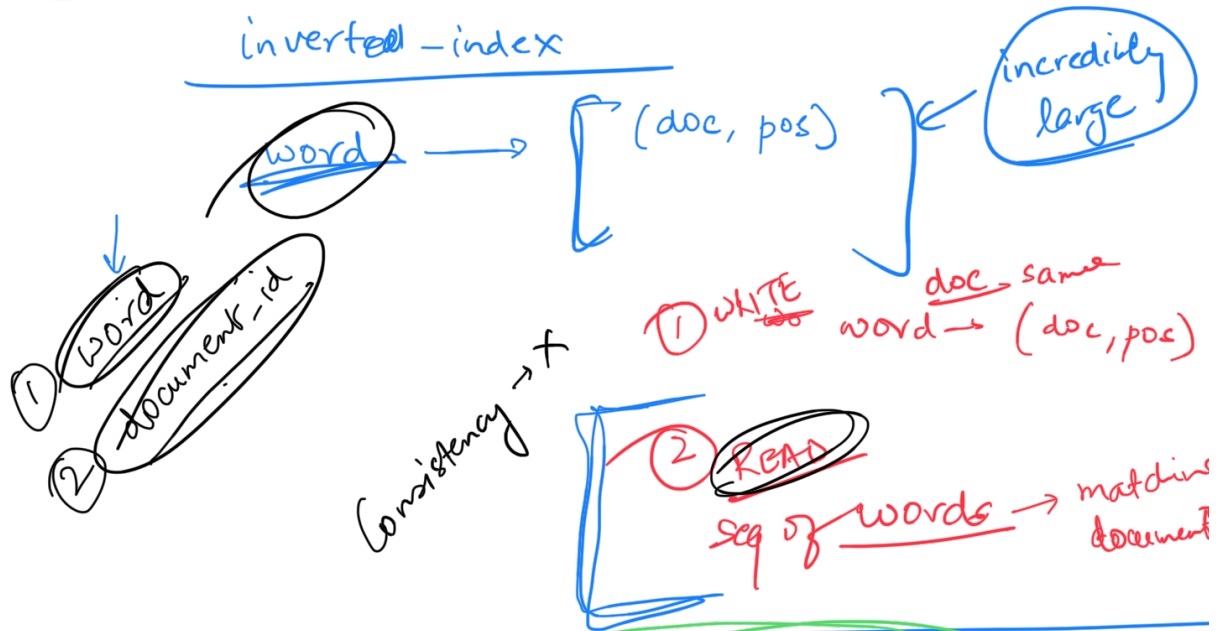
A2:

(1,15)	(3,51)
--------	--------

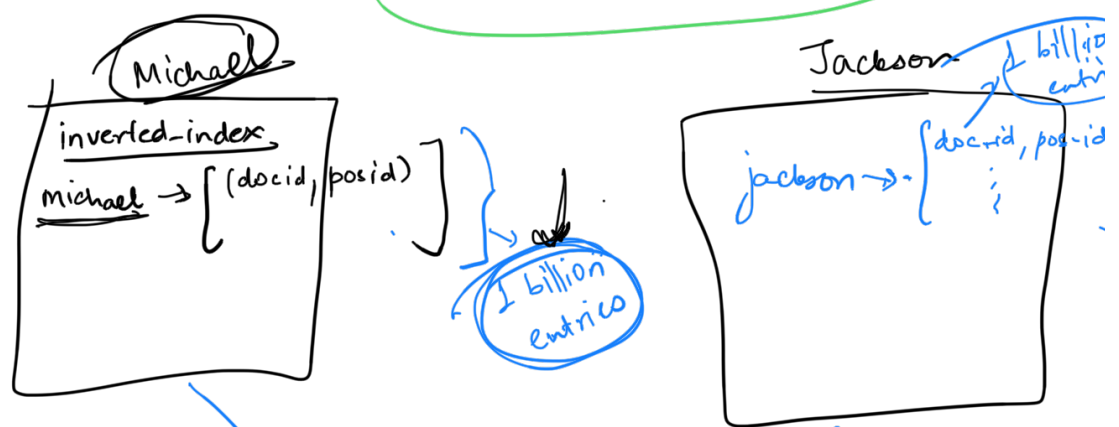
Handwritten notes below A2:
 $\$2$ \uparrow

