



RUNNING & SCALING LARGE ELASTICSEARCH CLUSTERS

FRED DE VILLAMIL, DIRECTOR OF
INFRASTRUCTURE

@FDEVILLAMIL

OCTOBER 2017

BACKGROUND

- FRED DE VILLAMIL, 39 ANS, TEAM COFFEE @SYNTHESIO,
- LINUX / (FREE)BSD USER SINCE 1996,
- OPEN SOURCE CONTRIBUTOR SINCE 1998,
- LOVES TENNIS, PHOTOGRAPHY, CUTE OTTERS, INAPPROPRIATE HUMOR AND ELASTICSEARCH CLUSTERS OF UNUSUAL SIZE.

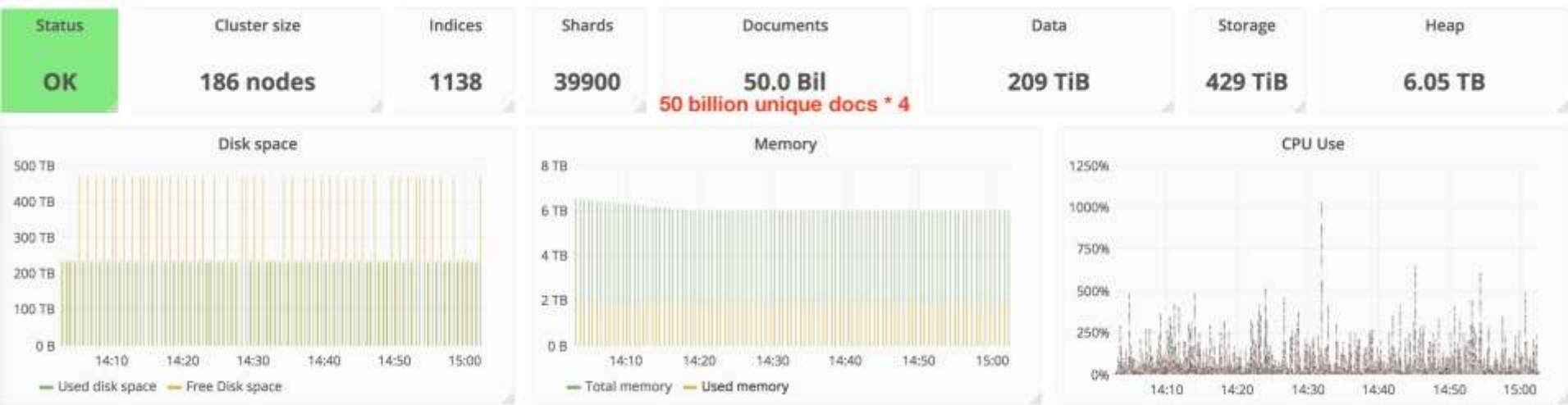
WRITES ABOUT ES AT
[HTTPS://THOUGHTS.T37.NET](https://thoughts.t37.net)



ELASTICSEARCH @SYNTHESIO

- 8 production clusters, 600 hosts, 1.7PB storage, 37.5TB RAM, average 15k writes / s, 800 search /s, some inputs > 200MB.
- Data nodes: 6 core Xeon E5v3, 64GB RAM, 4*800GB SSD RAID0. Sometimes bi Xeon E5-2687Wv4 12 core (160 watts!!!).
- We agregate data from various cold storage and make them searchable in a giffy.

AN ELASTICSEARCH CLUSTER OF UNUSUAL SIZE





ENSURING HIGH AVAILABILITY

NEVER GONNA GIVE YOU UP

- NEVER GONNA LET YOU DOWN,
- NEVER GONNA RUN AROUND AND DESERT YOU,
- NEVER GONNA MAKE YOU CRY,
- NEVER GONNA SAY GOODBYE,
- NEVER GONNA TELL A LIE & HURT YOU.



