

## Elastic Search

- full text search solution  
document database

Pragy

## Full Text Search

Amazon Reviews - find all reviews about "product quality."   
ambiguity

SQL

User-reviews

id | user-id | review-content | attachment-val | rating-of-5 | time

really durable sk  
would buy again

Product quality is  
filmy

:

Task:

find all reviews containing words "product" and "quality"  
intersection

select \* from user-reviews

where review-content like "% product %"  
and review-content like "% quality %"  
OR?

limit 10  
offset 100;

$n = \# \text{ rows}$   
 $\text{time} = O(n)$

↓  
in respective of  
another index is  
present or not

like "%product%" - can use index (prefix) ✓  
sorted list of strings

like "% product%" - index(col)  
index(revrev(col)) ↗ efficient

like "% product%" - neither a prefix nor a suffix query

Product was nice ①

awesom product

Poor quality ②

I didn't like product as it sucked

awesom product ③

Poor quality

I didn't like product as it sucked ④

Product was nice

index(col)

I didn't like product as it sucked

"Product%" prefix ①

Product was nice

"%product%" → ① ③  
④

awesom product

Poor quality

index (revrev(col))

"% Product"  
③

Full text search

↳ Search at any location.

Document - MongoDB / couchbase / cockroach db

Key Value - Redis / dynamodb / memcached

Wide col - Cassandra / scylladb / Hbase

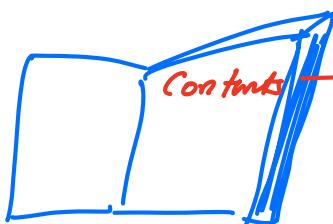
Graph - Neo4j

vector db

blob - S3 / git-lfs / HDFS

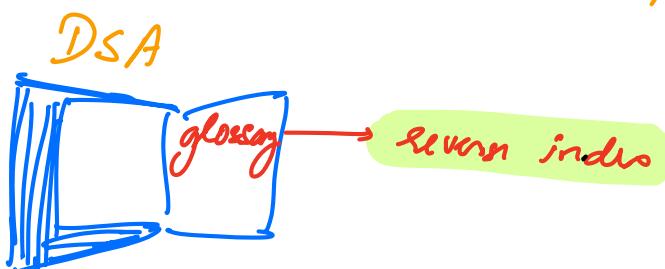
Solution - think of a datastructure not a database

## Inverse Index



index : given a page no what content does it have

Page → topic  
in            out



Word → Page no  
in            out

words      Page no

data : 3, 5, 7, 17, 53, ...

binary search : 23, 34, 57, ...

Product : 5, 11, 27, 29, 27, 23, ...

