

Agenda

- Full Text Search →
 - a bunch of docs
 - search query (text)
 - find all / top documents that match the query.
- Elastic Search
- Sharding

Problem Statement

- large number of documents
- search documents (text query → search)

→ Searching Logs

→ E-commerce Website

- ↳ Product title + description + attrs
- ↳ Search Services

→ Google Search X
Web Search Engine

→ Social media / User created content
Tweets / Quora answers / naukri.com
profits on jobs

→ Scala problem descriptions

→ Slack

→ Google Drive → Search all docs

→ Gmail Search

g / re/p → print
global { } Regular Expression

Amazon Reviews → Database

SQL

No SQL → Document DB

Product-reviews

id	product-id	Content
		↑ add index

Sorted by content
 B+ Tree
 ↳ only prefix queries

Select * from prod-rev

where content like query

N rows

Content → M chose on avg

$O(N \cdot M)$ ↴ → no index
 billion docs

find any string containing for phonew "value for money"

where content like "% value for money %"

index cannot be used.

find all docs that contain word [W]

tree / sorted \rightarrow prefer
no loss

yes/no

Inverted Index \rightarrow Map

- 1 \rightarrow "9/at/at/this/restaurant/ Food/was/really/nice"
- 2 \rightarrow "ambience/was/awesome"
- 3 \rightarrow "service/sucked"

Key

Query:
awesome

- 9 \rightarrow [(1, 0)]
- atc \rightarrow [(1, +)]
- at \rightarrow [(1, 2)]
- this \rightarrow

restaurant \rightarrow

- food \rightarrow [(1, 6), (2, 1)]
doc id *position*
- was \rightarrow [(1, 6), (2, 1)]

- really \rightarrow
- nic \rightarrow
- ambien \rightarrow
- aweso \rightarrow

\rightarrow all docs that contain the word

awesome.

Sucked

the : 18% of all English Text is the word 'the'

Index/Table of contents

Chapt 1 : \equiv

Chapt 2 : \equiv

load Bal : -----

hashing : -----



Apache Lucene / Elastic Search / Solr / ---

① Remove "STOP words" → useless words that appear often → e.g. "the" "a" "an"
"of" ---

1 → "eat/⁰/at/¹/for/²/restaurant/³/Food/⁴/really/⁵/nicer/⁶"
2 → "ambience/⁰/at/¹/awesome/²"
3 → "service/⁰/sucked"

↓
special char ↗ X

1 → eat restaurant food really nice
2 → ambience awesome
3 → service sucked.

② Stemming / Lemmatization

takes context into account
regular ML pipeline → super slow

"develop" → []
"developing" → []

-e
-ing
-s
-ed
-ion
-tion

word → root

escaped → escap
escaping → escap
escape → escap } stemming
} v. fast
Simple

ration → rat
 Caring → Car
 nation → na
 ↳ not a suffix
 lion → λ
 ↳ not a suffix

precipitate
 precipitation
 ↳ suffix

③ Other cleaning

- ↳ abusive words f***k
- ↳ sensitive words bomb kill
- ↳ hate speech feminist n***g*r

Inverted Index is now built.

Scara Query "The woman with the white dress"
 ↓ cleaning
 "Woman white dress"
 "data"
 "bcg"
 "CAT"

