# Elasticsearch

Jakub Podeszwik / infoShare ACADEMY / Yameo

# Agenda

- Full-text search
- Apache Lucene
- Elasticsearch / Elastic stack

# Full-text search

- Full-text database contains complete text of books, articles, e.t.c. They are stored as documents.
- examine all words in every stored document and return only documents that match
- serial scanning vs indexing

# Inverted index

1. Tom has a cat

2. Kate has a dog

3. Mike has an owl

| Term | Documents |
|------|-----------|
| a    | 1, 2      |
| an   | 3         |
| cat  | 1         |
| dog  | 2         |
| has  | 1, 2, 3   |
| kate | 2         |
| mike | 3         |
| owl  | 3         |
| tom  | 1         |

# Inverted index with fields

| | |
|---|---|
| 1. | |
| title: The cat | |
| text: The cat sits | |

| | |
|---|---|
| 2. | |
| title: The dog | |
| text: The dog stands | |

| Term | Documents |
|---|---|
| text:cat | 1 |
| text:dog | 2 |
| text:sits | 1 |
| text:stands | 2 |
| text:the | 1, 2 |
| title:cat | 1 |
| title:dog | 2 |
| title:the | 1,2 |

# Apache lucene

- high performance full-text search library written in Java
- open source (Apache license)
- initially released in 1999

# Term query

| Term | Documents |
|------|-----------|
| a | 1, 2 |
| cat | 1, 3 |
| dog | 2 |
| has | 1, 2 |
| kate | 2 |
| sits | 3 |
| the | 3 |
| tom | 1 |

{1, 3}

term: cat

1. Tom has a cat

2. Kate has a dog

3. The cat sits

# Boolean query

bool.should = or = sum

{1, 2, 3} ← {1, 3}
           {2}

bool.should:
  term: cat
  term: dog

| Term | Documents |
|------|-----------|
| a    | 1, 2      |
| cat  | 1, 3      |
| dog  | 2         |
| has  | 1, 2      |
| kate | 2         |
| sits | 3         |
| the  | 3         |
| tom  | 1         |

1. Tom has a cat

2. Kate has a dog

3. The cat sits

https://github.com/jpodeszwik/elasticsearch-workshop-01-2019

# Inverted index

1. Tom has a cat

2. Kate has a dog

3. Mike has an owl

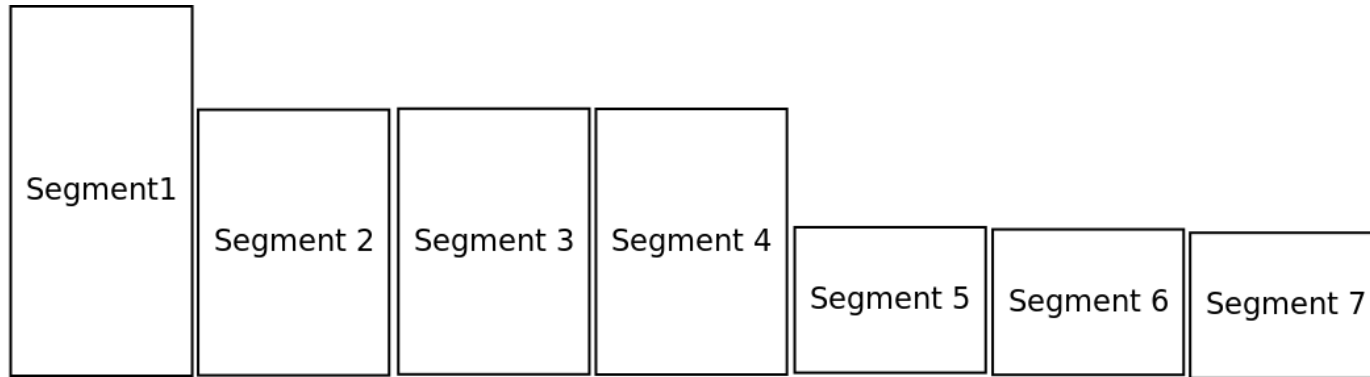| Term | Documents |
|------|-----------|
| a | 1, 2 |
| an | 3 |
| cat | 1 |
| dog | 2 |
| has | 1, 2, 3 |
| kate | 2 |
| mike | 3 |
| owl | 3 |
| tom | 1 |