quality : 6, 11, 25, 26, 27, 53 ...  
~~~~~

find all pages  
containing both  
'product' & 'quality'  
full-text

# dictionary

① What db would be ideal to store this kind of inverse index?

inverted / inverted / reverse



## Building the Inverted Index

NLP : Natural language processing.

~8% of all English words is "the"  
↳ filler word-

"to be able to understand ~~the~~ something"

Synonyms — multiple words mean the same thing!

lead  
— to lead someone  
— verb  
— material

# Cleaning before building the inverted index

## ① Tokenization

Split the text into "tokens" or words

## ② Stop word removal

remove any filler words like an a if of  
and ...

any bad words like fuck ---

I, this, for, my, who, the, is, He, a, is, at, to, from, ...

## ③ Stemming / Lemmatization (NLP Pipeline)

reduce a word to its "stem" in a dumb manner

Cards → card

if (word ends with (es)) {  
 remove s from the word  
}

-es -ing -s

-er -ion -ed -on

Cars → car ✓ good

Caring → car ✗ bad

lion - l ✗ bad

Make a large mapping of words to "roots"

bought → buy

husband → husband

husband

fener)  
Plurality/  
Synonyms

Shell chick.

Plays → play  
 piano → piano  
 Caes → piano  
 wonderful → good  
 fine → fine  
 hymns → song  
 :

↓  
husband

at the end of this card review is just a list of 'useful' words in sequence.

## Build Inverted index

- ① "My Ratty loves this brand of cat food"
- ② "My dog hates this food. It would touch it,"
- ③ "I love cats & dogs"
- ④ "I love cats but hate dogs"

(1, 2)

|                                                                                                    |                                                                                                                                                                                           |
