# Hadoop Monitoring, SOLR vs Elastic Search

- Paolo Nesi, Gianni Pantaleo, Imad Zaza

- **DISIT Lab**
- Dipartimento di Ingegneria dell'Informazione, DINFO
- Università degli Studi di Firenze
  Via S. Marta 3, 50139, Firenze, Italy
  Tel: +39-055-4796567,    fax: +39-055-4796363

- **http://www.disit.dinfo.unifi.it** *alias* http://www.disit.org

Ottobre 2020

# Agenda

- ## Monitoring

- Apache HBASE

- Apache Phoenix

- Case studio

- Apache Solr

- Elastic Search

# What really going on !?

- How many maps are issued?

- Where are computed the maps?

- Does my job executing normally or something is going wrong !?

- Does my cluster being underutilzied or overutilized !?

# Basic web interface

# Advanced web interface

# What really really going on !?

- The aforementioned application doesn't give specific information ...

- Unix tools

- Batch System tools

- Really need something that can provide a quick visual overview of the health and load on your cluster

# Ganglia

# Host level detail

# How does Ganglia work?

- Ganglia works through a small agent, *gmond*, on each node or machine to be monitored. You can distribute a single gmond instance to lots of machines at once. Gmonds communicate the state of their local node to a machine running a Master *gmetad* instance.

- The server uses RRDtool to store the data over time

- The Ganglia framework can be extended to monitor many parameters.

# Setup

The software can be downloaded from http://ganglia.sourceforge.net/

Computer A

Computer B

Computer C

Runs gmond

Runs gmond

Runs gmond

Clients just have to run gmond, which is configured by /etc/gmond.conf

Computer D

Runs gmetad

Server to collect the data runs gmetad.

It could also run gmond to monitor itself.

The web interface needs to run on a webserver.

Computer D

gmetad

gmond

httpd

# Client Setup

Computer A

Runs gmond

Computer B

Runs gmond

Computer C

Runs gmond

Apt-get install ganglia-gmond

edit config file

service gmond start

chkconfig gmond on

**/etc/gmond.conf extracts**

```
cluster {
  name = "LCG Workers"
}
/* Feel free to specify as many udp_send_channels as you like.  Gmond
   used to only support having a single channel */
udp_send_channel {
  mcast_join = 239.2.11.95
  port = 8649
}
```

```
udp_send_channel {
  port = 8649
  host = pplxconfig
}
```

```
/* You can specify as many udp_recv_channels as you like as well. */
udp_recv_channel {
  mcast_join = 239.2.11.95
  port = 8649
  bind = 239.2.11.95
}
```

```
/* You can specify as many tcp_accept_channels as you like to share
   an xml description of the state of the cluster */
tcp_accept_channel {
  port = 8649
}
```

# Server Setup

Computer D

gmetad

gmond

httpd

Aptitude install ganglia-gmond ganglia-gmetad ganglia-web

edit /etc/gmond.conf

edit /etc/gmetad.conf

**Extracts from /etc/gmetad.conf**

data_source "LCG Workers" 127.0.0.1:8655 computerA.physics.ox.ac.uk
ComputerB.physics.ox.ac.uk computerC.physics.ox.ac.uk:8655

data_source "LCG Servers" t2se01.physics.ox.ac.uk:8656
t2ce02.physics.ox.ac.uk:8656 gridlogger.physics.ox.ac.uk:8656

# Agenda

- Monitoring
- **Apache HBASE**
- Case studio
- Apache Phoenix
- Apache Solr
- Elastic Search

# Introduction

- Hadoop is a framework that supports operations on a large amount of data.

- Hadoop includes the Hadoop Distributed File System (HDFS)

- HDFS does a good job of storing large amounts of data, but lacks quick random read/write capability.

# Introduction (cont.)

- We need a tabular form to store our data
- The storing system must preserve the advantages of hadoop
  - Fault tolerance
  - High perfomance on large amount of data

- Limits of Relational Databases (RDBMS)
- OLTP vs OLAP
- Big Data storing paradigms…

# Introduction (cont.)

- HBase is a NoSQL (schema-less), column-oriented database.

- It is open source, sparse, consistently distributed, sorted map modeled after Google's BigTable.

- Developed as part of Apache's Hadoop project and runs on top of Hadoop Distributed File System.

- Horizontaly, linearly scalable

# Conceptual View

- A data row has a sortable row key

- Table is a collection of rows.

- Row is a collection of column families.

- Column family is a collection of columns.

- Column is a collection of key value pairs.

- A Time Stamp is designated automatically if not artificially.

- *<family>:<identifier>*

# Conceptual View

| Rowid | Column Family | | | Column Family | | | Column Family | | | Column Family | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | clo1 | col 2 | col 3 | col 1 | col 2 | col 3 | col1 | col2 | col3 | col1 | col2 | col3 |
| 1 | | | | | | | | | | | | |
| 2 | | | | | | | | | | | | |
| 3 | | | | | | | | | | | | |

**COLUMN FAMILIES**

| Row key | personal data | | professional data | |
|---|---|---|---|---|
| empid | name | city | designation | salary |
| 1 | raju | hyderabad | manager | 50,000 |
| 2 | ravi | chennai | sr.engineer | 30,000 |
| 3 | rajesh | delhi | jr.engineer | 25,000 |

# Hbase table

# Let's try it !

# Physical Storage View

- Physically, tables are stored on a per-column family basis.

- Empty cells are not stored in a column-oriented storage format.

- Each column family is managed by an HStore.

Memcache

Data MapFile
Index MapFile

Key/Value

Index key

# Hbase table



**HStore**

<family>:<label>

# Physical Storage View

- Each cell value of the table has a timestamp

- Subsequent column values are stored contiguously on the disk.