① all docs that contain ALL of the words

wom → }
 white → } intersection
 dress →

② all docs that cont ANY of them

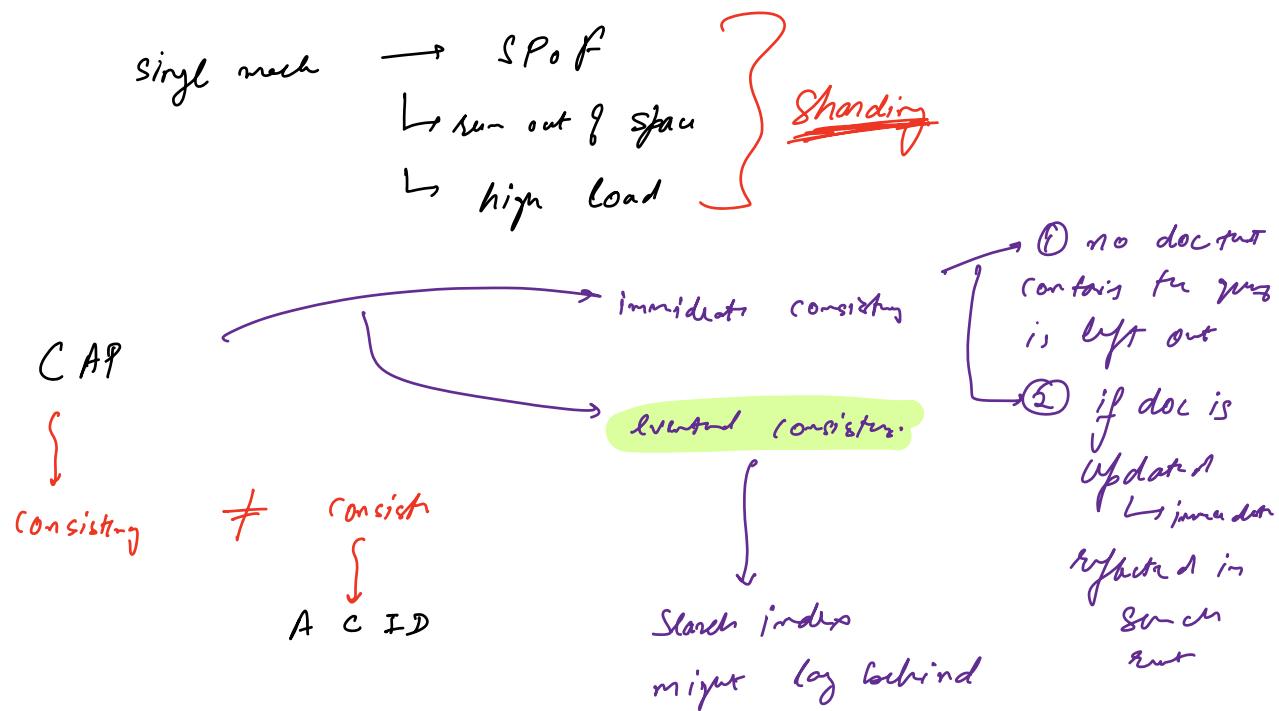
≡ } Union

③ all docs w/ all the words in this exact order

woman → i ?

addit. $\rightarrow i+1$
 dress $\rightarrow i+2$

\approx \approx
 \downarrow
 exact phon



- ① Shard by **document-id**
- ② Shard by word
- review-id
 → log-id
 → prod-id
 → post-id / tweet-id



insert / update / delete document → doc-id
 ↓
 shard
 ↓
 re-index doc content.

Search → go to all shards. → exact query is performed at every shard
 "woman white dress"

collect them all ← each shard will return a list of docs that match

Map - Reduce

items



for x in items:
Process (x)