|----------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| ① cat, love, brand, food<br>② dog, hate, food, touch<br>③ love, cat, dog<br>④ love, cat, hate, dog | cat → <u>(1, 1) (1, 6) (3, 2) (4, 2)</u><br>love → (1, 2) (3, 1) (4, 1)<br>brand → (1, 4)<br>food → (1, 2) (2, 4)<br>dog → (2, 1) (3, 4) (5, 5)<br>hate → (2, 2) (4, 4)<br>touch → (2, 7) |
|----------------------------------------------------------------------------------------------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|

## How to make queries?

- imagine the user has entered ~ few words in the query.
- what do we want?
  - all reviews that contain all the words in query?
  - " " " " some " " "
  - all reviews ranked by 'relevance'
  - all reviews that contain the exact subtext of the query.

## Tf-idf score

$$\text{Term Frequency} * \frac{1}{\text{document frequency}}$$
$$tf \quad - \quad \text{inverse df}$$

$(1, 2)$   
cat →  $(1, 1)$   $(1, 6)$   $(3, 2)$   $(4, 2)$   
love →  $(1, 2)$   $(3, 1)$   $(4, 1)$   
brand →  $(1, 4)$   
food →  $(1, 2)$   $(2, 4)$   
dog →  $(2, 1)$   $(3, 4)$   $(5, 5)$   
hate →  $(2, 2)$   $(4, 4)$   
touch →  $(2, 7)$

query : "love kitten"  
love cat

$$\sum_{w \in \text{Query}} \text{tf idf}(w, d)$$

① find all docs that match first query. any words

love → 1, 3, 5  
Cat → 1, 3, 4  
Union ⇒ 1, 2, 5

$$\text{doc 1} \Rightarrow \frac{tf_1(\text{cat}) * \frac{1}{df(\text{cat})}}{2 * \frac{1}{3}} + \frac{tf_1(\text{love}) * \frac{1}{df(\text{love})}}{1 * \frac{1}{3}}$$

$$\text{doc 2} \Rightarrow \frac{tf_2(\text{cat}) * \frac{1}{df(\text{cat})}}{1 * \frac{1}{3}} + \frac{tf_2(\text{love}) * \frac{1}{df(\text{love})}}{1 * \frac{1}{3}}$$

$$\text{doc 3} \Rightarrow \frac{tf_3(\text{cat}) * \frac{1}{df(\text{cat})}}{0.66 * \frac{1}{3}} + \frac{tf_3(\text{love}) * \frac{1}{df(\text{love})}}{1 * \frac{1}{3}}$$

doc 2  $\Rightarrow$

$tf_i(w) \Rightarrow$  how many times does 'w' occur in doc  $i$ ?

$df(w) \Rightarrow$  how many docs does this word appear in?

|        |   |        |
|--------|---|--------|
| apple  | / | myopia |
| common |   | rare   |

