# **Operating Elasticsearch at Scale!**

Vikram Ramakrishnan & Aishwarya Sankaravadivel June 2019 @ SRECON Asia/Pacific

#### Agenda

- Introductions
- Search Platform @ PayPal
- Elasticsearch Primer
- Making Elasticsearch truly Elastic
- Managing Elasticsearch @ Scale
- Elasticsearch in action!
- Q&A

#### Introductions

#### Vikram Ramakrishnan



Senior Manager, Technology Platforms & Experience

15 years of industry experience with specific focus in the areas of Large scale distributed systems, Enterprise Search & Big Data

vikramakrishnan@paypal.comhttps://www.linkedin.com/in/vikramakrishnan

#### Aishwarya Sankaravadivel



Senior Software Engineer Technology Platforms & Experience

Passionate technologist with deep expertise in managing Large scale Elasticsearch deployments & Reactive systems

asankaravadivel@paypal.comhttps://www.linkedin.com/in/aishwaryasankar/

#### Search Platform @ PayPal



Fully Managed Search offering for PayPal Inc.











20+ Billion Ingests / day

50+ Million Searches / day





#### What's in it for SREs?



Improved Observability



#### Realtime, Actionable Insights

Things to know before we dig deep...

- Cluster
- Node
- Document
- Index

- Field
- Mapping
- Shards
- Routing

Anatomy of an Elasticsearch cluster



Primary Shard Replica Shard

Anatomy of an Indexing request



#### Dissecting Indexing



#### doc 2

#### Inverted Index

quick	1
brown	1
fox	1
jumped	1
over	1
lazy	1, 2
dog	1, 2
two	2
dogs	2
slower	2
less	2
rover	2

Anatomy of a search request.



Dissecting a search request.



#### Inverted Index

Installation

\$ curl -L -O https://artifacts.elastic.co/downloads/elasticsearch/elasticsearch-{VERSION}.tar.gz

\$ tar -xvf elasticsearch-{VERSION}.tar.gz

\$ cd elasticsearch-{VERSION}/bin

\$ ./elasticsearch

Installation

```
[~]$ curl http://localhost:9200
  "name" : "LM-MAA-26500559_bloodseeker-masterdata",
  "cluster_name" : "bloodseeker",
  "cluster_uuid" : "gy1tSWy5RNCkAoIAxprpWA",
  "version" : {
   "number" : "6.2.4",
    "build_hash" : "ccec39f",
    "build_date" : "2018-04-12T20:37:28.497551Z",
    "build_snapshot" : false,
    "lucene_version" : "7.2.1",
    "minimum_wire_compatibility_version" : "5.6.0",
    "minimum_index_compatibility_version" : "5.0.0"
 },
  "tagline" : "You Know, for Search"
```

```
Now, how do we
```

```
scale this for the
```

```
enterprise?
```

# SCALABILITY?

# Scale Up Scale-up Scale-up Scale Up





Memory / CPU / IO / Network Issues

> Sub-optimal Queries

> Mapping Explosion

Undersized / Oversized shards



# Is Elasticsearch truly elastic?



Problem#1 : Low Ingestion Throughput

Need to process 'X' million / minute additionally.

#### Learnings

- ➢ \_routing for Bulk requests
- index.refresh\_interval
- > number\_of\_shards, total\_shards\_per\_node
- index.translog.durability:async

Problem #2 : Node instabilities due to Mapping Explosion and Sparse Fields

Problem #2 : Node instabilities due to Mapping Explosion and Sparse Fields

Mapping Explosion

	C 1	C 2	C 3	C 4	C 5	•	0	•	•	0	•	0	•	•	0	0	0	0	C 2000
R1																			
R2																			
R3																			

#### Problem #2 : Node instabilities due to Mapping Explosion and Sparse Fields

- Mapping Explosion
- Sparse fields

	C 1	C 2	C 3	C 4	C 5	C 6	С 7	C 8	C 9	C 10	C 11	C 12	C 13	C 14	C 15	C 16	C 17	C 18
R1	A 1			A 2		A 3	A 4			A 5		A 8	A 9		A 10			A 11
R2		A 1	2		A 2		-				A 3		A 4			A 5		
R3			A 1			A 2		A 3			A 4	A 5	A 6	A 7		A 9		

#### Problem #2 : Node instabilities due to Mapping Explosion and Sparse Fields

Learnings:

- > Keep a watch on mapping.
- > Split indices.
- > Date based indices works best for logs monitoring!