Billion
Items in
memory, in 1
go
large computers

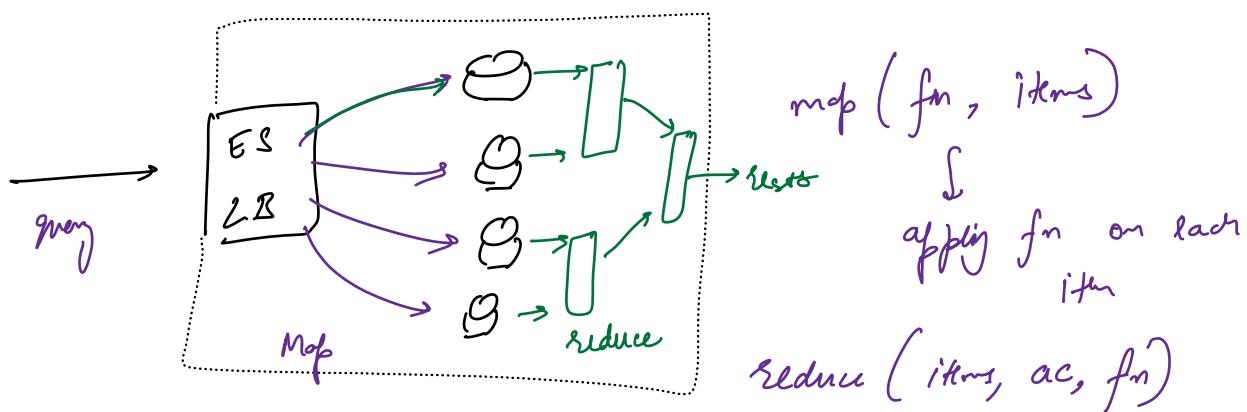
items



map → single

reduce → two items

easy to fit in
memory



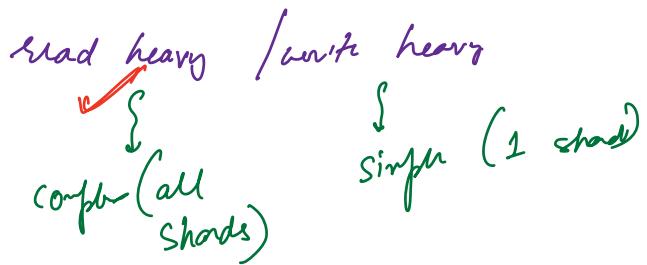
`reduce (nums, 0, $(a, b) \rightarrow a + b$)`

\downarrow reduce

`reduce (shards, [], $(a, b) \rightarrow [\underbrace{\dots a, \dots b}]$)`

\downarrow (concatenate) / interaction

Cons of doc-id sharding

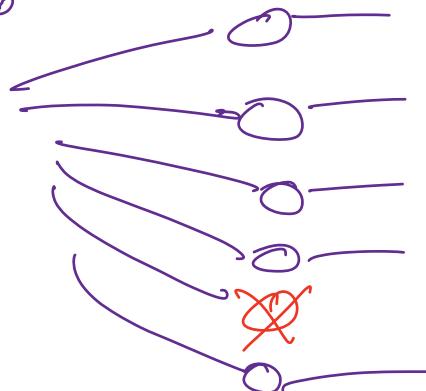


Pros of doc-id sharding

→ latency is highly predictable → ∵ response time don't depend on query

→ no need to do intersection

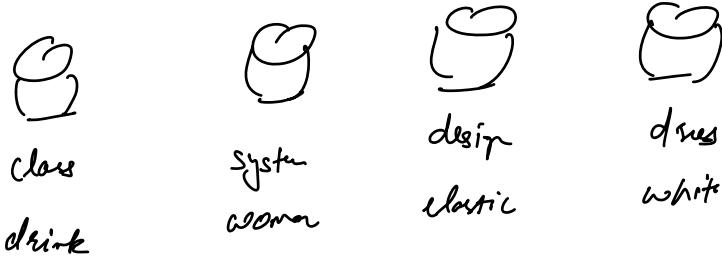
→ can skip unresponsive/slow shards



10:39 → 10:50

Open Source → Lucene / Solr / Open Search / PostgreSQL
Proprietary → Elastic Search / Algolia / Azure Search / SAP HANA - -

Word based Sharding

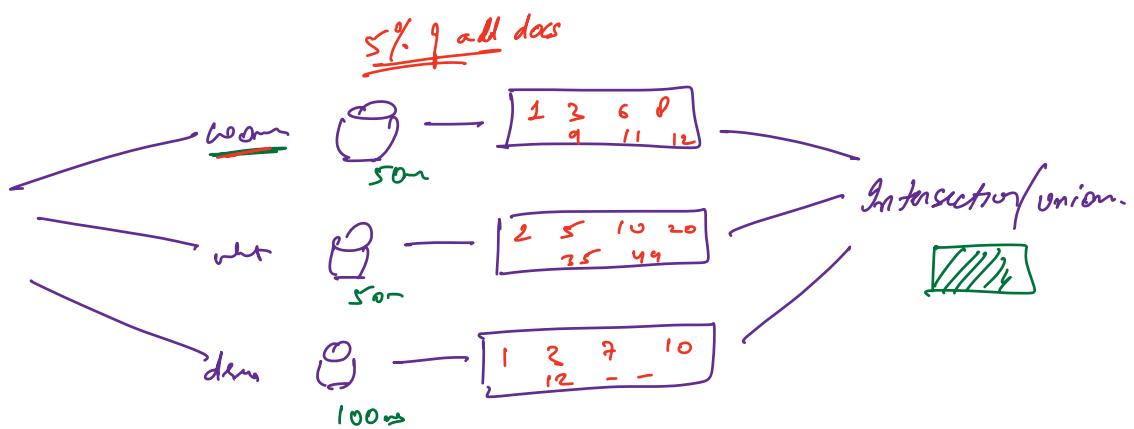


Insert / update / delete

$$\rightarrow \text{doc} = \left\{ \begin{array}{l} \text{The man is} \\ \text{a young man} \\ \text{boy dog.} \end{array} \right\}$$

go to all shards corresponding to the words in the doc → multiple shards

query = "woman white daves" → go to all shards for this query.
 ↴
3 words on avg. v.v. small (< 5)