# Lucene segment

immutable, stored on disk data structure constisting of:

- inverted index
- fielddata cache / doc_values
- _source
- live documents bitset (mutable)

# Lucene segments

Segment1

Segment 2

Segment 3

Segment 4

Segment 5

Segment 6

Segment 7

# Segment merging

# Elasticsearch

- full text search engine built on top of lucene
- log analytics
- horizontally scalable writes
- open source (Apache license)
- v 0.4 released in 2010
- v 1.0 released in 2014
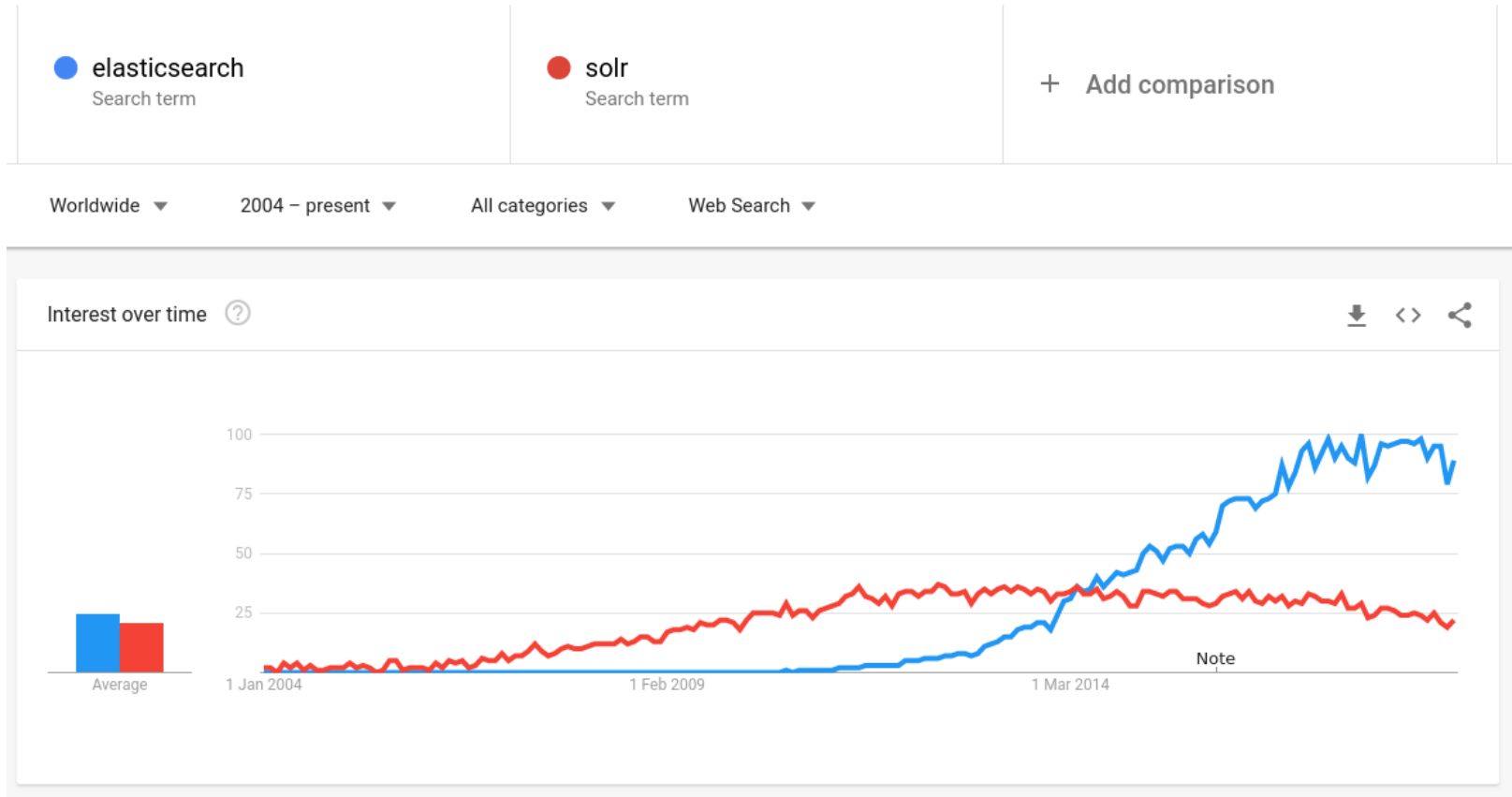- RESTful api
- resilence
- some features are paid