Example Queries:

#1 Group By operation on a ld field

```
"key": "a@xxx.com",
    "doc_count": 1
}, {
    "key": "b@yyy.com",
    "doc_count": 1
}, {
    "key": "c@zzz.com",
    "doc_count": 1
}......]
```

Example Queries:

#2 Retrieve all documents from Elasticsearch

```
_search?scroll=1m
"query": {
    "match_all": {}
}
```

```
"scroll": "1m",
"scroll_id" :
"cXVlcnlUaGVuRmV0Y2T="
```

Example Queries:

#3 Multiple levels of aggregations

Example Queries:

#4 Wildcard / Regex Queries



Learnings : Query Logger

Logging in Elasticsearch

Is index.search.slowlog.\* suffice?

Distinct Query Logger





#### Problem #3 : Expensive/sub optimal queries – A threat to availability!



Example: Time range > 48 hours | Nested Aggregations level is >10

#### Learnings : Query Optimizer



Learnings : Query Optimizer

```
{"index": ["INDEXNAME-*"]}
{"query":
    {"range":
        {"timestamp": {
            "gte": <timestamp>,
            "lte": <timestamp>}
        }..
    .}...
}
```

#### Problem #3 : Expensive/sub optimal queries – A threat to availability!



#### Problem #3 : Expensive/sub optimal queries – A threat to availability!



Learnings : Query Optimizer

- > Reduce the number of indices(shards)that are being queried.
- > Reorder the queries for reducing the scan margin \*



#### Problem #4: Oversized Shards

When?

- > Shards assigned per index is incorrect.
- Skewness in the values of the routing key.

#### Learnings

Anticipate the future and assign the shards.

> Split

> Re-index



- > Monitoring your clusters
- Securing Elasticsearch
- Upgrading Clusters



#### Monitoring your clusters!



"If you can't measure it, you can't manage it"

Peter Drucker



#### Collect Node stats from every node

Index	field data, docs, merge, refresh
> Shard	type, state, docs
Thread Pool	type, queue, rejected, max
MVL <	heap_used, gc collection count and time, buffer pools
Os & Process	CPU Load, Memory, File descriptors











#### Securing Elasticsearch!

#### Why do we need security?

In-House Guard Plugin

- End to End TLS/SSL
- Authentication & Authorization

Security features, free in Elasticsearch?

Audit Logging

Third party integrations with LDAP/SAML

It's tough only when you are in <5.X versions!

Must do checks:

- > Are queries backward compatible?
- Mapping changes between versions?

Approach 1: Parallel cluster with dual ingestion

Approach 2: Re-index from remote

# Demo

```
"hits": {
   "total": 11748840,
   "max score": 1,
 "hits": [
    T {
          " index": "demo",
          " type": "file",
          " id": "WwUcRWsBqsrfqUax3C-c",
            score": 1,
        v " source": {
              "correlationId": "FgPjzpga6r",
              "applicationName": "application-13962406",
              "api": "v1/tenant/demo/test-13962406"
       },
     Ψ.
           index": "demo",
          " type": "file",
           id": "XgUcRWsBqsrfqUax3C-c",
            score": 1,
        " source": {
              "correlationId": "711oefSR5U",
              "applicationName": "application-13962409",
              "api": "v1/tenant/demo/test-13962409"
       },
    W - {
          " index": "demo",
          " type": "file",
           id": "ZAUcRWsBqsrfqUax3C-c",
           score": 1,
        " source": {
              "correlationId": "mdOu4JwvgJ",
              "applicationName": "application-13962415",
              "api": "v1/tenant/demo/test-13962415"
```

# Q&A