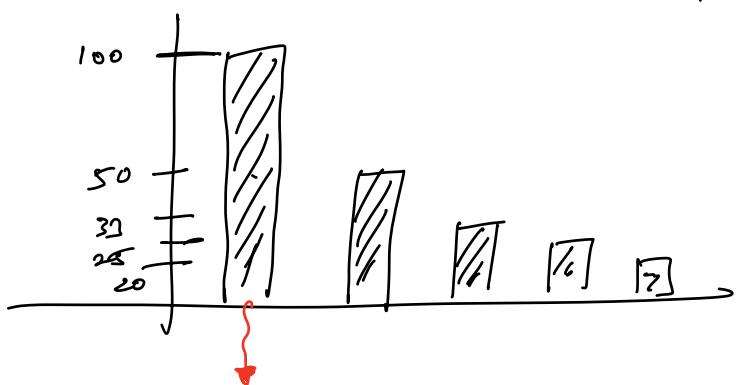


queries with simple/freq words → longer

you with more words \rightarrow longer.

we need to wait for all shards (involved in query) to respond.

Zipf's law - $\sqrt{\text{souce}}$



ith most used
word w/M
occur $\sim 1/i$ times

this shn is ordan

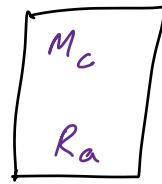
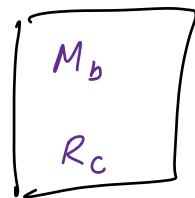
Replication

of shards Replication factor
 { }

N nodes \rightarrow each node is a server in ES.

M shards }
 R replication factor. } $M * R$ items

each server must hold $M * R / N$ shards/replicas

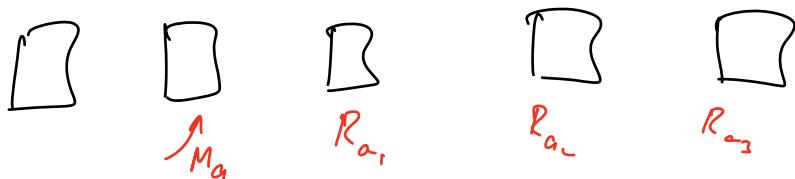


$N = 3$
 shards = 3
 $D_a \quad D_b \quad D_c$
 {
 set 1 \Rightarrow shard

ensure that 1 machine
 only has 1 copy of any shard
 at a time.

replication = 2

Master - Slave Replications



good / avcsn / great / --
 crack

All docs that match a certain query

}
 Ranking.

$TF \cdot IDF$ → Inverse document frequency

Term frequency

→ d_1 → quick fancy help quick → 100 → first quick pair
 → d_2 → gun bra fox TF = 1 → 1 n
 → d_3 → quick pain relief TF = 1
 d_4 → one piece
 d_5 → manato

