

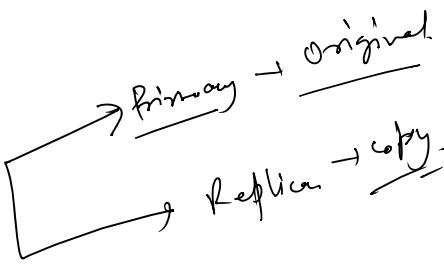
- ① Basis App/Arch → tool.
- ② Market Demand/Popularity.
- ③ Community Support.
- ④ License Cost.
- ⑤ Integration with App.
- ⑥ User friendly.
- ⑦ Scalability.
- ⑧ Customer Support
- ⑨ Analysis Engine.

ElasticSearch:-

index → logical division.  
 ↓  
 Data → json → (Key-value)  
 {  
 name: Vivek Arora.  
 } → Key      } → value.

Components of Index:-

- ① Document → json record.
- ② Fields → key-value.
- ③ Mapping → Analyzed & stored.
- ④ Shards → logical division of index for parallel processing.



Common Index operations:-

- (1) List all indices → GET \_cat?\_indices ?
- (2) Create an index → PUT / my-index
- (3) Delete an index → DELETE / my-index
- (4) View index setting → GET / my-index/\_settings
- (5) Add document → POST / my-index/\_doc/1 id=1
- (6) Search Document → GET / my-index/\_search ?q=user:john
- (7) View index mapping → GET / my-index/\_mapping

Create index :-

PUT | Product

Add document :-

POST | Product/\_doc/1

```
{
  "name": "dell",
  "price": 55000,
  "place": "B'k"
}
```

Search Document :- GET | Product/\_search

```
{
  "query": {
    "name": "dell"
  }
}
```

Forward Index.

Index

Inverted Index → Elastic Search.

Document ID

1

2

Content

"ElasticSearch is fast"  
"ElasticSearch scales horizontally"

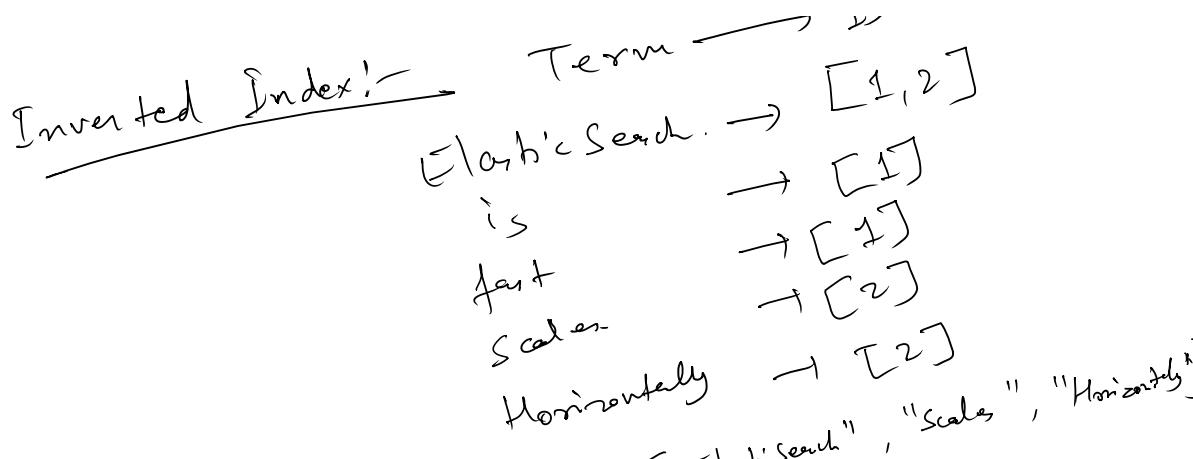
Forward Index :-

ID → Term.

1 → ElasticSearch is fast  
2 → ElasticSearch scales horizontally

Invert Index :-

Term → ID  
, [1, 2]



- How elasticSearch Builds:-
- ① Tokenization - Text is split into tokens (Word)
  - ② Normalization - Token into lowercase, stop word removed, stemming is applied.
  - ③ Inverted Index Creation - Map term with ID.

- "ElasticSearch is fast"
- ① Tokenization → ["ElasticSearch", "is", "fast"]
- ② Lowercasing → ["elasticsearch", "is", "fast"]
- ③ Stop-word Removal → [is, a, an, the, are, were, was]
- ["elasticsearch", "fast"]
- ④ Stemming: Reduce the word to the root form.
- ex) fastest → fast  
 run → run  
 running → run
- ["elasticsearch", "fast"]

- Benefits:-
- ① Reduce the index size → reduce noise.
  - ② Improve recalls → running → plan → route
  - ③ Improve precision → ignore irrelevant filler word.