# What we can do ?

- Mysql-like commands
  - Create table
  - Drop table
  - Insert data in table
  - Delete data
  - Query data
- Mysql-like interface (cli, API)

```
$ hbase shell
```

# Hbase commands

- Table manipulations
  - create 't1', 'f1', 'f2', 'f3'

    **create '<table-name>', '<column-family1>', '<column-family2>' ……**

  - alter 'tablename', NAME => '<new-column-family>'

  - Disable/enable 't1'

    **enable/disable '<table-name>'**

  - drop 't1' (but t1 must be disabled)

    **drop '<table-name>'**

# More commands

- put 't1', 'r1', 'c1', 'value1', ts1

> **put '<table-name>', 'row-key', 'columnfamily:columnname', 'value', 'timestamp'**

- delete 't1', 'r1', 'c1'

- delete 't1', 'r1', 'c1','ts1'

- get 't1', 'r1'            →        scan 't1'

- get 't1', 'r1', {TIMERANGE => [ts1, ts2]} →  scan 't1', {TIMERANGE => [ts1, ts2]}

- get 't1', 'r1', {COLUMN => 'c1'} → scan 't1', {COLUMN => 'c1'}

# More commands

- scan 'syslog', {COLUMNS => ['msg:body', 'msg:timestamp'], LIMIT => 100}

- scan 'syslog', {COLUMN => 'msg:body', FILTER => "ValueFilter( =, 'binaryprefix:*<STRING_TO_MATCH>*')"}

- scan 'syslog', {COLUMN => 'msg:body', FILTER => "ValueFilter( =,'regexstring:.*prova.*')"}

# Trace

- create 'students', cf=['a','b']
- put 'students', 'pippo', 'a:name', 'pippo'
- put 'students', 'pippo', 'a:age', 20
- put 'students', 'pluto, 'a:name', 'pluto'
- scan 'students'

# HBASE Client

- ## Cli : hbase

- ## Java API
  - HTable hTable = new HTable(config, "students");
  - Put p = new Put(Bytes.toBytes("row1"));
  - p.add(Bytes.toBytes("personal"),Bytes.toBytes("name"),Bytes.toBytes("PIPPO"));
  - p.add(Bytes.toBytes("personal"),Bytes.toBytes("age"),Bytes.toBytes("21"));
  - hTable.put(p); hTable.close();

  - Get g = new Get(Bytes.toBytes("row1")); Result result = hTable.get(g);
  - byte [] name = result.getValue(Bytes.toBytes("personal"),Bytes.toBytes("name"));

  - Delete delete = new Delete(Bytes.toBytes("row1"));
  - delete.deleteColumn(Bytes.toBytes("personal"), Bytes.toBytes("age"));
  - hTable.delete(delete);

# Hbase Mapreduce Flavor

- hbase org.apache.hadoop.hbase.mapreduce.RowCounter <tablename>

# Hbase Mapreduce Program

```
Configuration config = HBaseConfiguration.create();
Job job = new Job(config, "ExampleRead");
job.setJarByClass(MyReadJob.class);    // class that contains mapper

Scan scan = new Scan();
scan.setCaching(500);       // 1 is the default in Scan, which will be bad for MapReduce jobs
scan.setCacheBlocks(false);  // don't set to true for MR jobs
// set other scan attrs
...

TableMapReduceUtil.initTableMapperJob(
  tableName,       // input HBase table name
  scan,           // Scan instance to control CF and attribute selection
  MyMapper.class,   // mapper
  null,          // mapper output key
  null,          // mapper output value
  job);
job.setOutputFormatClass(NullOutputFormat.class);   // because we aren't emitting anything from mapper

boolean b = job.waitForCompletion(true);
if (!b) {
  throw new IOException("error with job!");
}
```

# Hbase Mapreduce Program Mapper

```
public static class MyMapper extends
TableMapper<Text, Text> {


  public void map(ImmutableBytesWritable row,
Result value, Context context) throws
InterruptedException, IOException {
    // process data for the row from the Result
instance.
    }
}
```

# Hbase Mapreduce Program Sink and source

Count the number of distinct instances of a value in a table and write those summarized counts in another table.

```
Configuration config = HBaseConfiguration.create();
Job job = new Job(config,"ExampleSummary");
job.setJarByClass(MySummaryJob.class);     // class that contains mapper and reducer

Scan scan = new Scan();
scan.setCaching(500);        // 1 is the default in Scan, which will be bad for MapReduce jobs
scan.setCacheBlocks(false);  // don't set to true for MR jobs
// set other scan attrs

TableMapReduceUtil.initTableMapperJob(
            sourceTable,       // input table
            scan,             // Scan instance to control CF and attribute selection
            MyMapper.class,    // mapper class
            Text.class,        // mapper output key
            IntWritable.class, // mapper output value
            job);
TableMapReduceUtil.initTableReducerJob(
            targetTable,       // output table
            MyTableReducer.class,   // reducer class
            job);
job.setNumReduceTasks(1);   // at least one, adjust as required

boolean b = job.waitForCompletion(true);
if (!b) {
            throw new IOException("error with job!");
}
```

# Hbase Mapreduce Program Sink and source: Mapper

A column with a String-value is chosen as the value to summarize upon. This value is used as the key to emit from the mapper, and an IntWritable represents an instance counter.

```
public static class MyMapper extends TableMapper<Text, IntWritable> {

        private final IntWritable ONE = new IntWritable(1);
        private Text text = new Text();

        public void map(ImmutableBytesWritable row, Result value, Context context)
throws IOException, InterruptedException {
        String val = new String(value.getValue(Bytes.toBytes("cf"),
Bytes.toBytes("attr1")));
        text.set(val);    // we can only emit Writables...

        context.write(text, ONE);
        }
}
```

# Hbase Mapreduce Program Sink and source: Reducer

In the reducer, the "ones" are counted (just like any other MR example that does this), and then emits a Put.

```
public static class MyTableReducer extends TableReducer<Text, IntWritable,
ImmutableBytesWritable>  {

        public void reduce(Text key, Iterable<IntWritable> values, Context context)
throws IOException, InterruptedException {
                int i = 0;
                for (IntWritable val : values) {
                        i += val.get();
                }
                Put put = new Put(Bytes.toBytes(key.toString()));
                put.add(Bytes.toBytes("cf"), Bytes.toBytes("count"),
Bytes.toBytes(i));

                context.write(null, put);
        }
}
```

# Let's try it !

# Missing !?

- Mysql is sql complaint
  - To search data we issue commands like
    - «select columns where conditions >>
  - Generally  we take advantgae of joining table


- HBASE is essentially an hashmap
  - No join of tables
  - No select

# Filters

- Filters are Java classes restricting matches;

- Filter list: combines multiple filters with AND and OR

- Compare values of one or multiple columns
  - Smaller, equal, greater, substring, prefix,

- Compare metadata: column family and qualifier    Qualifier prefix filter: Return (first few) matching columns
  - Column range filter: return a slice of columns (e.g. bb-bz)

- Compare names of rows Note: it is preferable to use scan options

# Scan advanced

- scan 'student',{ FILTER => "KeyOnlyFilter()"}

- scan 'student',{ FILTER => "(PrefixFilter ('pi')) AND MultipleColumnPrefixFilter('a','b') AND ColumnCountGetFilter(2)" }

- scan 'student',{ COLUMNS=>['a:age'], FILTER => "SingleColumnValueFilter('a','age',>,'binary:19')"}

# HBase Components

- **Region**
  - A subset of a table's rows, like horizontal range partitioning
  - Automatically done
- **RegionServer (many slaves)**
  - Manages data regions
  - Serves data for reads and writes (*using a log*)
- **Master**
  - Responsible for coordinating the slaves
  - Assigns regions, detects failures
  - Admin functions

# Big Picture

# ZooKeeper

- HBase depends on ZooKeeper
- By default HBase manages the ZooKeeper instance
  - E.g., starts and stops ZooKeeper
- HMaster and HRegionServers register themselves with ZooKeeper

# HBaseMaster

- Assign regions to HRegionServers.