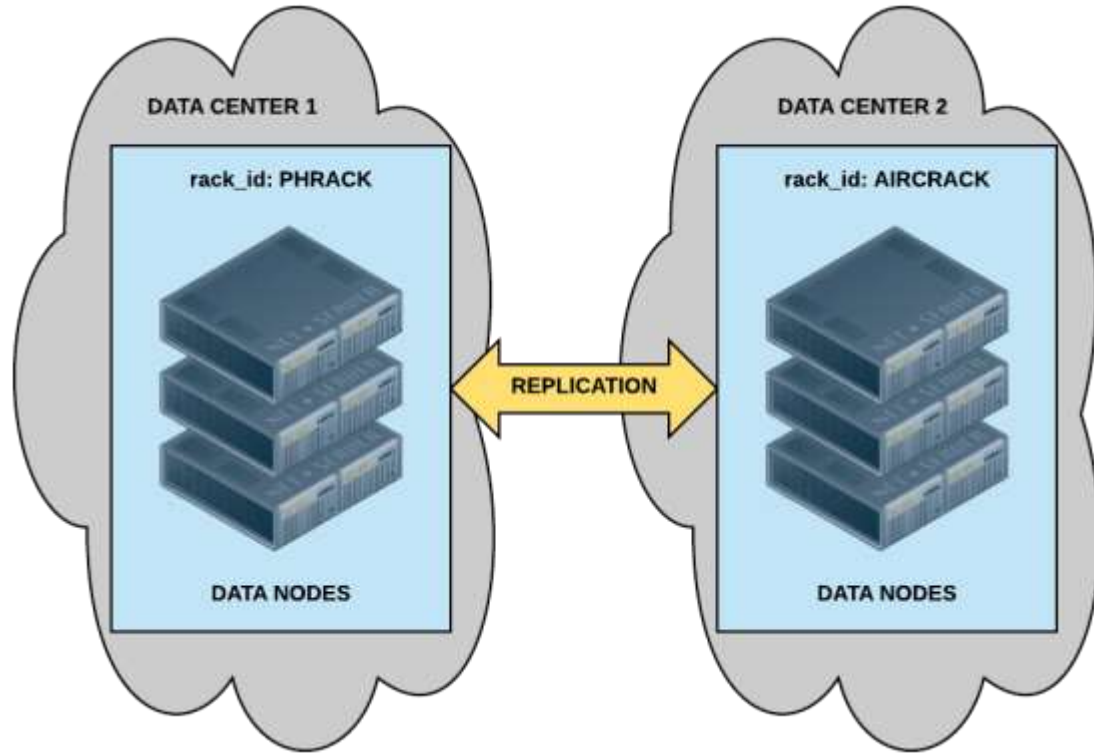
AVOIDING DOWNTIME & SPLIT BRAINS

- RUN AT LEAST 3 MASTER NODES INTO 3 DIFFERENT LOCATIONS.
- NEVER RUN BULK QUERIES ON THE MASTER NODES.
- ACTUALLY NEVER RUN ANYTHING BUT ADMINISTRATIVE TASKS ON THE MASTER NODES.
- SPREAD YOUR DATA INTO 2 DIFFERENT LOCATION WITH AT LEAST A REPLICATION FACTOR OF 1 (1 PRIMARY, 1 REPLICA).

RACK AWARENESS

ALLOCATE A
RACK_ID TO THE
DATA NODES FOR
EVEN REPLICATION.