|                  |        |
|------------------|--------|
| $d1 \rightarrow$ | apple  |
| $d2 \rightarrow$ | myopia |

$w_1 \rightarrow$  appears in literally all docs

$w_2 \rightarrow$  appears in exactly 1 doc

Query  $(w_1 \ w_2)$

The more common the word is the less information it contains

less interesting it is

$$tf * \log\left(\frac{1}{df}\right)$$

$\hookrightarrow df$  are very large

All of the above is happening within 1 server

Apache Lucene → tf idf algorithm  
+ inverted index Data struct } or  
} 1 server.

Amazon has over 1 trillion reviews

& people search these reviews billion of times / day!

large data → "Sharding"

10:18 → 10:25

---

Single server running Apache Lucene

- ① Single point of failure
- ② it will not be able to store trillion docs
- ③ it will not be able to handle billion searches / day.

Sharding → how to partition across servers  
= horizontal partitioning across servers

Replication

# Sharding in Elastic Search

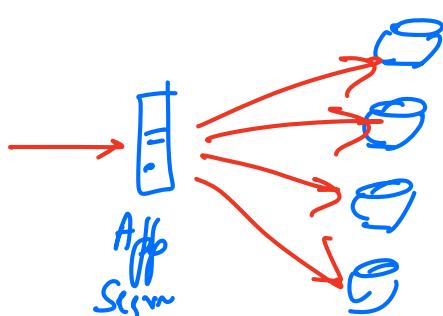
Question : Determine the sharding key

- Sharding key should be part of the query.
- data & load should be evenly distributed
- freq queries should hit 1 or almost a few shards  
not fan-out many shards
- Cardinality of Key should be high  
 $\#$  distinct values  
never choose 'gender' as a sharding key  
 $\therefore \text{Max } \# \text{ users} \leq \text{Cardinality of Sharding Key}$



if you're using user-id as sharding key  
2B users  
 $\Rightarrow$  2B shards? 2B shards  $\times$  No

1000 shards  $\rightarrow$  each shard will store 2M users



What data am we sharding?

① the text documents

Comments  
reviews  
posts  
megs  
answers

② Inverted index (hashmap)

Query

① given doc id  
fetch doc

② given some text  
find

matching doc

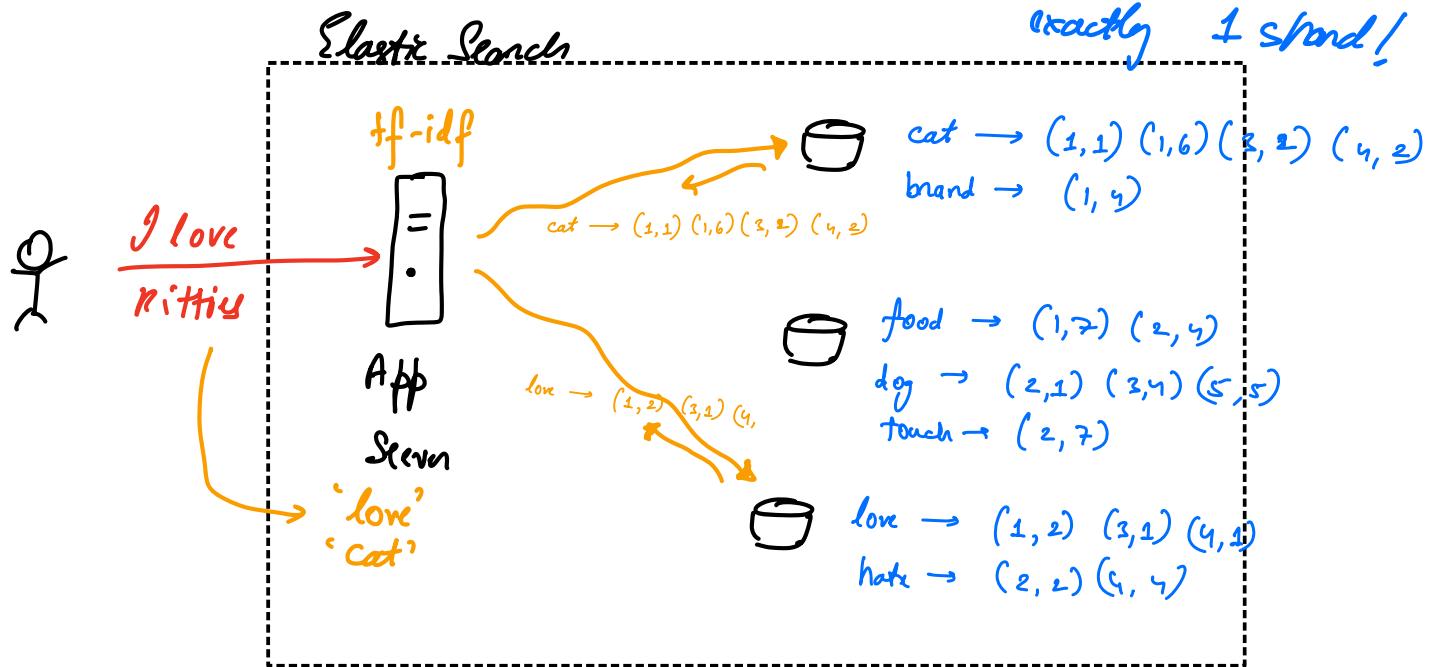
③ given a word find its df

④ given a word find which docs it appears in

① for text docs best shard key  $\Rightarrow$  doc id