# Apache Solr

- in 2010 Elasticsearch's competetive project apache Solr joined lucene as a subproject
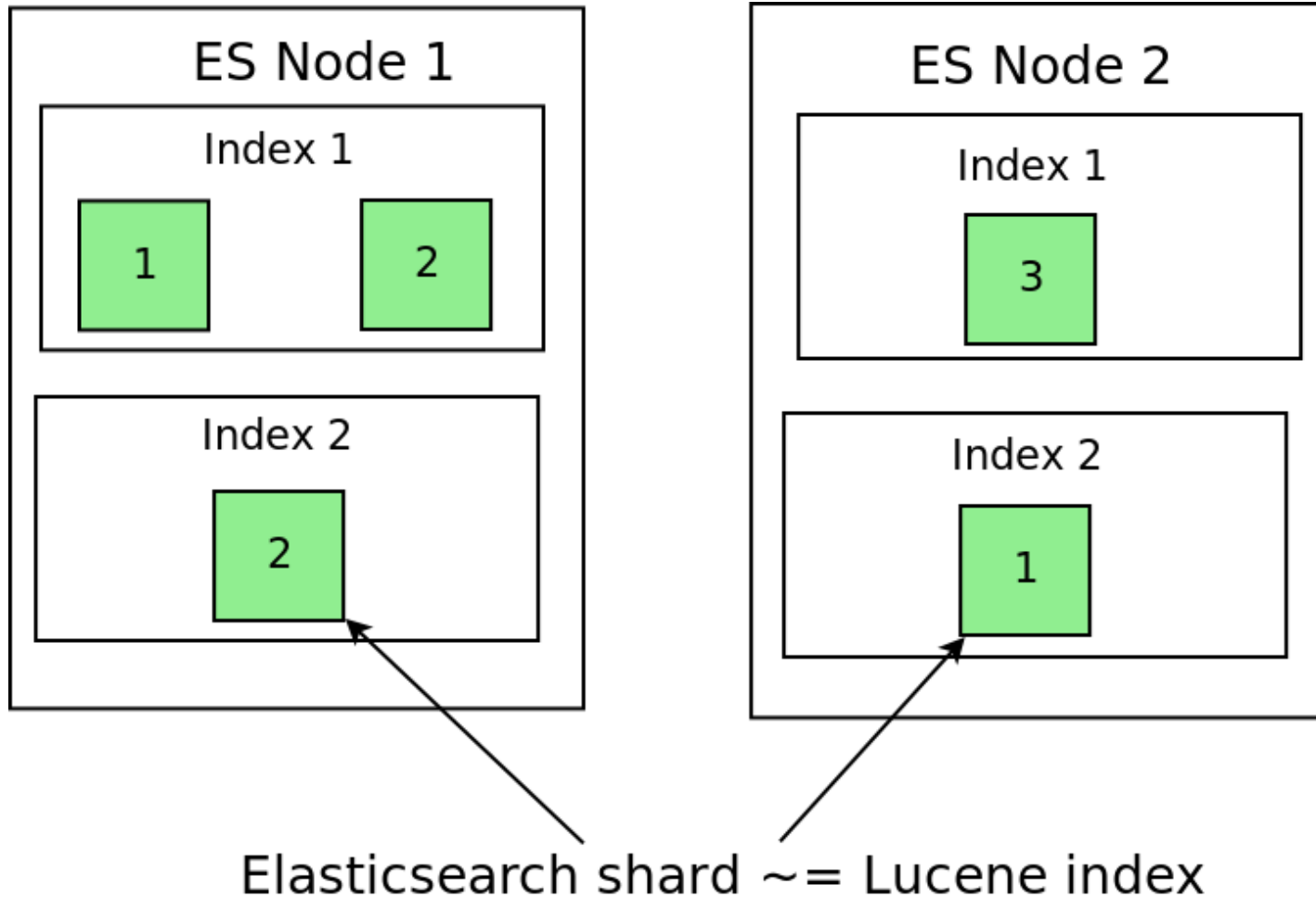- features are similar
- elasticsearch is more popular