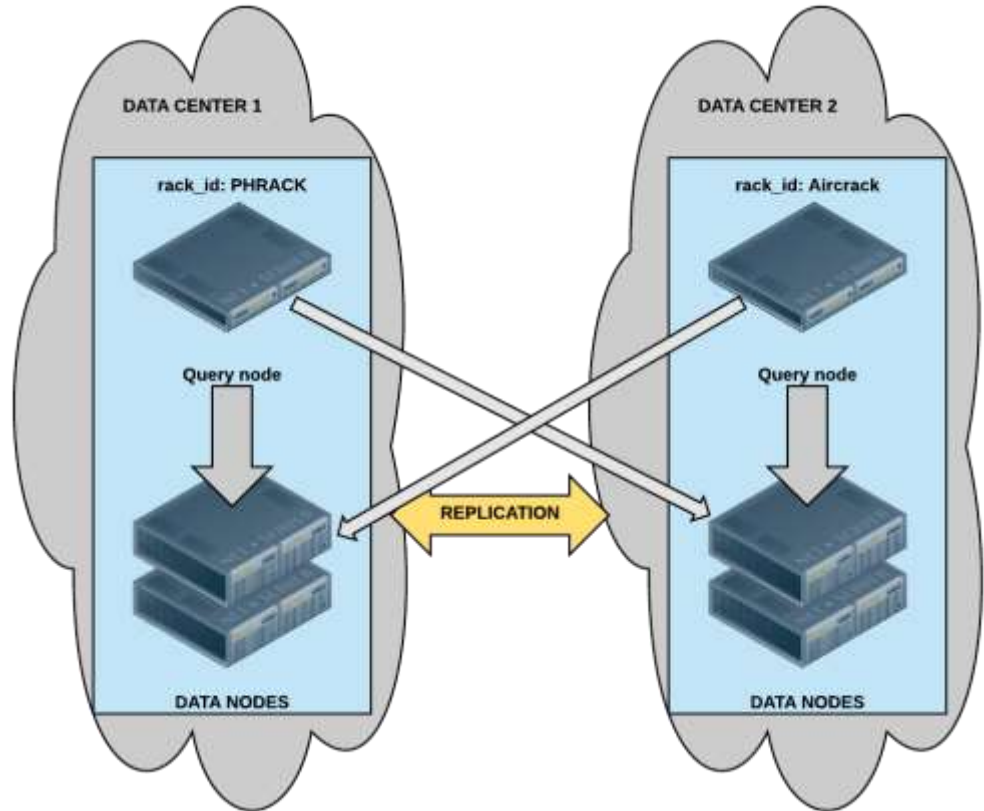
RESTART A WHOLE
DATA CENTER
@ONCE WITHOUT
DOWNTIME.



RACK AWARENESS + QUERY NODES == MAGIC

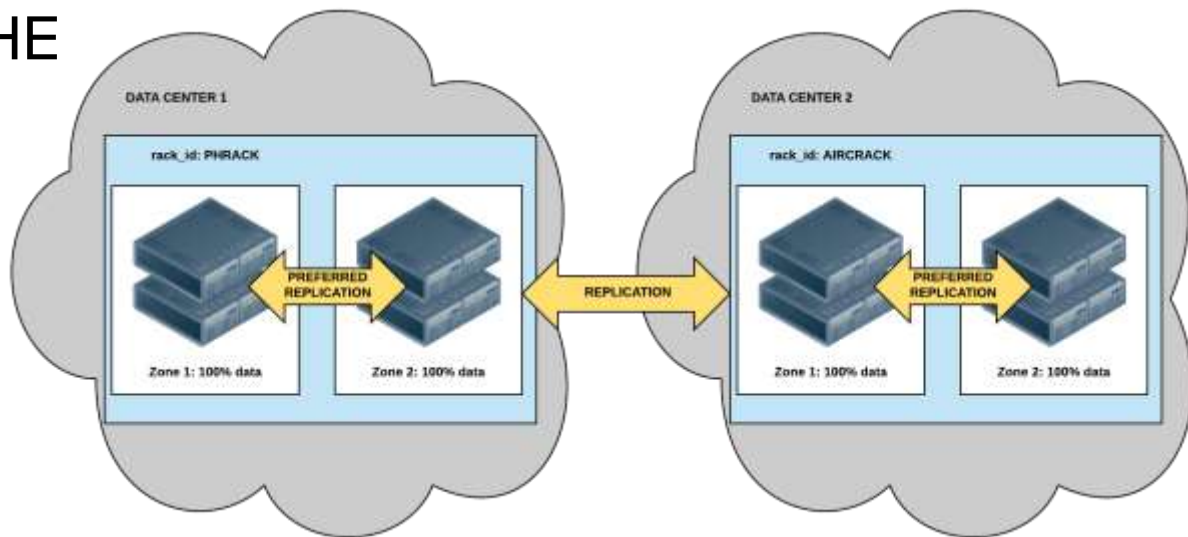
ES PRIVILEGES THE
DATA NODES WITH
THE SAME RACK_ID
AS THE QUERY.