② for inverted index best shard key  $\Rightarrow$  ~~word~~ (key)  
Key-value  
bad!!

Attempt 1: shard the inverted index by word (key)  
each word lies in exactly 1 shard!



cat → (1, 1) (1, 6) (3, 2) (4, 2)  
love → (1, 2) (3, 1) (4, 1)  
brand → (1, 4)  
food → (1, 2) (2, 4)  
dog → (2, 1) (3, 4) (5, 5)  
hate → (2, 2) (4, 4)  
touch → (2, 7)

- ① give query process it
- ② for each word  
go to shard  
get data for that word
- ③ collect all this in

## Issues

- for common words entry list is v. large

1 trillion docs  
word  $\approx 10\%$  of all docs

100 R docs  
word  $\rightarrow$  [8 bytes + bytes  
[doc id, freq], [doc id, freq] ...]  
100 R entries

v.v. large n/w overhead

1200 R bytes  
 $\Rightarrow 1.2 \text{ TB}$  data

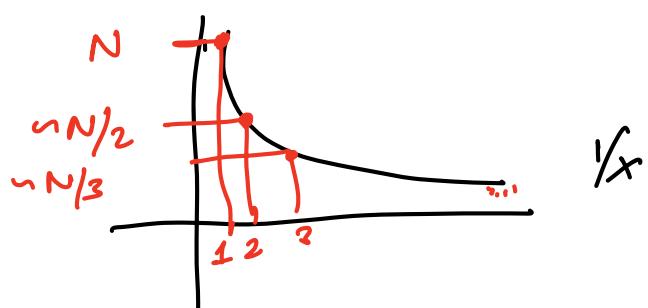
- all this massive amount of data needs to be processed in opp system  
calculate tf idf  
Sort by score

v. Compute intensive

- Zipf's law

Some words are very heavy  
if a shard unfortunately has multiple heavy words

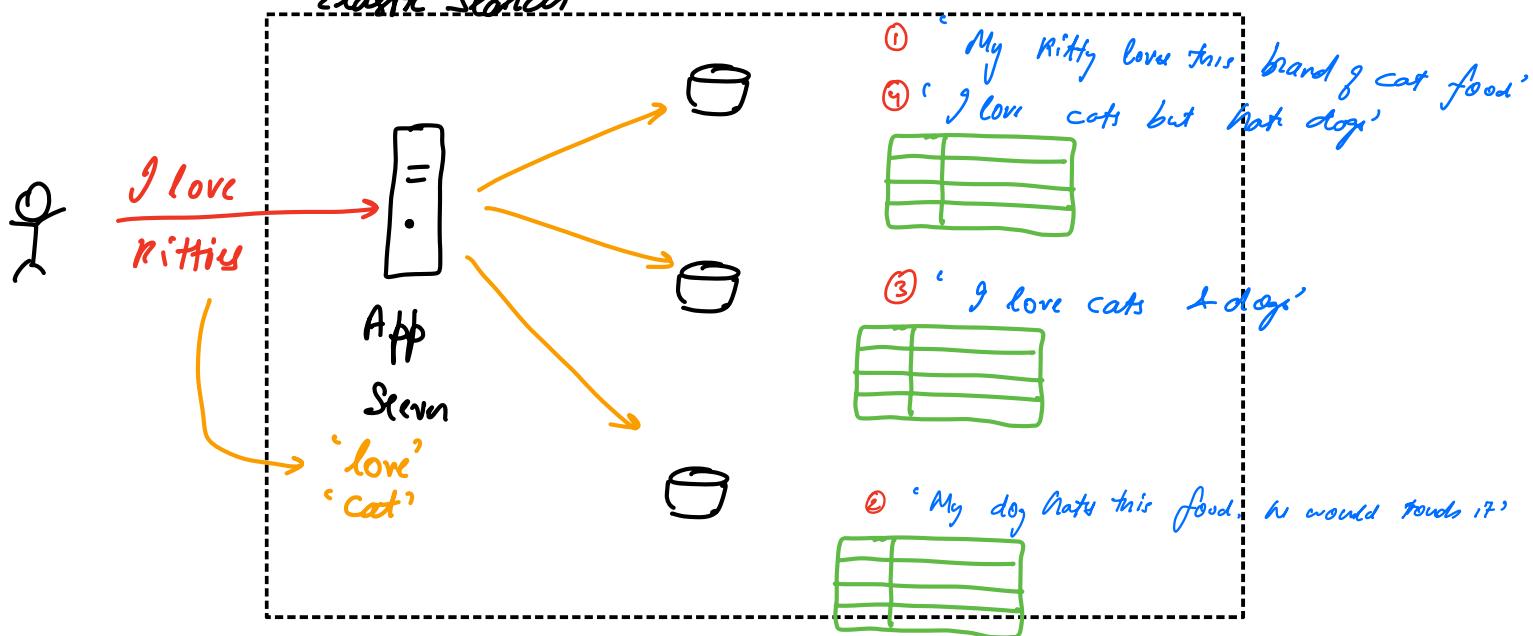
↓  
uneven data dist.



Attempt 2 - shard by the doc id  
inverted index

1 trillion docs  
1000 servers  
each server has  $\approx 1B$  docs

### Elastic Search



① "My <sup>0</sup> <sub>1</sub> Ratty <sub>2</sub> loves <sub>3</sub> this <sub>4</sub> brand <sub>5</sub> <sub>6</sub> <sub>7</sub> of cat food"

② "My <sup>0</sup> <sub>1</sub> dog <sub>2</sub> <sub>3</sub> <sub>4</sub> <sub>5</sub> <sub>6</sub> <sub>7</sub> <sub>8</sub> nati this food. <sub>9</sub> would <sub>10</sub> foods <sub>11</sub> <sub>12</sub> <sub>13</sub> <sub>14</sub> <sub>15</sub> <sub>16</sub> <sub>17</sub>"

③ "I <sup>0</sup> <sub>1</sub> love <sub>2</sub> cats <sub>3</sub> <sub>4</sub> today"

④ "I <sup>0</sup> <sub>1</sub> love <sub>2</sub> cats <sub>3</sub> <sub>4</sub> but <sub>5</sub> <sub>6</sub> <sub>7</sub> hate dogs"

Q1 will the dict (list of words in inverted index) be different for different shards?