1. ROOT region locates all the META regions.

2. META region maps a number of user regions.

3. Assign user regions to the HRegionServers.

- Enable/Disable table and change table schema

- Monitor the health of each Server

Master

2 META Region

2 META Region

2 META Region

2 META Region

1 ROOT Region

Server  Server  Server  Server  Server

USER Region

META Region

ROOT Region

USER Region

META Region

USER Region

# ROOT/META Table

- Each row in the ROOT and META tables is approximately 1KB in size. At the default size of 256MB.

$$1\ ROOT\ table$$
$$= 2^{18}\ META\ regions$$
$$= 2^{18} \times 2^{18} USER\ regions$$
$$= 2^{54} KB = 2^{64} bytes$$

$2^{24}$TB

# HRegionServer

- **Write Requests**
- Read Requests
- Cache Flushes
- Compactions
- Region Splits

write

HLog

| Row key | Time Stamp | Column "contents:" | Column "anchor:" | |
|---|---|---|---|---|
| "com.apache.www" | t12 | "\<html\>…" | | |
| | t11 | "\<html\>…" | | |
| | t10 | | "anchor:apache.com" | "APACHE" |
| "com.cnn.www" | t9 | | "anchor:cnnsi.com" | "CNN" |
| | t8 | | "anchor:my.look.ca" | "CNN.com" |
| | t6 | "\<html\>…" | | |
| | t5 | "\<html\>…" | | |
| | t3 | "\<html\>…" | | |

Mapfile1.1
Mapfile1.2

Memcache2

Memcache1

**Hstore1**

**Hstore2**

October 2020

45

# HRegionServer

Read

- Write Requests

- Read Requests

- Cache Flushes

- Compactions

- Region Splits

| Row key | Time Stamp | Column "contents:" | Column "anchor:" | |
|---|---|---|---|---|
| "com.apache.www" | t12 | "<html>…" | | |
| | t11 | "<html>…" | | |
| | t10 | | "anchor:apache.com" | "APACHE" |
| "com.cnn.www" | t9 | | "anchor:cnnsi.com" | "CNN" |
| | t8 | | "anchor:my.look.ca" | "CNN.com" |
| | t6 | "<html>…" | | |
| | t5 | "<html>…" | | |
| | t3 | "<html>…" | | |

Mapfile1.1
Mapfile1.2

Memcache1

**Hstore1**

October 2020

46

# HRegionServer

Cache Flushes

HLog

- Write Requests
- Read Requests
- Cache Flushes
- Compactions
- Region Splits

| Row key | Time Stamp | Column "contents:" | Column "anchor:" | |
|---|---|---|---|---|
| "com.apache.www" | t12 | "<html>…" | | |
| | t11 | "<html>…" | | |
| | t10 | | "anchor:apache.com" | "APACHE" |
| "com.cnn.www" | t9 | | "anchor:cnnsi.com" | "CNN" |
| | t8 | | "anchor:my.look.ca" | "CNN.com" |
| | t6 | "<html>…" | | |
| | t5 | "<html>…" | | |
| | t3 | "<html>…" | | |

Mapfile1.1
Mapfile1.2
Mapfile1.3

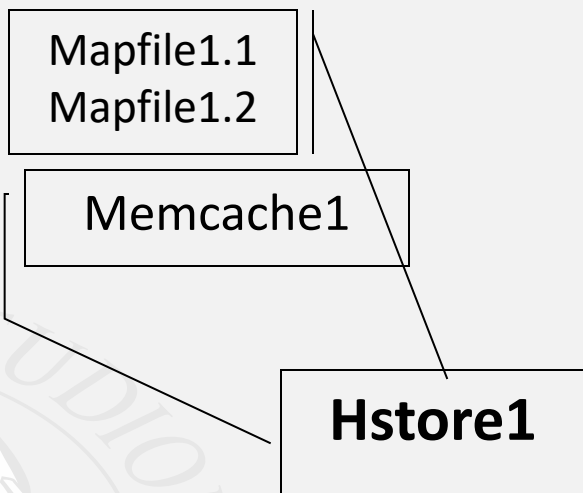Mapfile1.1
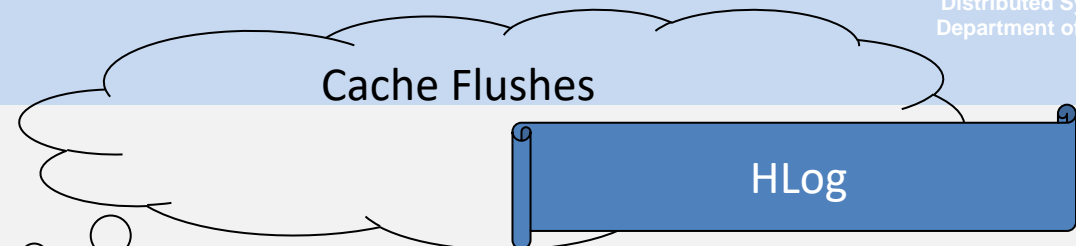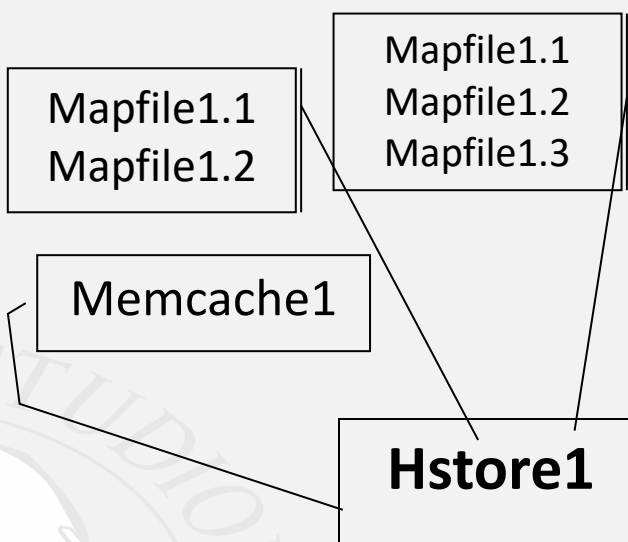Mapfile1.2