REDUCES LATENCY
AND BALANCES THE
LOAD.



RACK AWARENESS + QUERY NODES + ZONE == MAGIC + FUN

ADD ZONES INTO THE
SAME RACK FOR
EVEN REPLICATION
WITH HIGHER
FACTOR.

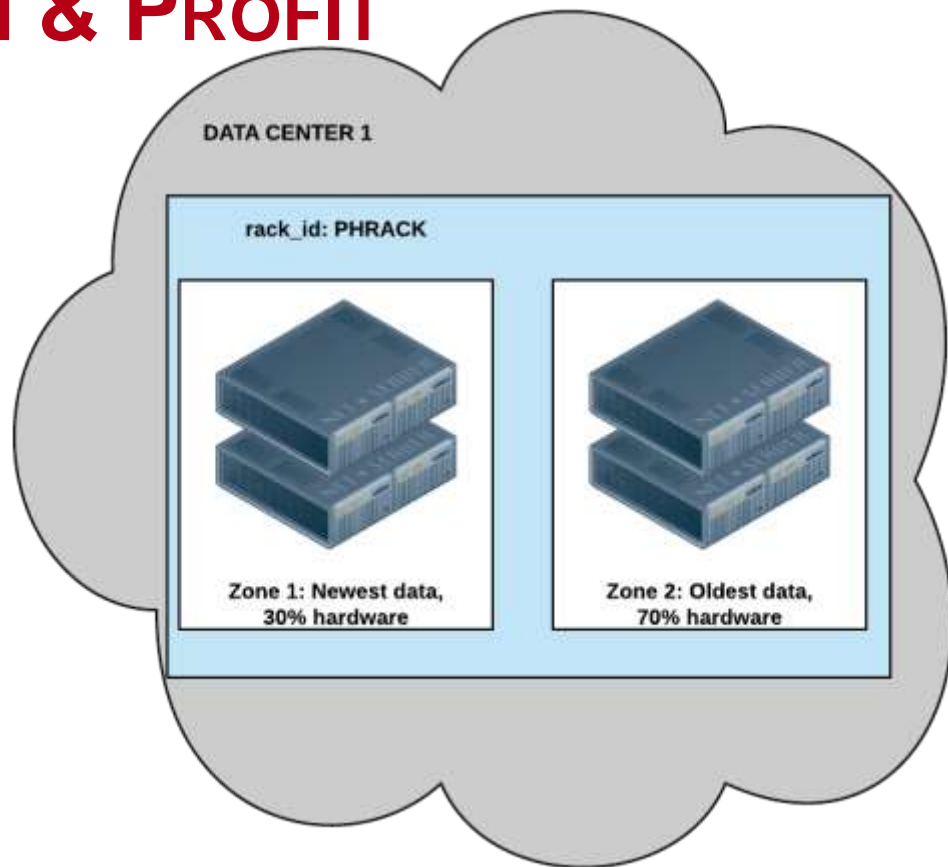


USING ZONES FOR FUN & PROFIT

ALLOWING EVEN REPLICATION
WITH A HIGHER FACTOR WITHIN
THE SAME RACK.

ALLOWING MORE RESOURCES
TO THE MOST FREQUENTLY
ACCESSED INDEXES.

...





AVOIDING MEMORY NIGHTMARE

HOW ELASTICSEARCH USES THE MEMORY

- Starts with allocating memory for Java heap.
- The Java heap contains all Elasticsearch buffers and caches + a few other things.
- Each Java thread maps a system thread: +128kB off heap.
- Elected master uses 250kB to store each shard information inside the cluster.

ALLOCATING MEMORY

- Never allocate more than 31GB heap to avoid the compressed pointers issue.
- Use 1/2 of your memory up to 31GB.
- Feed your master and query nodes, the more the better (including CPU).



MEMORY LOCK

- Use `memory_lock: true` at startup.
- Requires `ulimit -l unlimited`.
- Allocates the whole heap at once.
- Uses contiguous memory regions.
- Avoids swapping (you should disable swap anyway).