18 systems in ranking, January 2019

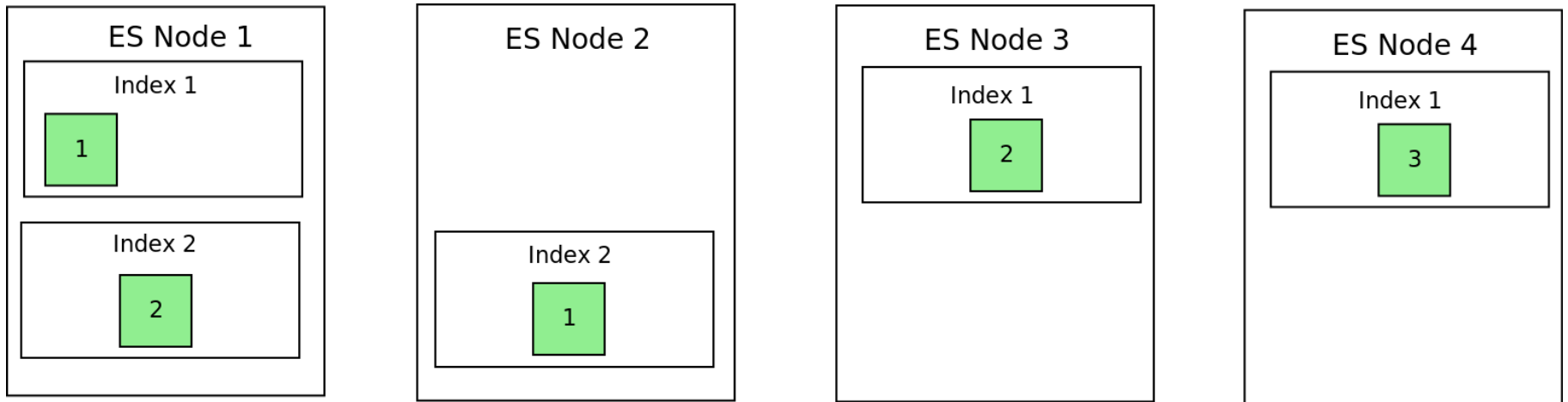| Rank | | | DBMS | Database Model | Score | | |
|---|---|---|---|---|---|---|---|
| Jan 2019 | Dec 2018 | Jan 2018 | | | Jan 2019 | Dec 2018 | Jan 2018 |
| 1. | 1. | 1. | Elasticsearch ➕ | Search engine | 143.44 | -1.26 | +20.89 |
| 2. | 2. | ⬆3. | Splunk | Search engine | 81.43 | -0.76 | +17.42 |
| 3. | 3. | ⬇2. | Solr | Search engine | 61.48 | +0.13 | -2.89 |
| 4. | 4. | 4. | MarkLogic ➕ | Multi-model 🔢 | 14.26 | -0.02 | +3.05 |