Memcache1

**Hstore1**

# HRegionServer

Compactions

- Write Requests
- Read Requests
- Cache Flushes
- Compactions
- Region Splits

| Row key | Time Stamp | Column "contents:" | Column "anchor:" | |
|---|---|---|---|---|
| "com.apache.www" | t12 | "<html>…" | | |
| | t11 | "<html>…" | | |
| | t10 | | "anchor:apache.com" | "APACHE" |
| "com.cnn.www" | t9 | | "anchor:cnnsi.com" | "CNN" |
| | t8 | | "anchor:my.look.ca" | "CNN.com" |
| | t6 | "<html>…" | | |
| | t5 | "<html>…" | | |
| | t3 | "<html>…" | | |

Mapfile1.1
Mapfile1.2

Mapfile1

Memcache1

**Hstore1**

# HRegionServer

Region Splits

- Write Requests
- Read Requests
- Cache Flushes
- Compactions
- **Region Splits**

Mapfile1

Memcache1

**Hstore1**

| Row key | Time Stamp | Column "contents:" | Column "anchor:" | |
|---------|------------|--------------------|------------------|---|
| "com.apache.www" | t12 | "<html>…" | | |
| | t11 | "<html>…" | | |
| | t10 | | "anchor:apache.com" | "APACHE" |
| "com.cnn.www" | t9 | | "anchor:cnnsi.com" | "CNN" |
| | t8 | | "anchor:my.look.ca" | "CNN.com" |
| | t6 | "<html>…" | | |
| | t5 | "<html>…" | | |
| | t3 | "<html>…" | | |

# Agenda

- Monitoring
- Apache HBASE
- **Case studio**
- Apache Phoenix
- Apache Solr
- Elastic Search

# Solution

- ## Migrating to a distribuited architcure
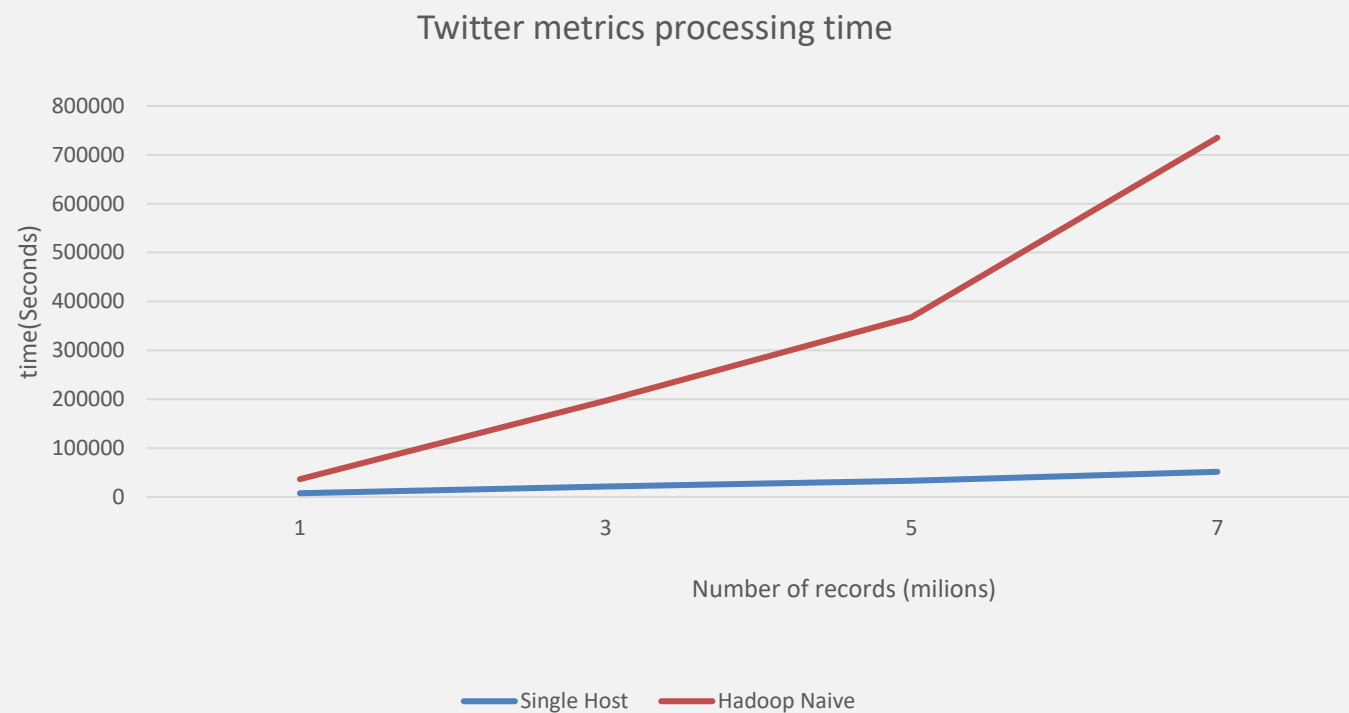  - We have to rethink the problem in parallel way

# Sequential minded approach

- For each search execute a MR job
  - If a MR job takes 1 minute (ideally) time and i have to compute for 1000
    -> i have to wait 1000 minute !