No!  $\because$  each shard has  $\approx 1B$  docs  
every word is in each shard!

Q2 Which shard contains data relevant to a query?

All of them!

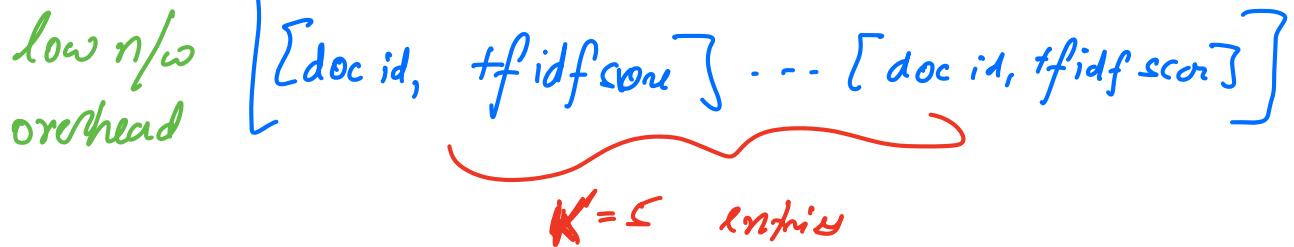
for any read we should send it to all shards!??  
fan-out reads?

① Why are fan out reads / writes bad?

- $Q$  queries get converted into  $N \cdot Q$  queries  
load gets multiplied  $\hookrightarrow N$  shards hit
- we need to wait for all shards to respond  
we are limited by the slowest shard  
latencies are super high!
- what if some read / write fails?  
fail / retry / rollback  
latencies are v. high!  
 $\therefore$  it is v.v. likely for at least 1 shard to  
be down at any given time all queries  
will always fail!

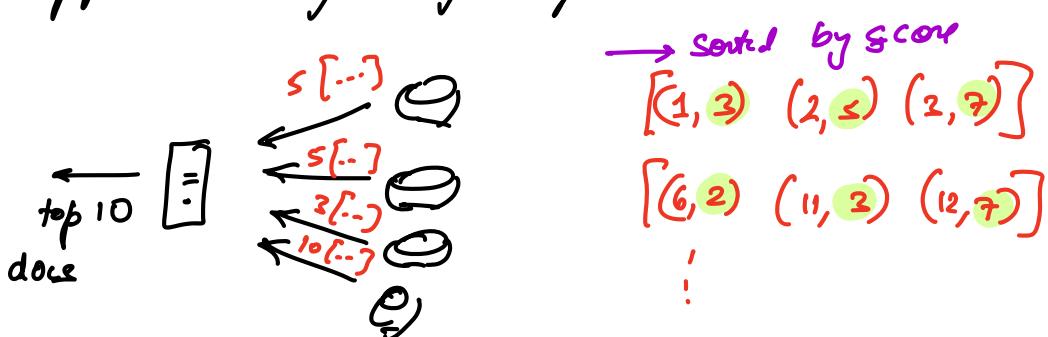
① when a query is forwarded to a shard what does shard return?  
sharding key  $\rightarrow$  doc id

top 'k' most relevant docs that matches this query  
from this shard!



② When is the major computation happening?  $\xrightarrow{\text{App db}}$  inside the shard

③ How does app server finally respond?



fast : K-way merge (min heap)  
 $\because$  tiny data ( $K \cdot N$  entries)

$\downarrow$   $\hookrightarrow$  # 8 shards  
 $\# 7$  docs returned by each shard

$$K = 5 \\ N = 1000 \\ K \cdot N = 5000$$

④ Do we necessarily need to hit every shard??

stands "diabolical cracker bar"

10B facets that 'match' this query.

User wants highly relevant results

↳ do not care about exact ranking  
 ↳ care about a handful of relevant results  
 if a relevant result is not shown

that's okay as long as whatever was shown is also super relevant

Configure for each query → how many shards to hit!!

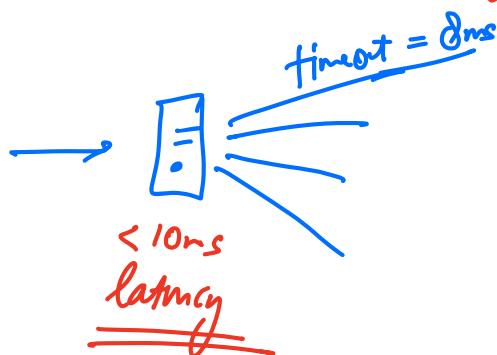
⑤ if we send query to 20 shards

15 respond quickly

5 respond slowly

do we need to wait for slow shards? No!

We can control the latency 'exactly' !!



Sharding → doc → doc id  
→ inverted index → doc id

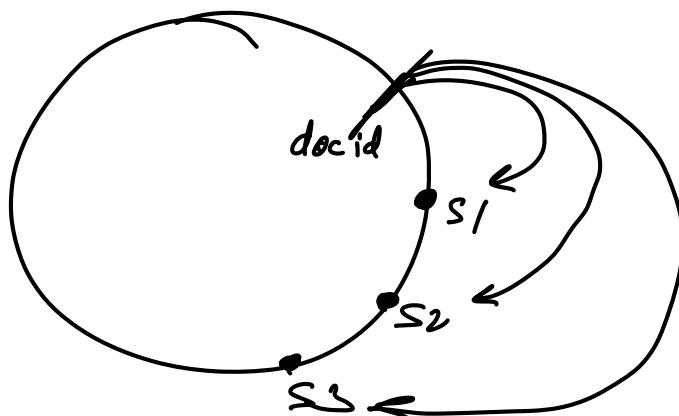