# Apache Solr
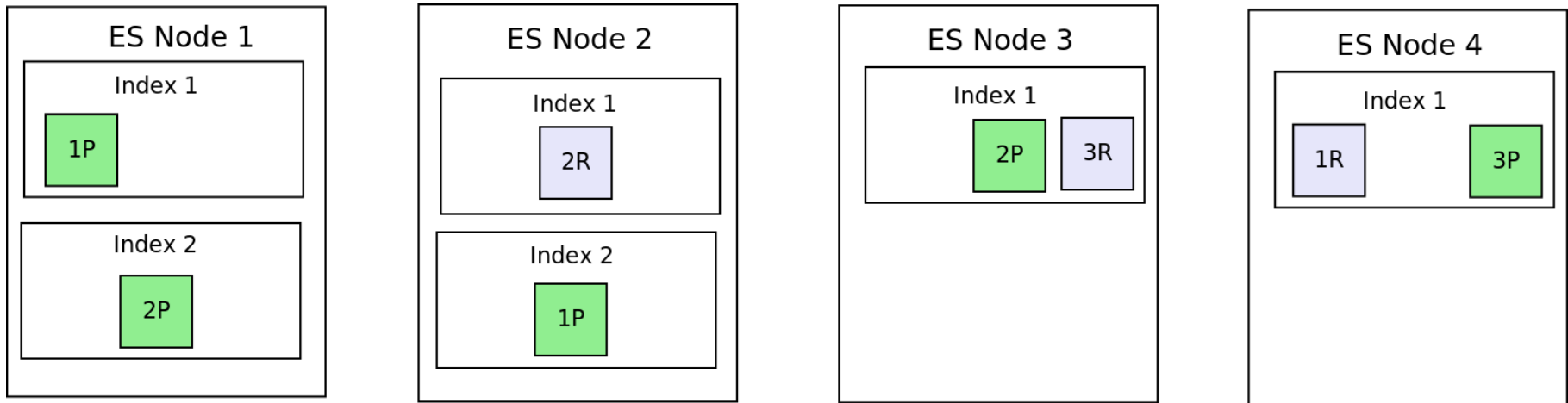
# Elasticsearch index



Elasticsearch shard ~= Lucene index

# Elasticsearch index

# Elasticsearch index

# Kibana

- dev tools
- admin interface
- visualizations
- clickable dashboards
- managing elasticsearch

# Kibana

# Kibana

# Logstash

- Log forwarder
- many inputs (tcp, file, kafka)
- filters for modifying / enriching the data
- outputs to elasticsearch and some other

# Logstash



Analysis

Archiving

Monitoring

Alerting

# Query types

- term
- bool
- prefix
- wildcard
- fuzzy
- match, match_all
- query_string

# Term query

Term | Documents
--- | ---
a | 1, 2
cat | 1, 3
dog | 2
has | 1, 2
kate | 2
sits | 3
the | 3
tom | 1

{1, 3}

term: cat

1. Tom has a cat

2. Kate has a dog

3. The cat sits

# Bool query

bool.should = or = sum

{1, 2, 3} ← {1, 3}
　　　　　　{2}

bool.should:
　term: cat
　term: dog

| Term | Documents |
|------|-----------|
| a | 1, 2 |
| cat | 1, 3 |
| dog | 2 |
| has | 1, 2 |
| kate | 2 |
| sits | 3 |
| the | 3 |
| tom | 1 |

1. Tom has a cat

2. Kate has a dog

3. The cat sits

# Wildcard query

| Term | Documents |
|------|-----------|
| a | 1, 2 |
| cat | 1, 3 |
| dog | 2 |
| has | 1, 2 |
| kate | 2 |
| sits | 3 |
| the | 3 |
| tom | 1 |

{1, 3} ← {3}
{1}

wildcard: t*

1. Tom has a cat

2. Kate has a dog

3. The cat sits

# Wildcard query

| Term | Documents |
|------|-----------|
| a | 1, 2 |
| cat | 1, 3 |
| dog | 2 |
| has | 1, 2 |
| kate | 2 |
| sits | 3 |
| the | 3 |
| tom | 1 |

{2, 3} ← {3}
{2}

wildcard: *e

1. Tom has a cat

2. Kate has a dog

3. The cat sits

# Fuzzy query

- Levenshtein distance
- fuzziness value: 1, 2, AUTO

# Match query

sum (should)

{1, 2, 3} ← {1, 3}
        {2}

| Term | Documents |
|------|-----------|
| a | 1, 2 |
| cat | 1, 3 |
| dog | 2 |
| has | 1, 2 |
| kate | 2 |
| sits | 3 |
| the | 3 |
| tom | 1 |

analysis

match: cat dog → bool.should:
term: cat
term: dog

1. Tom has a cat

2. Kate has a dog

3. The cat sits

# Analyzers

- whitespace
- standard - default analyzer for text fields
- keyword - noop

# Custom analyzers

- char_filters - removes / changes / adds characters before text goes to tokenizer
- tokenizer - splits text into tokens
- filters - removes / changes / adds tokens

# Custom analyzer

Tom has a cat, and Kate has two dogs.

| Strip puntuation char filter |
|:---:|

Tom has a cat and Kate has two dogs

| Whitespace tokenizer |
|:---:|

[Tom, has, a, cat, and, Kate, has, two, dogs]

| Lowercase filter |
|:---:|

[tom, has, a, cat, and, kate, has, two, dogs]

| English stopwords filter |
|:---:|

[tom, cat, kate, two, dogs]

| English stem filter |
|:---:|

[tom, cat, kate, two, dog]