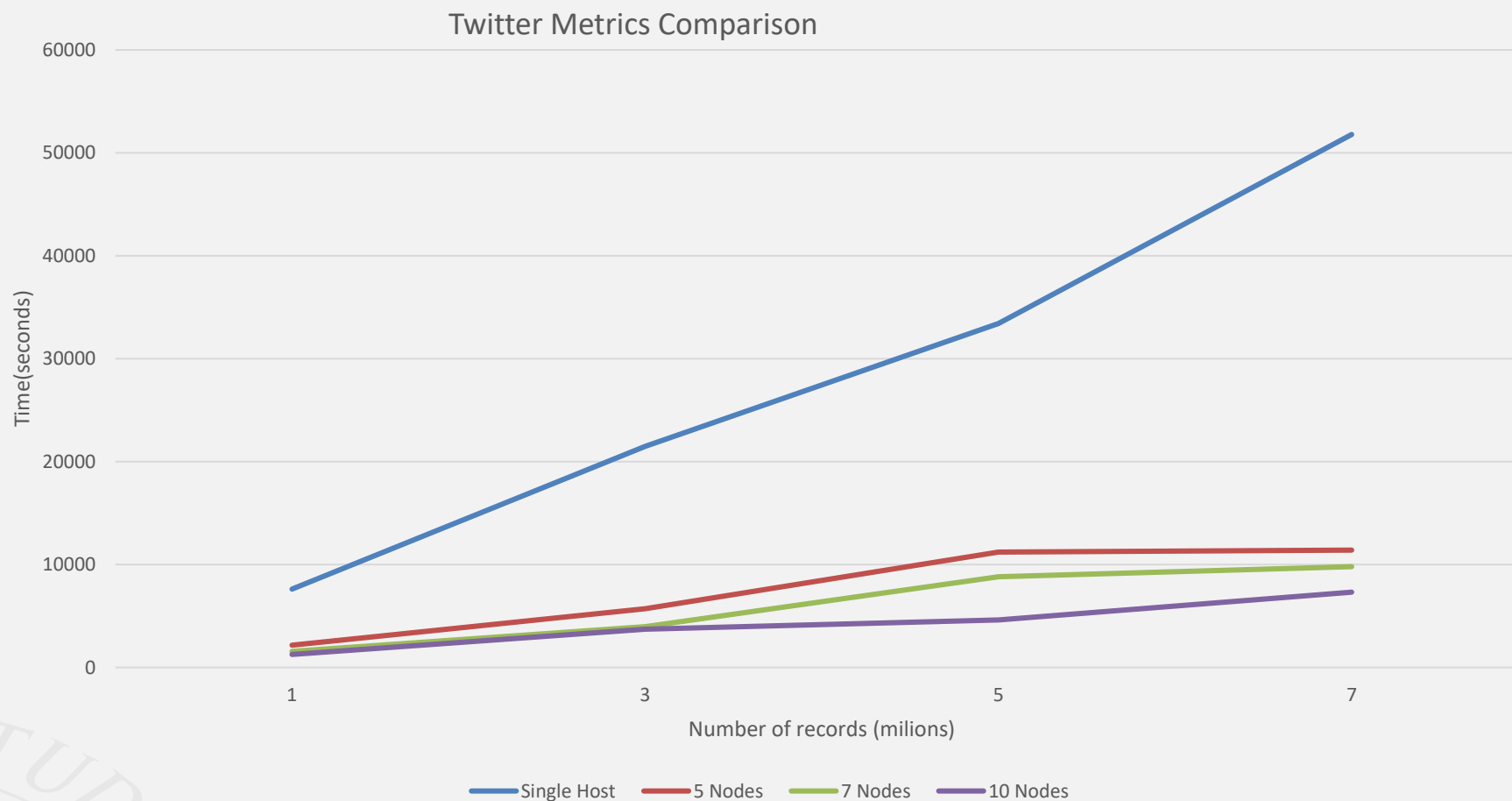# Single Host vs Hadoop



Twitter metrics processing time

# Parallel-computing minded aproach

- Scan all records and counts each word sorting per day

- As final step consider only the words involded in our search

# Single Host vs Hadoop



Twitter Metrics Comparison

# *All Quiet on the Western Front …*

| Google calls it: | Hadoop equivalent: |
|---|---|
| MapReduce | Hadoop |
| GFS | HDFS |
| Bigtable | HBase |
| Chubby | Zookeeper |

- 1979: First commercial SQL RDBMs

-  1990: Transaction processing on SQL now popular

- 2006: Hadoop and other "big data" technologies

- 2008: NoSQL

-  2011: SQL on Hadoop

-  2014: Interactive analytics on Hadoop and NoSQL with SQL (Phoenix)

# *Some References*

I. Nesi, Paolo, Gianni Pantaleo, and Gianmarco Sanesi. "A hadoop based platform for natural language processing of web pages and documents." *Journal of Visual Languages & Computing* 31 (2015): 130-138.

II. Michael, Maged, et al. "Scale-up x scale-out: A case study using nutch/lucene." *2007 IEEE International Parallel and Distributed Processing Symposium*. IEEE, 2007.

III. Appuswamy, Raja, et al. "Scale-up vs scale-out for hadoop: Time to rethink? " *Proceedings of the 4th annual Symposium on Cloud Computing*. ACM, 2013.

- Apache Hadoop ver. 3.0.0

- Apache Hbase ver 1.2.6

- Apache Phoenix ver 4.14.0

# Agenda

- Monitoring

- Apache HBASE

- Case studio

- **Apache Phoenix**

- Apache Solr
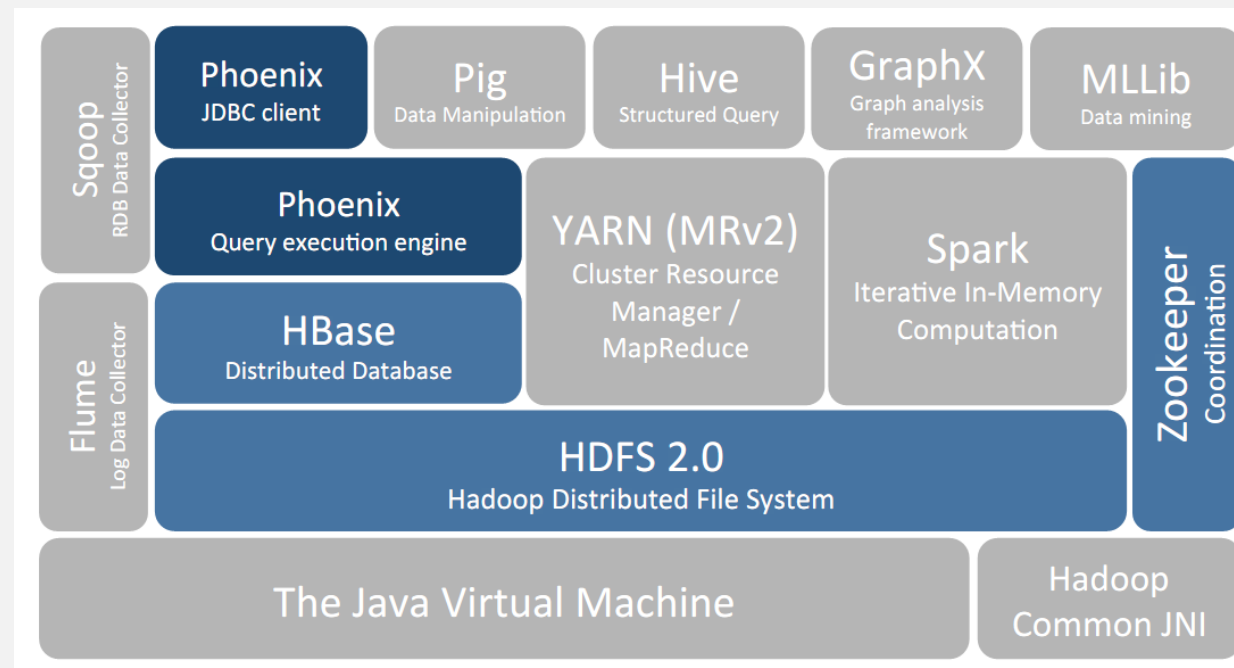
- Elastic Search

# Why

- Hbase is Ok BUT
  - It takes too much expertise to write an application
  - It takes too much code to do anything
  - Your application is tied too closely with your data model