Special-case $[docid] = 1$
 \rightarrow relevant $[docid] = 1$

inverted-index



word based sharding **SLOWER**

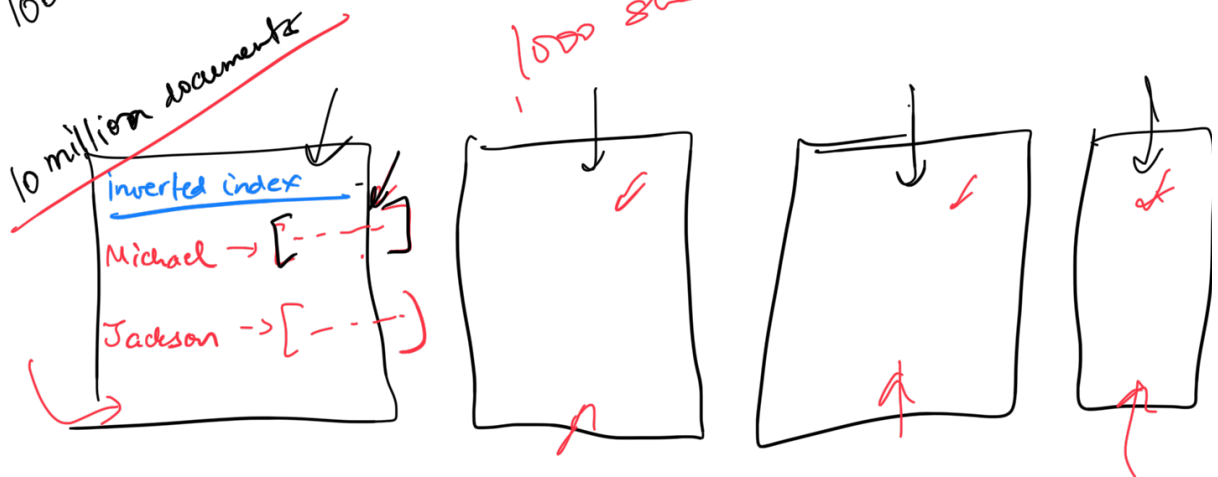


① Go to "Michael" shard and fetch value for ...

Common docs - does with conc

② key = mic
 Go to Jackson doc
 and fetch value for
 key = Jackson for
 step 1 and step 2 complete:
 - 2 pointer iterate on
 both arrays
 and find intersection

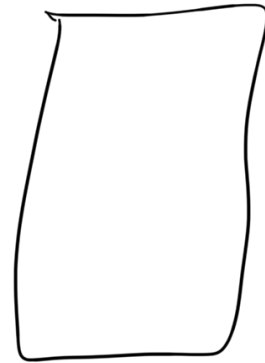
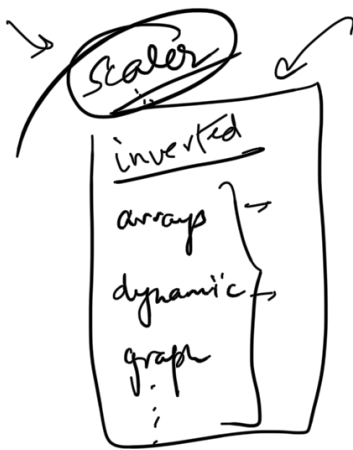
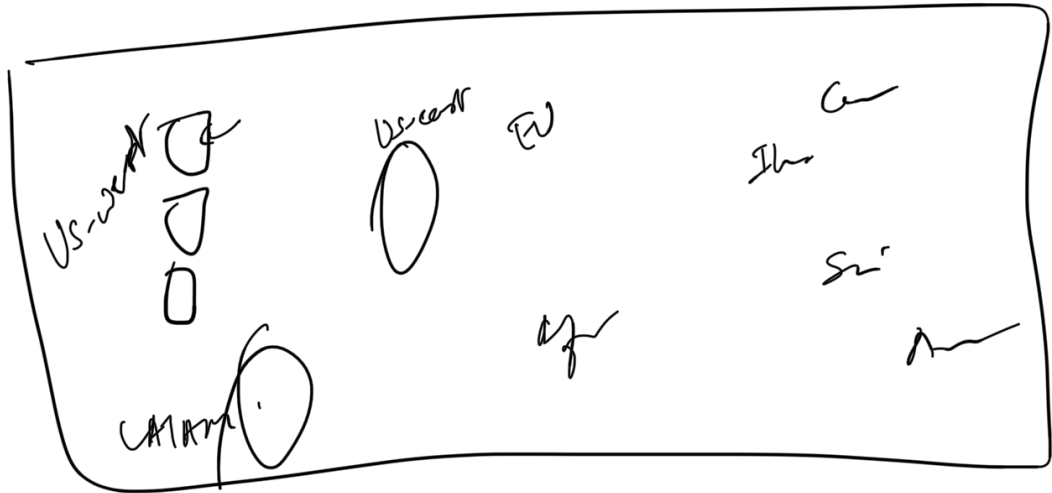
1000 shards Doc based sharding FASTER