# Scoring

$$score(q, d) =$$

$$queryNorm(q) * coord(q, d)*$$

$$\sum (tf(t \ in \ d) * idf(t)^2 * t.getBoost() * norm(t, d))(t \ in \ q)$$

# Scoring

- score(q,d) - the relevance score of document d for query q
- queryNorm(q) - query normalization factor
- coord(q,d) is the coordination factor (rewards for higher percentage of query terms contained in document)
- tf(t in d) - term frequency for term t in document d,
- idf(t) - inverse document frequency for term
- t.getBoost() - boost that has been applied to the query,
- norm(t,d) - field-length norm, combined with the index-time field-level boost

# Aggregations

- aggregate results into buckets
- count metrics for each of buckets
- can be nested
- fielddata vs doc_values

# Kibana dashboards

# Nested documents

- can be used to model one to many relations
- each nested document is treated as separate document

# Routing

- possibility to split documents into shards basing on a value
- documents with the same routing value will be stored in the same shard
- documents with different routing value might or might not be stored in different shards
- in low cardinality fields can lead to putting all documents into single shard.

# Parent child

- another way to store one to many relations
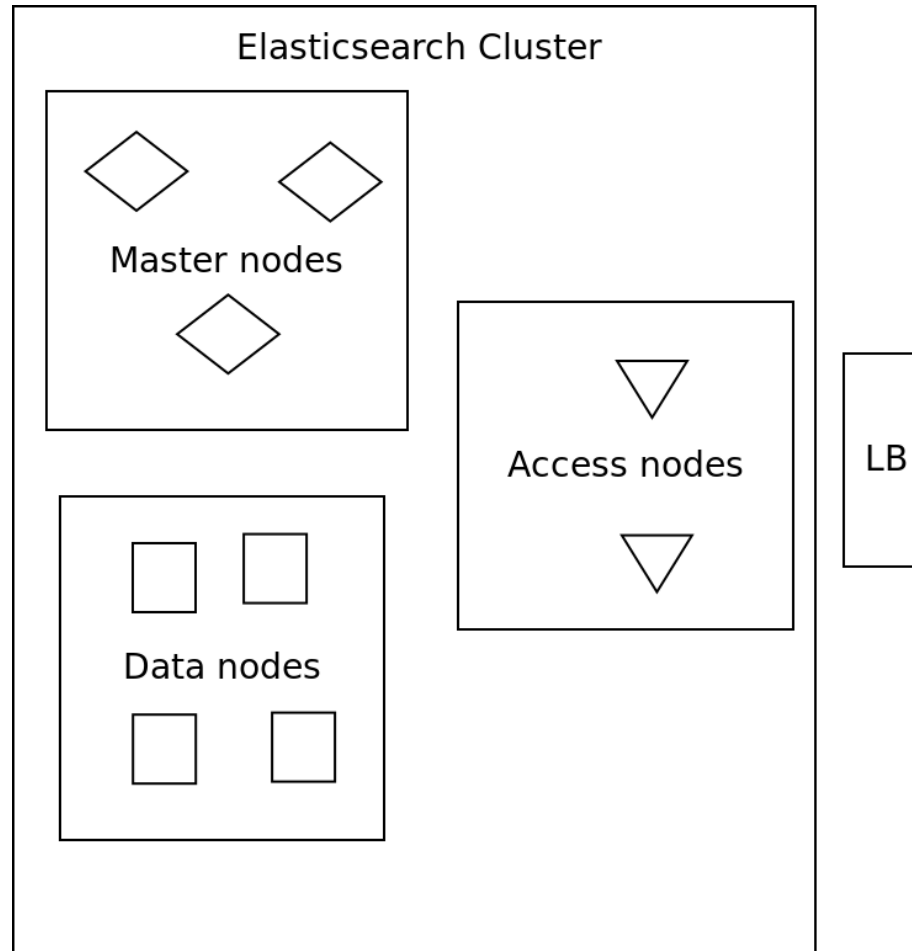- search performance is significantly lower

# Cluster

- Master node - keeping, updating, broadcasting state of the cluster (list of nodes, mappings)
- Data node - storing data and executing searches on shard
- Access node - forwarding requests to data nodes and merging results
- Ingest node - preprocess documents before indexing them

# Leader election

- zen discovery
- quorum

# Cluster

# _cat api

- indices
- shards
- nodes
- master
- ?v - legend

# Premium features

- security (alternative search guard)
- machine learning
- alerting
- graph exploration api