# Apache Phoenix

- **What**: A relational database layer for Apache Hbase

  - Low Latency Query Model & SQL support over HBase API
  - SQL query compiled it into a series of HBase scans
  - Metadata is stored in an HBase table and versioned
  - Pushes Computation to the HBase Region Servers

- **How: Coprocessors** (Server-Side) **Minimize Data Transfer** Uses **Native HBase APIs** Not Map/Reduce framework of Hbase
- A JDBC driver

# Apache Phoenix

# Performance

- Pushes down computatioon to region servers
  - Start/stop key range(s)
  - Time range min/max
  - Predicates
  - Aggregation
  - Sort
  - Limit
- Parallelizes query from client
  - Intra-region through staBsBcs collecBon
- Supports secondary indexes
  - Global & co-located

# Agenda

- Monitoring
- Apache HBASE
- Case studio
- Apache Phoenix
- **Apache Solr**
- Elastic Search

# Overview

- Introduction to Lucene & Solr

  – Getting started

  – Indexing using Solr

  – Updating & deleting files

  – Searching using Solr

# What is Lucene?

- Lucene is An open source Java-based IR library enabling text based search

# What is Solr

- Solr is:
  - An open source enterprise search server
  - Based on the Lucene Java search library
  - A web based application that processes HTTP request and returns HTTP responses
  - completed with XML/HTTP APIs, caching, replication, and web administration interface.

# Why solr

- Using many Lucene best practices
- Easy setup, configuration and Easy to extend
- Providing faceted navigation, spell checking, highlighting, clustering, grouping, & any other search features
- Supporting clients in:
  - HTTP
  - Java
  - Python
  - PHP
  - Ruby
  - JSON

# Why Solr

- ## Some reasons of using Solr:
  - ### Good indexing performance
  - ### SolrCloud feature (Solr 4.x above)
  - ### Multi-Faceted searches
  - ### Geospatial searches
  - ### Who uses Lucene/Solr?
    - Cisco, ebay, Boeing, AT& T, Ford
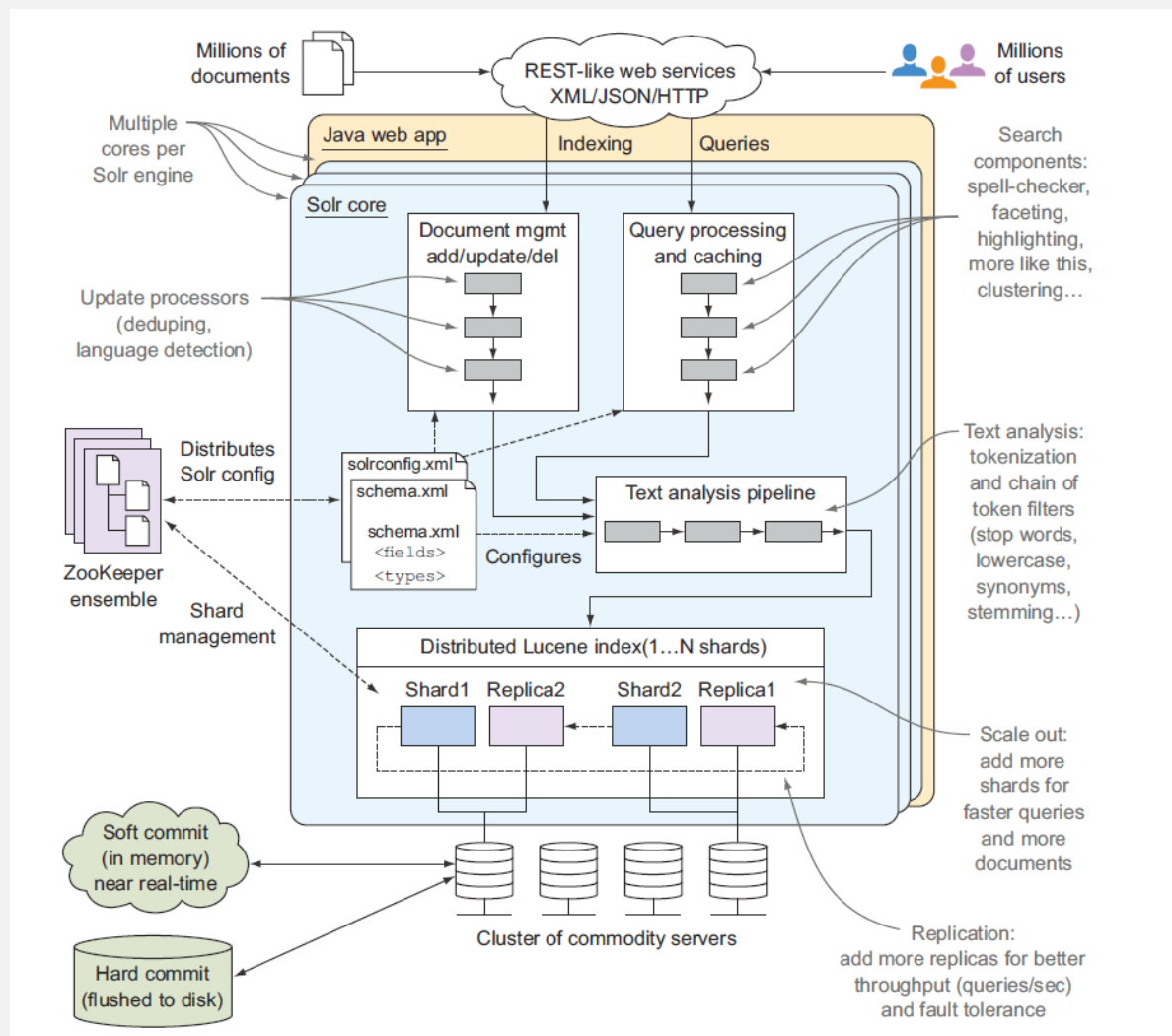      and many, many others...!