CHOSING THE RIGHT GARBAGE COLLECTOR

- ES runs with CMS as a default garbage collector.
- CMS was designed for heaps < 4GB.
- Stop the world garbage collection last too long & blocks the cluster.
- Solution: switching to G1GC (default in Java9, unsupported).



CMS VS G1GC

- CMS: SHARED CPU TIME WITH THE APPLICATION. “STOPS THE WORLD” WHEN TOO MANY MEMORY TO CLEAN UNTIL IT SENDS AN `OutOfMemoryError`.
- G1GC: SHORT, MORE FREQUENT, PAUSES. WON'T STOP A NODE UNTIL IT LEAVES THE CLUSTER.
- ELASTIC SAYS: **DON'T USE G1GC** FOR REASONS, SO READ THE DOC.

G1GC OPTIONS

+USEG1GC: ACTIVATES G1GC

MAXGC_PAUSEMILLIS: TARGET FOR MAX GC PAUSE TIME.

GC_PAUSEINTERVALMILLIS: TARGET FOR COLLECTION TIME SPACE

INITIATING_HEAP_OCCUPANCY_PERCENT: WHEN TO START COLLECTING?

CHOOSING THE RIGHT STORAGE

- `MMapFS` : MAPS LUCENE FILES ON THE VIRTUAL MEMORY USING `mmap`. NEEDS AS MUCH MEMORY AS THE FILE BEING MAPPED TO AVOID ISSUES.
- `NIOFS` : APPLIES A SHARED LOCK ON LUCENE FILES AND RELIES ON VFS CACHE.



BUFFERS AND CACHES

- ELASTICSEARCH HAS MULTIPLE CACHES & BUFFERS, WITH DEFAULT VALUES, KNOW THEM!
- BUFFERS + CACHE MUST BE < TOTAL JAVA HEAP (OBVIOUS BUT...).
- AUTOMATED EVICTION ON THE CACHE, BUT FORCING IT CAN SAVE YOUR LIFE WITH A SMALL OVERHEAD.
- IF YOU HAVE OOM ISSUES, DISABLE THE CACHES!
- FROM A USER POV, CIRCUIT BREAKERS ARE A NO GO!



MANAGING LARGE INDEXES

INDEX DESIGN

- VERSION YOUR INDEX BY MAPPING: 1_*, 2_* ETC.
- THE MORE SHARDS, THE BETTER ELASTICITY, BUT THE MORE CPU AND MEMORY USED ON THE MASTERS.
- PROVISIONNING 10GB PER SHARDS ALLOWS A FASTER RECOVERY & REALLOCATION.

REPLICATION TRICKS

- NUMBER OF REPLICAS MUST BE 0 OR ODD.
CONSISTENCY QUORUM: $\text{INT} ((\text{PRIMARY} + \text{NUMBER_OF_REPLICAS}) / 2) + 1$.
- RAISE THE REPLICATION FACTOR TO SCALE FOR READING
UP TO 100% OF THE DATA / DATA NODE.

ALIASES

- ACCESS MULTIPLE INDICES AT ONCE.
- READ MULTIPLE, WRITE ONLY ONE.

EXAMPLE ON TIMESTAMPED INDICES:

```
"18_20171020": { "aliases": { "2017": {}, "201710": {}, "20171020": {} } }
```

```
"18_20171021": { "aliases": { "2017": {}, "201710": {}, "20171021": {} } }
```

Queries:

```
/2017/_search
```

```
/201710/_search
```

AFTER A MAPPING CHANGE & REINDEX, CHANGE THE ALIAS:

ROLLOVER

- CREATE A NEW INDEX WHEN TOO OLD OR TOO BIG.
- SUPPORT DATE MATH: DAILY INDEX CREATION.
- USE ALIASES TO QUERY ALL ROLLOVERED INDEXES.

```
PUT "logs-000001" { "aliases": { "logs": {} } }
```

```
POST /logs/_rollover { "conditions": { "max_docs": 10000000 } }
```

- API ONLY!!! YOU NEED A CRON JOB TO





DAILY OPERATIONS

CONFIGURATION CHANGES

- PREFER CONFIGURATION FILE UPDATES TO API CALL FOR PERMANENT CHANGES.
- VERSION YOUR CONFIGURATION CHANGES SO YOU CAN ROLLBACK, ES REQUIRES LOTS OF FINE TUNING.
- WHEN USING _SETTINGS API, PREFER TRANSIENT TO PERSISTENT, THEY'RE EASIER TO GET RID OF.