IF for quick = 3
 $DF(\text{quick}) = 3^1$

$$DF(\text{quick}) = 3^1$$

$$\frac{1}{DF(\text{quick})} = \frac{1}{3^1}$$

$$3^{-1/3}$$

$$1^{-1/3}$$

$$1^1/1^1$$

$$1_2$$

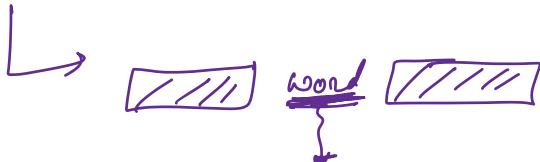
open past → filter → now ~
can → be ~

~~The song supercolifragilisticexpialidocious~~
↓ ↓ ↓
~~stop frag~~

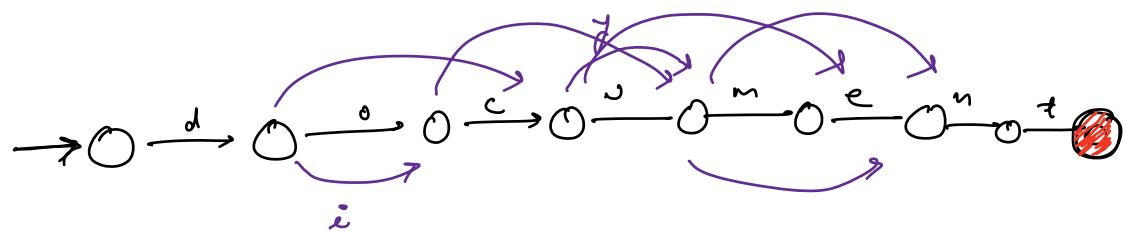
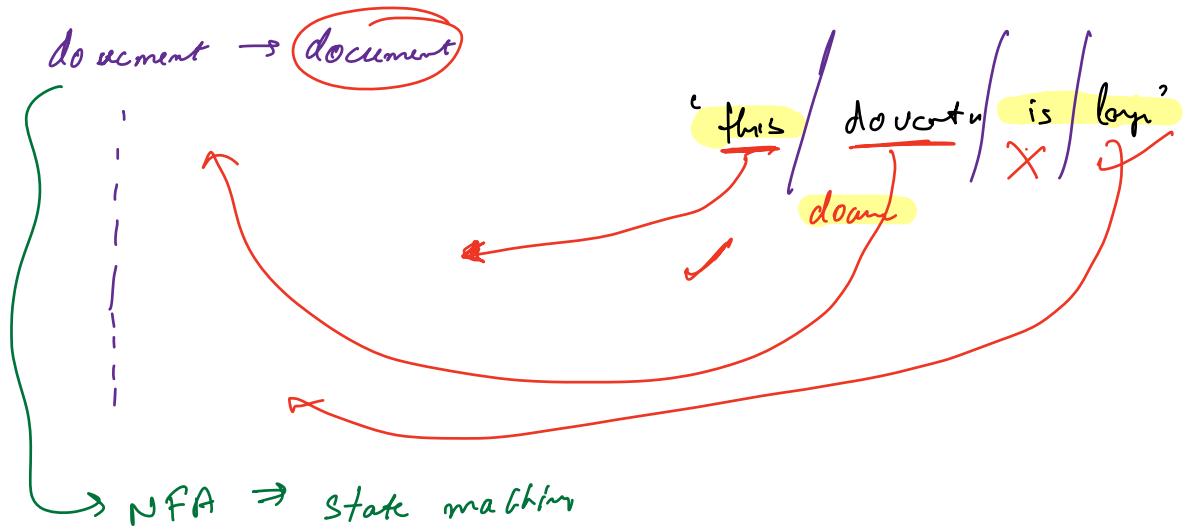
Stemming → quick & dirty & dumb process

Stemming → Stem ✓
Caring → car X] Case in which dumb process
lion → l X Stemming fails

Lemmatization → Smart sister of stemming.