Elastic Search  
is sharded by  
doc id

∴ Elastic Search is a  
doc DB!!

## Replication

Each doc must lie in more than 1 shard

$$X = 3$$



now if we choose shards at random for query it is less likely to miss important doc  
∴ each doc is in multiples shards!

## Availability vs Consistency

### ① Problem

- (1) think of what it means to be eventually consistent in this case
- (2) Availability / consistency

### DS1

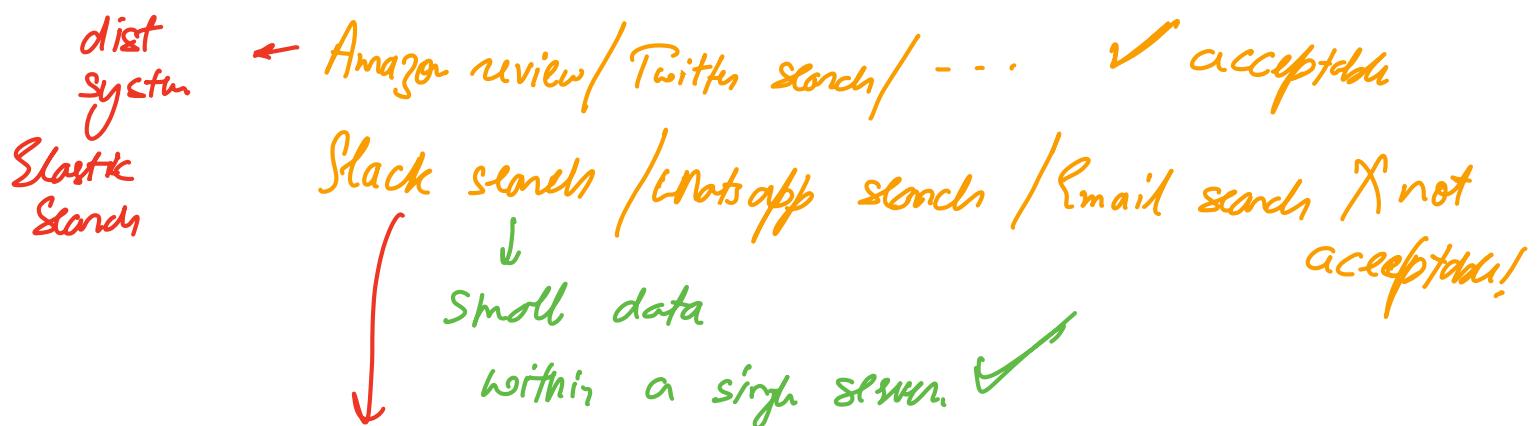
- ① problem  
② but for  
③ observation → magic?  
④ optimized

Consistency

- ① No consistency → data can be lost after commit!  
Analytics / trends  
view count on video
- ② Eventual consistency  
lack of finality  
↳ for some time we can have stale reads  
↳ eventually updated data will show up
- ③ Immediate  
↳ up-to-date data always

# What does Eventually Consistent in context of Search?

- ① all relevant docs might not show up in 1 result  
but if we do the same search multiple times all  
relevant docs will eventually show up in  
some result.



## Apache Lucene

1 server with multiple  
read replicas  
write → all replicas  
to ensure consistency!

100,000 user

each msg is 1 KB

5 TB you can have  $\frac{5 \times 10^12}{10^2}$  msgs

$$= 5 \times 10^9 \text{ msgs}$$

$$\frac{5 \times 10^9}{10^5} \text{ msg/user} = 5 \times 10^4$$

Company with 100,000+ emp  
where each emp has set > 50,000  
msg?

as of Feb 2024

1 server that fits  
on your table

Typical servers

Max servers that many  
can buy.

CPU 2-16

SSD 512 GB - 8 TB

RAM 16 GB - 64 GB

N/W 100 Mbps - 10 Gbps

50,000 Rs - 2 lakh  
Rs

CPU ≈ 400 cores

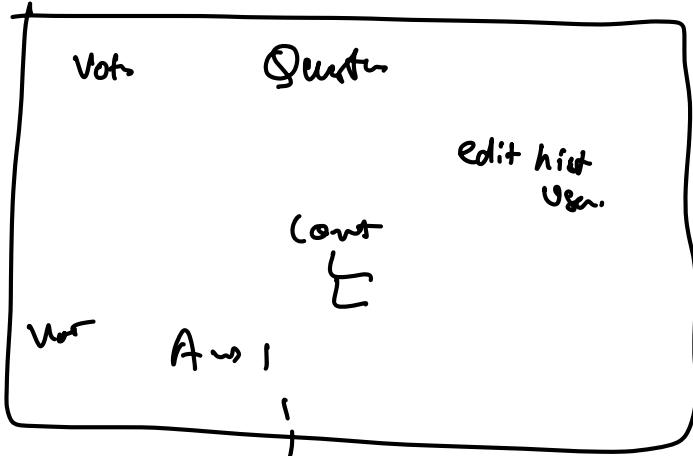
SSD ≈ 4 PB (4000 TB)

RAM ≈ 12 TB

N/W ≈ 10 Tb/s

200,000 \$  
1.5 crore Rs

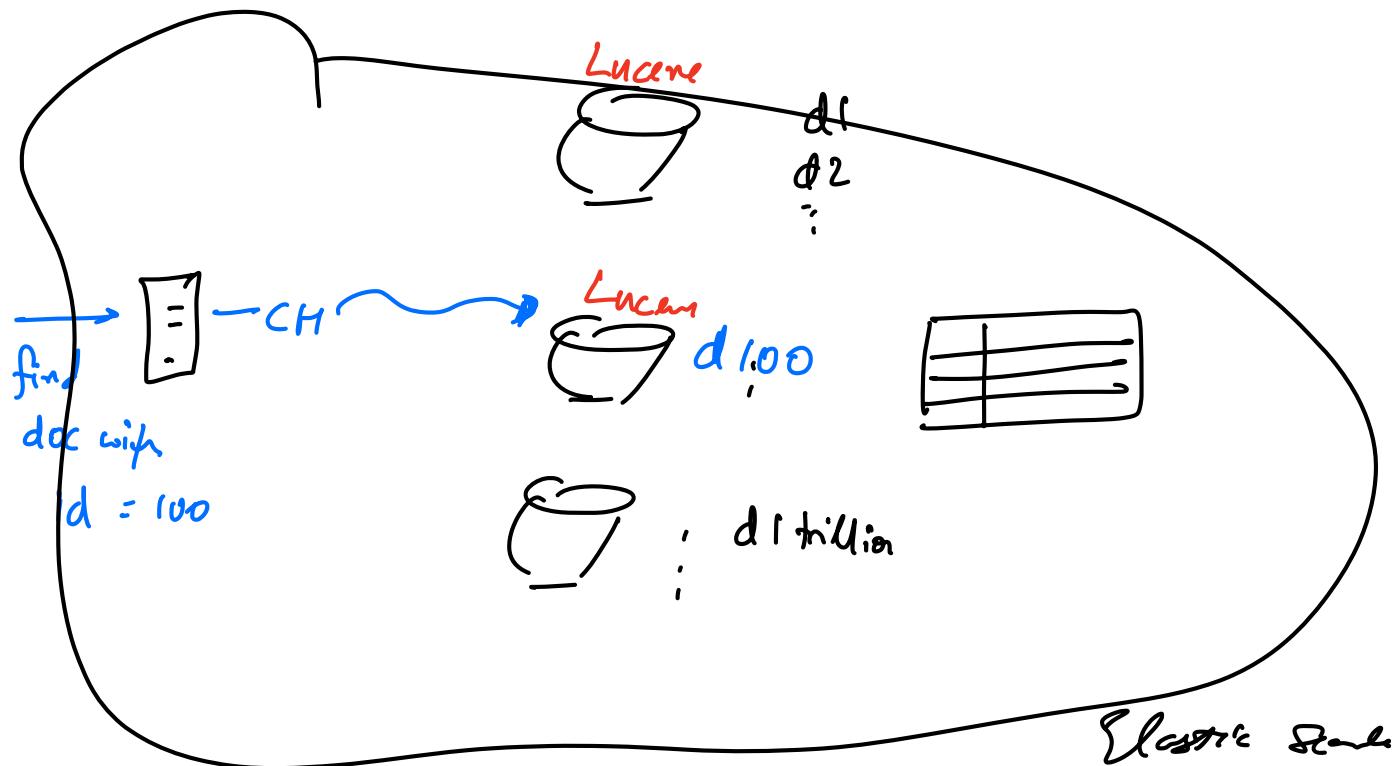
Stack overflow → SQL



100 TB RAM based  
cache

denormalized.

IT



## Partitioning

Split the data into multiple 'parts'

### Users

| id  | Name | Phone | Password | Email | Age | Gender | Pic URL | Nickname |