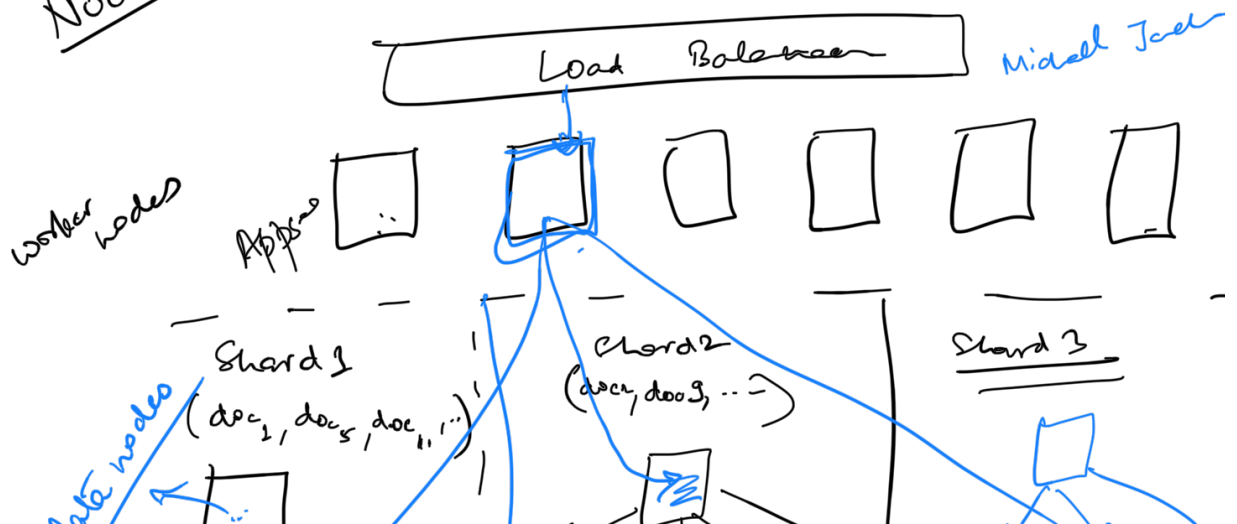
① In every shard, in parallel, 2 pointer iterate on value arr and find intersection 0.1-0.2 sec