RECONFIGURING THE WHOLE CLUSTER

LOCK SHARD REALLOCATION & RECOVERY:

```
"cluster.routing.allocation.enable" : "none"
```

OPTIMIZE FOR RECOVERY:

```
"cluster.routing.allocation.node_initial primaries_recoveries": 50
```

```
"indices.recovery.max_bytes_per_sec": "2048mb"
```

RESTART A FULL RACK, WAIT FOR NODES TO COME
BACK:

THE REINDEX API

- IN CLUSTER AND CLUSTER TO CLUSTER REINDEX API.
- ALLOWS CROSS VERSION INDEXING: 1.7 TO 5.1...
- SLICED SCROLLS ONLY AVAILABLE STARTING 6.0.
- ACCEPT ES QUERIES TO FILTER THE DATA TO REINDEX.
- MERGE MULTIPLE INDEXES INTO 1.

BULK INDEXING TRICKS

LIMIT REBALANCE:

```
"cluster.routing.allocation.cluster_concurrent_rebalance": 1
```

```
"cluster.routing.allocation.balance.shard": "0.15f"
```

```
"cluster.routing.allocation.balance.threshold": "10.0f"
```

DISABLE REFRESH:

```
"index.refresh_interval": "0"
```

NO REPLICA:

```
"index.number_of_replicas": "0" // having replica index n times in Lucene, adding one just "rsync" the data.
```

ALLOCATE ON DEDICATED HARDWARE:



OPTIMIZING FOR SPACE & PERFORMANCES

- LUCENE SEGMENTS ARE IMMUTABLE, THE MORE YOU WRITE, THE MORE SEGMENTS YOU GET.
- DELETING DOCUMENTS DOES COPY ON WRITE SO NO REAL DELETE.

`index.merge.scheduler.max_thread_count`: default CPU/2 with min 4

`POST /_force_merge?only_expunge_deletes`: faster, only merge segments with deleted

`POST /_force_merge?max_num_segments`: don't use on indexes you write on!

WARNING: `__FORCE_MERGE` HAS A COST IN CPU AND I/Os. 

MINOR VERSION UPGRADES

- CHECK YOUR PLUGINS COMPATIBILITY, PLUGINS MUST BE COMPILED FOR YOUR MINOR VERSION.
- START UPGRADING THE MASTER NODES.
- UPGRADE THE DATA NODES ON A WHOLE RACK AT ONCE.

OS LEVEL UPGRADES

- ENSURE THE WHOLE CLUSTER RUNS THE SAME JAVA VERSION.
- WHEN UPGRADING JAVA, CHECK IF YOU DON'T HAVE TO UPGRADE THE KERNEL.
- PER NODE JAVA / KERNEL VERSION AVAILABLE IN THE `_STATS` API.

MONITORING

CAPTAIN OBVIOUS, YOU'RE MY ONLY HOPE!

- GOOD MONITORING IS BUSINESS ORIENTED MONITORING.
- GOOD ALERTING IS ACTIONABLE ALERTING.
- DON'T MONITORE THE CLUSTER ONLY, BUT THE WHOLE PROCESSING CHAIN.
- USELESS METRICS ARE USELESS.
- LOSING A DATACENTER: OK. LOSING DATA: NOT OK!

MONITORING TOOLING

- ELASTICSEARCH X-PACK,
- GRAFANA...