10:41 → 10:43

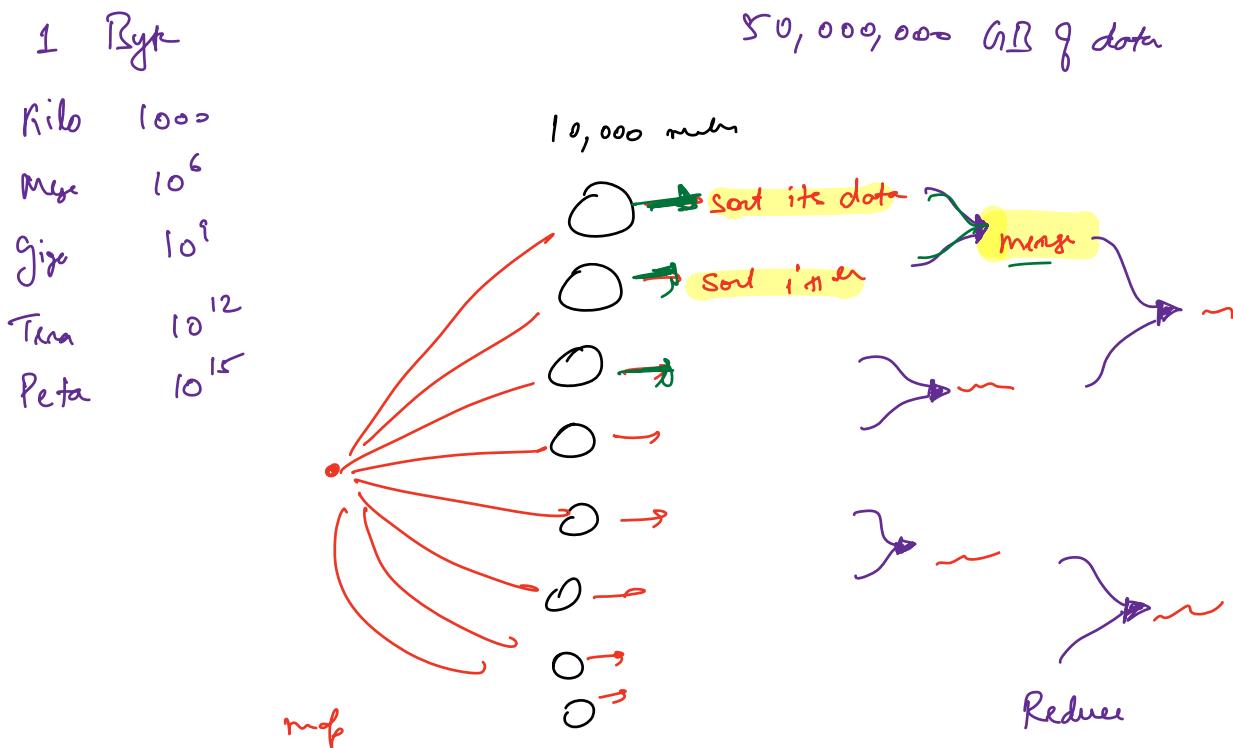


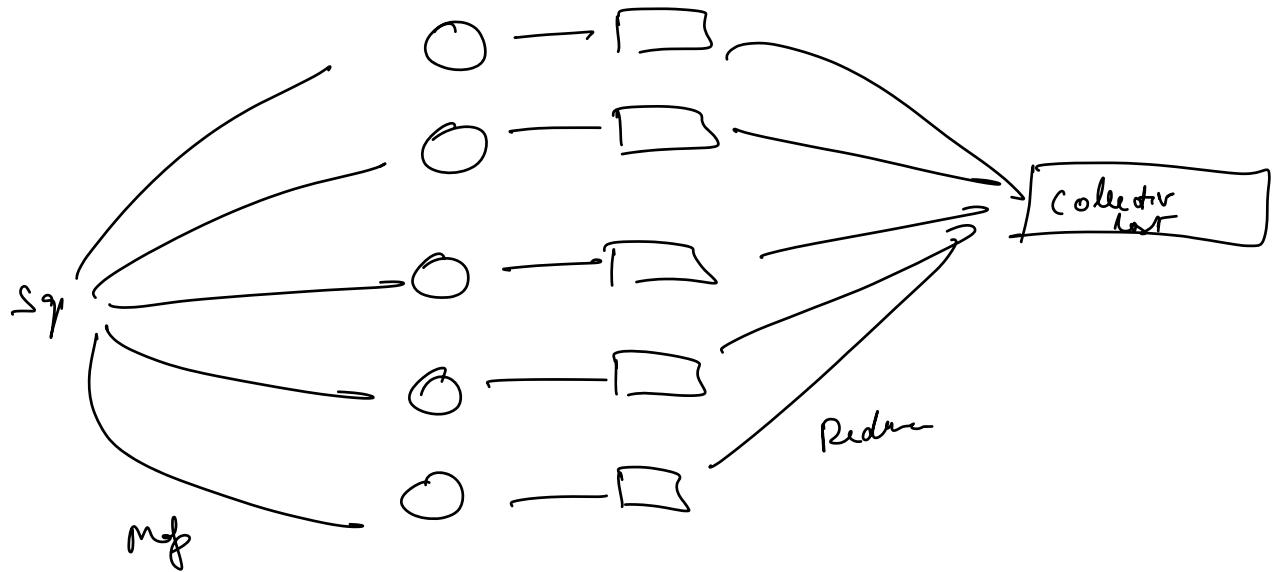
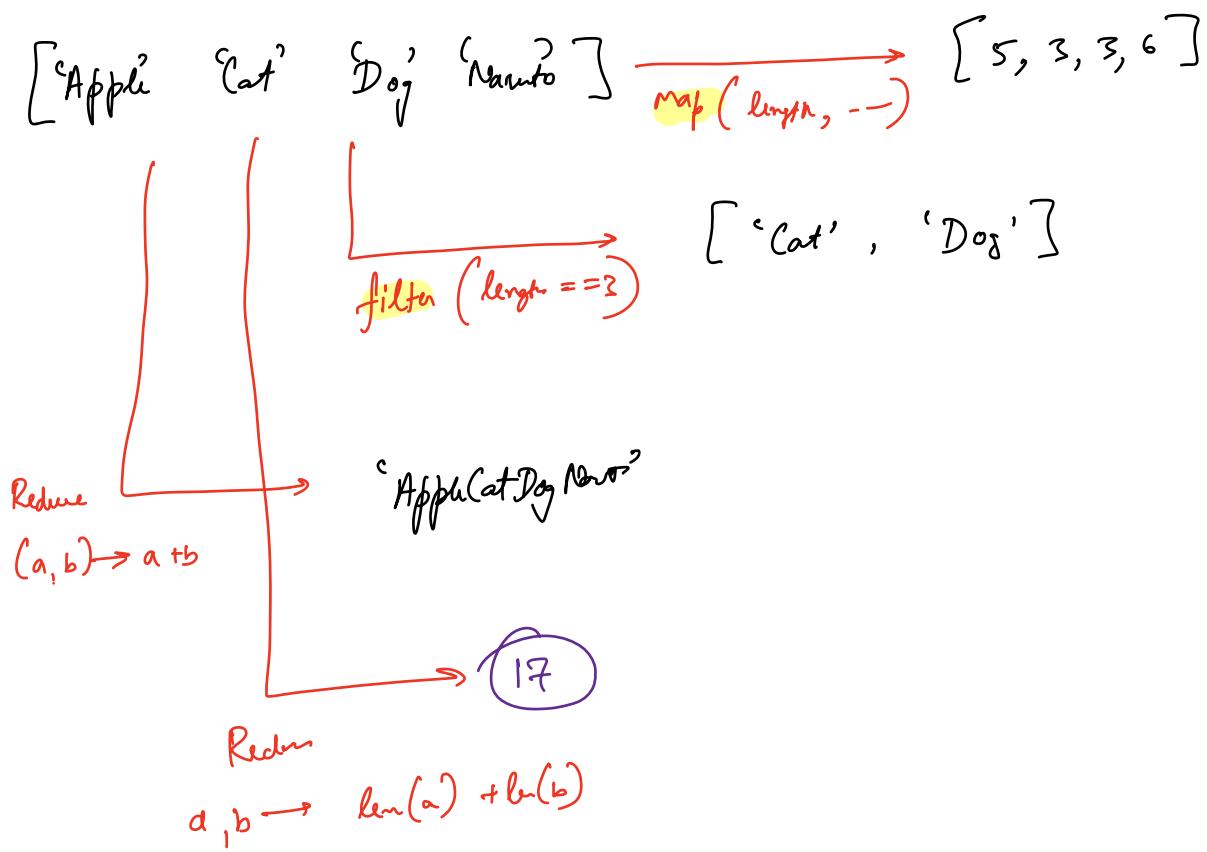
Non-deterministic Finite Automata

$\text{map}(\text{fn}, \text{items}) \rightarrow \text{apply this fn on each}$
 $\text{item } \in \text{items}$

$\text{reduce}(\text{items}, \text{acc}, \text{fn})$
 $\xrightarrow{\text{start value}}$ $\xrightarrow{\text{reducer}}$
 $\xrightarrow{\quad\quad\quad}$ $(a, b) \rightarrow \text{single value}$

given a list of text lines \rightarrow Sort in lexicographic order
 $\xrightarrow{\quad\quad\quad}$ 50 PB of data





LR + Queue

seen-id



d1 = "the quick brown fox jumps over the lazy dog"

d2 = "every strand has its own invited index"

d3 = "

doc { id
field 1 : _____
for 2 : _____
}

word : [(id, field-id, post), - - - ()]

Immediata consistency → ES.