- (1) Improve Precision → Ignore irrelevant terms
- (2) Improve Precision → Ignore irrelevant terms

Filler word [a, an, the, are, was, were, is, am, from, at, as]

- ① TF - Term frequency,
- ② IDF - Inverse Document frequency.

① Term frequency (TF) :- How often the term occurs in the document.

$$TF(t, d) = \frac{\text{Number of times term } t \text{ occurs in document } d}{\text{Total no. of terms in document } d}$$

② IDF (Inverted document frequency) :-

$$IDF = \frac{\text{Total no. of documents}}{\text{No. of documents that contain the term } t}$$

$$\text{Combined Score} = TF \times IDF$$

"ElasticSearch is fast"

$D_1$  = "Elasticsearch scales horizontally"

$D_2$  = "Fast System Scales"

$D_3$  = "Fast System Scales"

$TF(D_1)$

$1$

$1$

ElasticSearch

$\approx$

Fast

$TF(D_2)$

$1$

$0$

$\underline{\quad}$

$TF(D_3)$

$0$

$1$

$\underline{\quad}$

$(IDF)$

$3/2 - 1.5$

$3/2 - 1.5$

What is benefit?

①

TF - IDF

"Signal to the noise ratio"

TF - IDF → Combined Score  
↓

Strong & clear picture

Elastic Search Scaling

Components

- ① Cluster - collection of nodes is called cluster.
- ② Node - single instance where elastic search is run.
- ③ Index - Database
- ④ Shards - A physical piece of index.
- ⑤ Replica - copy the data from primary index.

Ex -

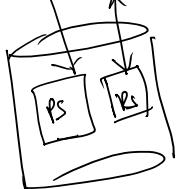
Index = VK-movies.

Primary shards

replica = 1

write feed.

Parallel write & read concept.



N ↴

## How Query Scales?

①

Coordinating node that will receive query.



② Query sent to all the relevant shard.

③ Execute query - fetch result.

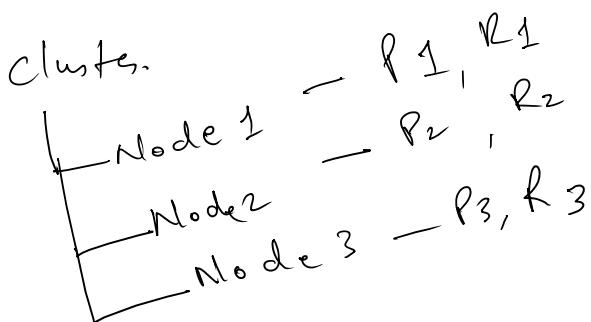
④ Combine all → give result.

⑤

## Cluster Benefits

⑥ Query speed.

- ① Availability
- ② Handle huge data.
- ③ Performance



Index - logical division of data.

Physical division of data into smaller

Shards - unit.

↳ fully functional Lucene index.

↳ Search library that powers the search engine like elasticsearch, solr

Why?

etc

## ④ Performance :-

① Scalability

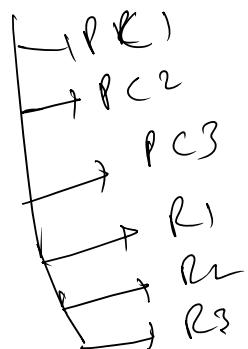
② Availability

③ High Availability

Ex:- PUT | Customer - data.

{ "setting": {  
    "no. of shards": 3,  
    "no. of replica": 1  
}}

Index = Customer - data.



\* ILM (Index lifecycle management) :-

Automate the index creation - deletion  
on the basis of index age, size & no. of documents.

① Recovery, which date for which index to

② Define which date to retain

③ Adhere company policy.

④ free - up disk space.

⑤ free - up read/write

→ ⑥ Hot - Read .... data

- ④ ↑
- 4 Phases :-
- ① Hot - Read
  - ② Warm - Read
  - ③ Cold. - Rarely query data
  - ④ Delete - Deletion of data will happen.
- ⑤ Attach the template with index.

Practical:-

- ① Policy.
- ② Template.
- ③ Index

Analyzers:-

Process the text during indexing & searching time.

Components:-

- ① character filters.
- ② Tokenizer.
- ③ Token filter.
- ④ language
- ⑤ stop
- ⑥ Keyword

Types of Built-in Analyzers:-

- ① Standard.
- ② Simple.
- ③ whitespace

```
POST my-index/_analyze
{
  "analyzer": "English",
  "text": "electric search scales horizontally"
}
```