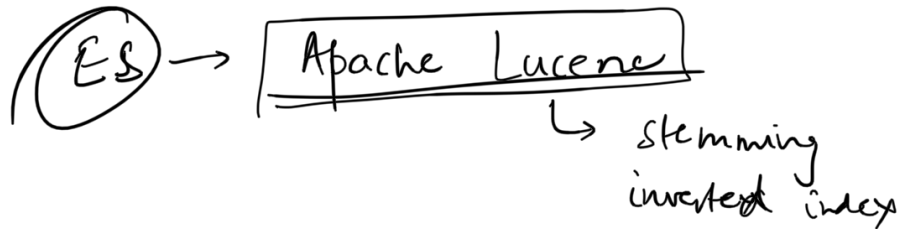
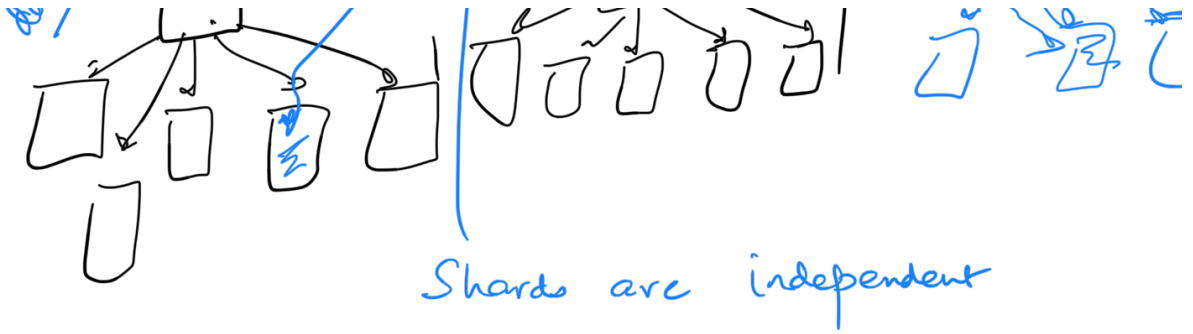
② Wait for shards 0.5 seconds
 → (doc-id, x, y, z)



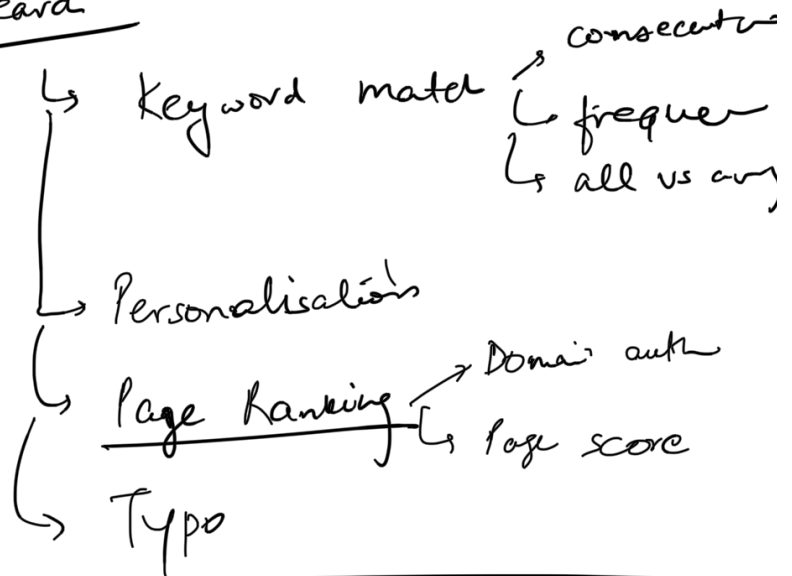


Nodes





Google Search

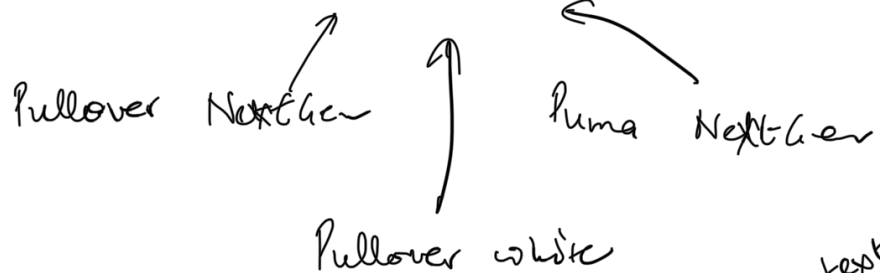


User → resumes

Text Search

Keywords → resumes
Users

"Puma Men XL Pullover NextGen white"



Product name LIKE "Pullover"

Am pr Like "y. Noel"

ES → { product title { screws
product door { table
product tags { small

List of meta docs of matching score

score > (450)

text search
string

Image Text