# Comparison to Database Technology

- The most important comparison to make is the data model
  - Data model is the organizational structure of data
- RDBMS:
  - Its data model is based on multiple tables with lookup keys between them
  - A join capability for querying across tables
  - A flexible data model
- Lucene Solr:
  - Has a document oriented data model
    - Analogous to a single table without join possibilities
    - Document-oriented databases have a rich nested structure similar to XML/JSON → MongoDB (NoSQL)
- Has a flat document structure
  - Supporting multi-valued fields with an array of values

# Solr Architecture

# Solr Architecture



Source: «Solr in Action», (T. Greinger, T. Potter) – Manning Ed. (2014)

# Solr Interface & Functionalities

# Solr Interface & Functionalities

# Solr API

**COLLECTIONS API**

create -c NewC

curl http://localhost:8983/solr/admin/cores?action=CREATE&name=NewC&instanceDir=NewC

**INSERT / UPDATE API**

```
curl http://localhost:8983/solr/NewC/update -d "[
{          \"id\" : \"book1\",
           \"title_t\" : \"Solr In Action\",
           \"author_s\" : \"Timothy Potter\"
}]"
```

**SEARCH API**

http://localhost:8983/solr/<COLLECTION>/select?q=<QUERY>

curl http://localhost:8983/solr/techproducts/select?q=<KEYWORD>

curl http://localhost:8983/solr/techproducts/select?q=cat:electronics

# Solr API

http://localhost:8983/solr/techproducts/select?q=cat:electronics&fl=id,name,price

## Solr Quert Parameters



- **qt** – Query handler for the request. Standard query handler is used if not specified.
- **q** – It is used to specify the query event.
- **fq** – Used to specify filter queries.
- **sort** – Used to sort the results in ascending or descending order.
- **start, rows** – start specifies the staring number of the result set. By default it is zero. rows specify the number of records to return.
- **fl** – Used to return selective fields.
- **wt** – Specifies the response format. Default is XML.
- **indent** – Setting to true makes the response more readable.
- **debugQuery** – Setting the parameter to true gives the debugging information as part of response.
- **dismax** –  To specify the dismax parser.
- **edismax** – To specify the edismax parser.
- **facet** – Setting to true enables the faceting.
- **spatial** – Used for geospatial searches.
- **spellcheck** – Setting to true help in searching similar terms.

# Solr Faceted Search

http://localhost:8983/solr/techproducts/select?q=price:[0 TO 400]&fl=id,name,price&facet=true&facet.field=cat

# Agenda

- Monitoring
- Apache HBASE
- Apache Phoenix
- Case studio
- Apache Solr
- **Elastic Search**

# Elasticsearch

- real time
- Search & Analytics Engine
- Distributed
- Scales massively
- High availability
- Restful api
- Json over HTTP
- Schema free
- Multi tenancy
- Open source
- Lucene based

# API

➢curl –X GET localhost:9200/?pretty

➢Input Data

➢Retrieve Data

➢Update Data

➢Delete Data

# Input Data

PUT /index/type/id

```
PUT /myapp/tweet/1 -d '
{ "tweet": "A tweet text",
"nick": "@twitternick",
"name": "Pippo",
'
```

# Retrieve Data

- GET /myapp/tweet/1

```
{
 "_index": "myapp",
 "_type": "tweet",
 "_id": "1",
 "_version": 1,
 "exists": true,
 "_source": { ...OUR
TWEET... }
}
```

# Update Data

- PUT /myapp/tweet/1 -d '

```
{
"tweet": "I know #elasticsearch is AWESOME",
"nick": "@clintongormley",
"name": "Clinton Gormley",
"date": "2013-06-03",
"rt": 5,
"loc": {
  "lat": 13.4,
  "lon": 52.5
}
}
'
```

```
{
 "_index": "myapp",
 "_type": "tweet",
 "_id": "1",
 "_version": 2,
 "ok": true
}


# atomic delete and put
```

# Delete Data

- DELETE /myapp/tweet/1

```
{
 "_index": "myapp",
 "_type": "tweet",
 "_id": "1",
 "_version": 3,
 "ok": true,
 "found": true
}
```

# RDBMS comparison

SQL Server manage Databases i.e. Tables i.e. Columns/Rows

# RDBMS comparison

Elastic Search manage  Indices   i.e. Types i.e. Documents with Properties

An Elastic Search cluster can contain multiple **Indices** (databases), which in turn contain multiple **Types** (tables). These types hold multiple **Documents** (rows), and each document has **Properties**(columns).

# Glossary

- Node: A node is a running instance of elasticsearch which belongs to a *cluster.*

- Shard: A shard is a single Lucene instance. It is a low-level "worker" unit which is managed automatically by elasticsearch. An index is a logical namespace which points to primary and replica shards.

- Primary Shard: Each document is stored in a single *primary shard*. When you index a document, it is *indexed* first on the *primary shard,* then on all *replicas* of the primary shard.

- Replica Shard: Each *primary shard* can have zero or more *replicas*. A replica is a copy of the primary shard, and has two purposes: a) increase fail over b) increase performance

# Glossary

- Index: An index is like a *database* in a relational database.
- Type:  A type is like a *table* in a relational database. Each type has a list of fields that can be specified for documents of that type.
- Document: JSON document which is stored in elasticsearch. It is like a row in a table in a relational database.
- Field: A *document* contains a list of fields, or key-value pairs. The value can be a simple (scalar) value (eg a string, integer, date), or a nested structure like an array or an object. A field is similar to a column in a table in a relational database.
- Mapping:  mapping is like a *schema definition* in a relational database. The mapping defines how each field in the document is *analyzed*.
- Routing: When you index a document, it is stored on a single primary shard. That shard is chosen by hashing the *routing* value. By default the *routing* value is derived from the ID of the document.

# Core Field Types

- Strings:                string

- Datetimes:               date

- Whole numbers:      byte, short, integer, long

- Floats:            float, double

- Booleans:                boolean

- Objects:        object


- Also: multi_field, ip, geo_point, geo_shape,

# Auto Detection of Field

- "foo bar"                    string
- "2018-10-09"              date
- 10                byte, short, integer, long
- 10.0                    float, double
- true                    boolean
- { foo: "bar" }          object

- ["foo","bar"]        No special mapping. Any field can have multi-values

# Some more Glossary

- Term: A term is an exact value that is indexed in elasticsearch. The terms foo, Foo, FOO are NOT equivalent.

- Text: Text (or full text) is ordinary unstructured text, such as this paragraph. By default, text will be analyzed into terms, which is what is actually stored in the index. Text fields need to be analyzed at index time in order to be searchable as full text, and keywords in full text queries must be analyzed at search time to produce (and search for) the same terms that were generated at index time.

- Analysis: Analysis is the process of converting *full text* to *terms*. Depending on which analyzer is used, these phrases: FOO BAR, Foo-Bar, foo,bar will probably all result in the terms foo and bar. These terms are what is actually stored in the index.

# Some more glossary (Cont.)

- Tokenizer: Tokenizers are used to break a string down into a stream of terms or tokens. A simple tokenizer might split the string up into terms wherever it encounters whitespace or punctuation.

- Facets: They enable you to calculate and summarize data about the current query on-the-fly. They can be used for all sorts of tasks such as dynamic counting of result values or even distribution histograms. Facets only perform their calculations one-level deep, and they cant be easily combined.