LIFE, DEATH & _CLUSTER/HEALTH

- A RED CLUSTER MEANS AT LEAST 1 INDEX HAS MISSING DATA. DON'T PANIC!
- USING `LEVEL={ INDEX, SHARD }` AND AN INDEX ID PROVIDES SPECIFIC INFORMATION.
- LOTS OF PENDING TASKS MEANS YOUR CLUSTER IS UNDER HEAVY LOAD AND SOME NODES CAN'T PROCESS THEM FAST ENOUGH.
- LONG WAITING TASKS MEANS YOU HAVE A CAPACITY PLANNING PROBLEM.

```
{  
  "cluster_name": "blackhole",  
  "status": "green",  
  "timed_out": false,  
  "number_of_nodes": 78,  
  "number_of_data_nodes": 68,  
  "active_primary_shards": 13790,  
  "active_shards": 27580,  
  "relocating_shards": 18,  
  "initializing_shards": 0,  
  "unassigned_shards": 0,  
  "delayed_unassigned_shards": 0,  
  "number_of_pending_tasks": 2,  
  "number_of_in_flight_fetch": 0,  
  "task_max_waiting_in_queue_millis": 6519,  
  "active_shards_percent_as_number": 100  
}
```

USE THE `_CAT` API

- PROVIDES GENERAL INFORMATION ABOUT YOUR NODES, SHARDS, INDICES AND THREAD POOLS.
- HIT THE WHOLE CLUSTER, WHEN IT TIMEOUTS YOU'VE PROBABLY HAVING A NODE STUCK IN GARBAGE COLLECTION.



MONITORING AT THE CLUSTER LEVEL

- USE THE `__STATS` API FOR PRECISE INFORMATION.
- MONITORE THE SHARDS REALLOCATION, TOO MANY MEANS A DESIGN PROBLEM.
- MONITORE THE WRITES AND CLUSTER WIDE, IF THEY FALL TO 0 AND IT'S UNUSUAL, A NODE IS STUCK IN GC.

MONITORING AT THE NODE LEVEL

- USE THE `_NODES / {NODE} , _NODES / {NODE} / STATS` AND `_CAT / THREAD_POOL` API.
- THE GARBAGE COLLECTION DURATION & FREQUENCY IS A GOOD METRIC OF YOUR NODE HEALTH.
- CACHE AND BUFFERS ARE MONITORED ON A NODE LEVEL.
- MONITORING I/Os, SPACE, OPEN FILES & CPU IS CRITICAL.

MONITORING AT THE INDEX LEVEL

- USE THE `{ INDEX } / _STATS` API.
- MONITORE THE DOCUMENTS / SHARD RATIO.
- MONITORE THE MERGES, QUERY TIME.
- TOO MANY EVICTIONS MEANS YOU HAVE A CACHE CONFIGURATION LEVEL.

TROUBLESHOOTING

WHAT'S REALLY GOING ON IN YOUR CLUSTER?

- THE `__NODES / {NODE} / HOT_THREADS` API TELLS WHAT HAPPENS ON THE HOST.
- THE ELECTED MASTER NODES TELLS YOU MOST THING YOU NEED TO KNOW.
- ENABLE THE SLOW LOGS TO UNDERSTAND YOUR BOTTLENECK & OPTIMIZE THE QUERIES. DISABLE THE SLOW LOGS WHEN YOU'RE DONE!!!
- WHEN NOT ENOUGH, MEMORY PROFILING IS YOUR FRIEND.

MEMORY PROFILING

- LIVE MEMORY OR HPROF FILE AFTER A CRASH.
- ALLOWS YOU TO TO KNOW WHAT IS / WAS IN YOUR BUFFERS AND CACHES.
- YOURKIT JAVA PROFILER AS A TOOL.

TRACING

- KNOW WHAT'S REALLY HAPPENING IN YOUR JVM.
- LINUX 4.X PROVIDES GREAT PERF TOOLS, LINUX 4.9 EVEN BETTER:
 - LINUX-PERF,
 - JAVA PERF MAP.
- VECTOR BY NETFLIX (NOT VEKTOR THE TRASH METAL BAND).

QUESTIONS
?

@FDEVILLAMI
L

@SYNTHESIO

SLIDES: [HTTP://BIT.DO/ELASTICSEARCH-SYSADMIN-201](http://bit.do/elasticsearch-sysadmin-201)