|-----|------|-------|----------|-------|-----|--------|---------|----------|
| 1   |      |       |          |       |     |        |         |          |
| 2   |      |       |          |       |     |        |         |          |
| 3   |      |       |          |       |     |        |         |          |
| 4   |      |       |          |       |     |        |         |          |
| :   |      |       |          |       |     |        |         |          |
| 100 |      |       |          |       |     |        |         |          |

## ① Vertical Partitioning (Normalization)

## Users

id password email

## Profiles

user-id name phone age gender pic\_url nickname

## ② Horizontal Partitioning.

India

### Users - India

| id | name | phone | password | email | age | gender | pic_url | nickname |
|----|------|-------|----------|-------|-----|--------|---------|----------|
| 1  |      |       |          |       |     |        |         |          |
| 2  |      |       |          |       |     |        |         |          |
| :  |      |       |          |       |     |        |         |          |

$\lg(n)$  is small but not on disk  
disk seek  $\approx 10\text{ms}$

Horizontal partition

USA

### Users - USA

|     |  |
|-----|--|
| 3   |  |
| 100 |  |
| 4   |  |

When you horizontally partition data & you distribute  
some rows in 1 part  
other rows in other part  
across multiple servers

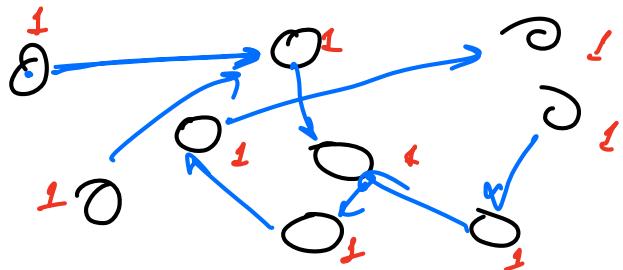
||

Sharding.

cleaning / inverted index / tf - idf

Information  
retrieval

- ① AdSense integration
- ② Page rank algorithm



- ① full blown NLP pipeline ML
- ② N-grams

'This brand of cat food is awesome, my kitty adores it,'

1-gram

This  
brand  
of :

2-gram

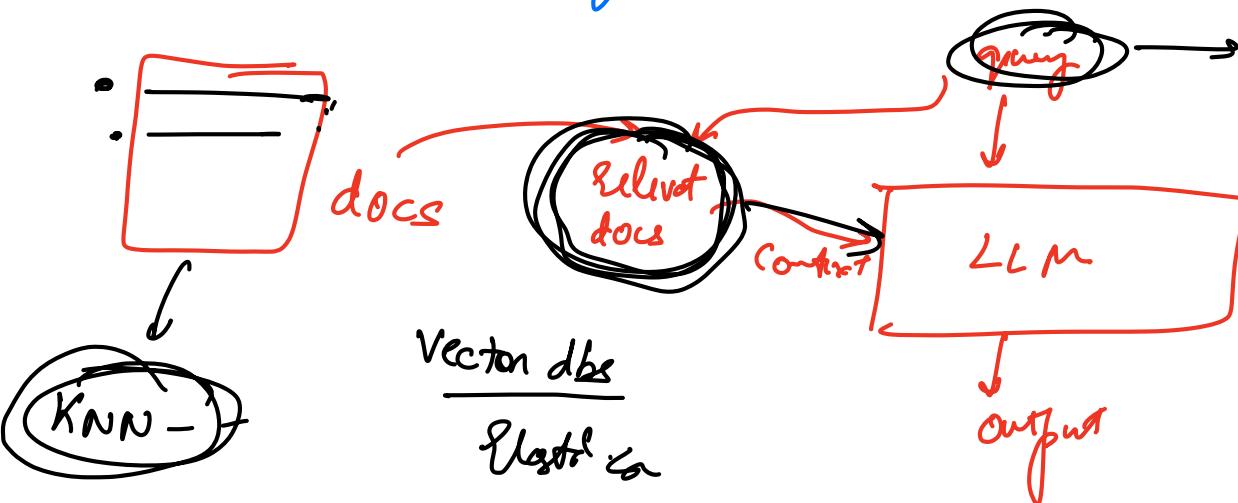
This brand  
brand of  
of cat

3-gram

This brand of  
brand of cat  
of cat food

# Retrieval Augmented Generation (RAG)

with knowledge Gen AI



7351769221