- Aggregations: Aggregations are similar to facets in many ways, and overcome the limitations of facets. Indeed, aggregations are meant to eventually replace facets altogether. Facets are and should be considered deprecated and will likely be removed in one of the future major releases. One of the major limitations of facets is that you can't have facets of facets. Which is to say, facets cannot be nested. The ability to nest aggregations therefore brings a great deal of power that was missing in facets.
  - The two broad families of aggregations are metrics aggregations and bucket aggregations. Metrics aggregations calculate some value (like an average) over a set of documents, and bucket aggregations group documents into buckets.
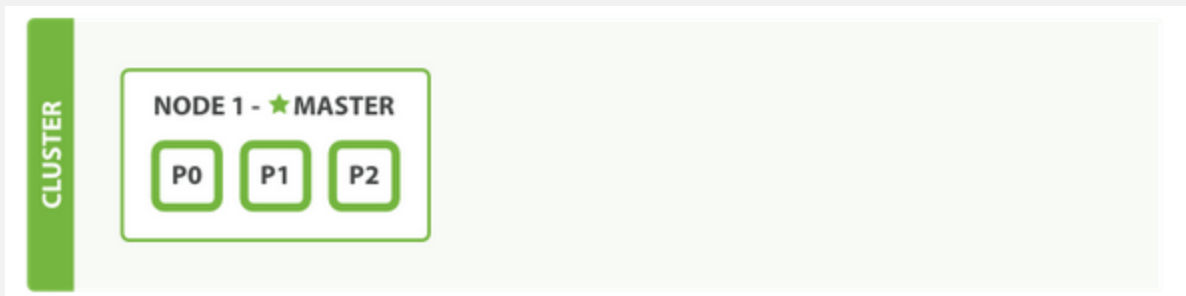
- Schemaless, Document Oriented

- No need to configure schema upfront

- No need for slow ALTER TABLE – like operations

- Define mapping (schema) to customize the indexing process
  - Require fields to be of certain type
  - If you want text fields that should not be analyzed

# Distributed & Highly Available

- Multiple nodes running in a cluster
  - Acting as a single service
  - Nodes in cluster that store data or nodes that just help in speeding up search queries
- Sharding
  - Indices are sharded (#shards are configurable)
  - Each Shard can have zero or more replicas
    - Replicas on different servers for failover
- Master
  - Automatic Master detection + failover
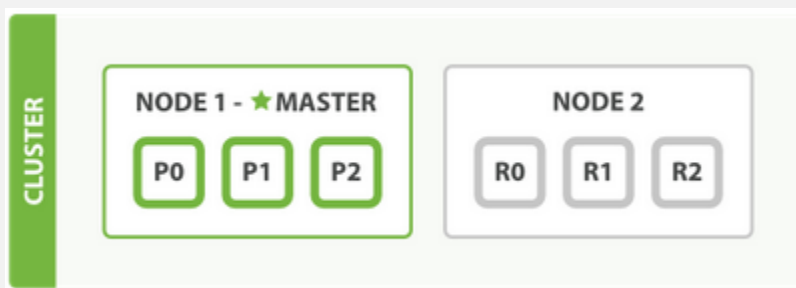  - Responsible for distribution/balancing of shards

# A Single Node Cluster with An Index

- All 3 Primary Shards allocated to Node1

- No replication Nodes

- A single node means single point of failure
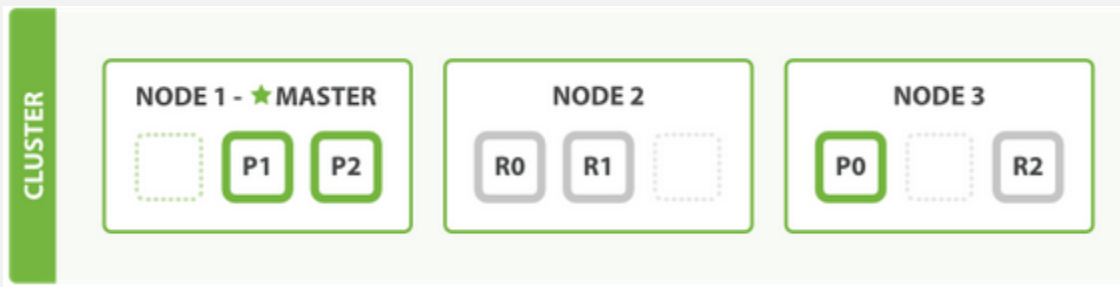
- Health of the Cluster: Yellow

# Add Failover

- Add one Node to cluster by configuring the cluster name.

- 3 replica shards have been allocated.

-  Cluster Health : Green.

- Now 6 Shards. There is redundancy

# Scale horizontally

- A 3 node cluster

- One shard each from Node 1 and Node 2 have moved to Node 3

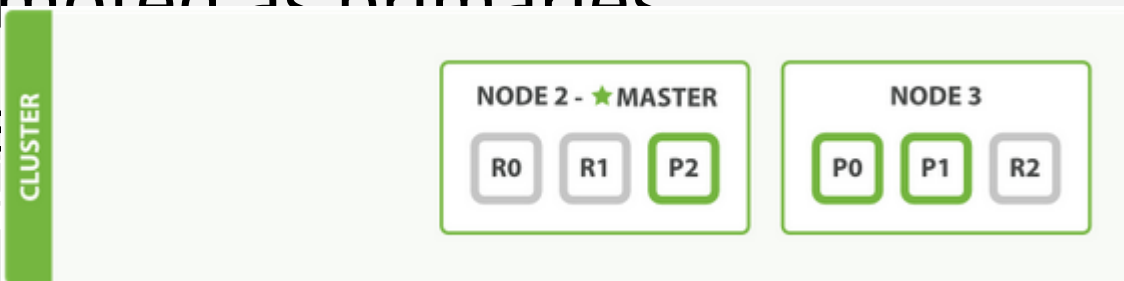- Better performance as hardware resouces (CPU,RAM, I/O) are shared

# Scale some more

- More  Nodes can be added
- More replicas can be added
- This will allow faster searches
- Allows better redundancy
- However the number of **primary shards is fixed** at the moment an index is created.
- Effectively, the maximum amount of data that can be stored in the index is defined by this number.
- Is this a limitation …..

# Coping with Node Failure

- Kill Master Node

- Elect a New Master (Node 2)

- Primary Shard 1 and 2 were lost

- Cluster Health : Red

- Node 2 & 3 have Replicas of these shards, which are now promoted as primaries

- Clus

# Beauty Of Elastic Search

- In Elasticsearch, **all data in every field** is **indexed by default**. That is, every field has a dedicated inverted index for fast retrieval. And, unlike most other databases, it can use all of those inverted indices **in the same query**, to return results at breathtaking speed