



elastic



INNOVATION  
SOFTWARE  
LEARN. EMPOWER. INNOVATE.

┌ search as  
you type ┐



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# search- as-you-type

```
curl -XGET '127.0.0.1:9200/movies/movie/_search?pretty' -d '{
  "query": {
    "match_phrase_prefix": {
      "title": {
        "query": "star trek", "slop": 10
      }
    }
  }
}'
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# index with N-grams

“star”:

unigram:	[ s, t, a, r ]
bigram:	[ st, ta, ar ]
trigram:	[ sta, tar ]
4-gram:	[ star ]

*edge n-grams* are built only on the beginning of each term.

# indexing n-grams

Create an  
“autocomplete”  
analyzer

filter for edge n-grams  
min = 1  
max = 20

custom analyzer, in addition to  
standard lowercase filter also  
has autocomplete filter (n-  
grams)

```
curl -XPUT '127.0.0.1:9200/movies?pretty' -d '{
  "settings": {
    "analysis": {
      "filter": {
        "autocomplete_filter": {
          "type": "edge_ngram",
          "min_gram": 1,
          "max_gram": 20
        }
      },
      "analyzer": {
        "autocomplete": {
          "type": "custom", "tokenizer": "standard",
          "filter": ["lowercase",
                    "autocomplete_filter"]
        }
      }
    }
  }
}
```



I N

LEARN. EMPOWER. INNOVATE.

# map your field

```
curl -XPUT '127.0.0.1:9200/movies/_mapping/movie?pretty' -d '{
  "movie": {
    "properties": {
      "title": {
        "type": "string",
        "analyzer": "autocomplete"
      }
    }
  }
}'
```

title is of type "string" and uses our custom analyzer "autocomplete"

# n-grams only on index

Use n-grams only on the index side or query will also get split into n-grams, and we'll get results for everything that matches 's', 't', 'a', 'st', etc.

```
curl -XGET 127.0.0.1:9200/movies/movie/_search?pretty -d '{
  "query": {
    "match": {
      "title": {
        "query": "sta",
        "analyzer": "standard"
      }
    }
  }
}'
```

analyzer = standard so don't split up what was typed into n-grams.

# completion suggester

You can also upload a list of all possible completions ahead of time using **completion suggester**.

- Most customizable
- Reliable results
- Most control

Suggesters

Completion – auto-complete/search-as-you-type

Term – spell correction

Phrase – did-you-mean

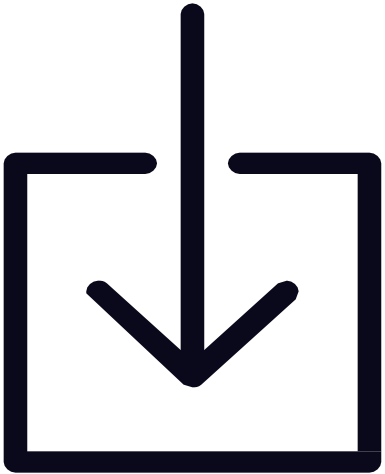


# lab: Query tips

- Lab 10: Pagination, sorting, filtering and fuzzy matching
- Lab 11: Prefix/wildcard and auto-completion

┌ importing  
data └

# importing data



stand-alone **scripts** can submit bulk documents via REST API

**logstash** and **beats** can stream data from logs, S3, databases, and more

AWS systems can stream in data via **lambda** or **kinesis firehose**

**kafka**, **spark**, and more have Elasticsearch integration add-ons

┌ importing  
via script / json └



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# python import

- read in data from some distributed filesystem
- transform it into JSON bulkinserts
- submit via HTTP / REST to your elasticsearch cluster

```
import csv
import re

csvfile = open('ml-latest-small/movies.csv', 'r')

reader = csv.DictReader( csvfile )
for movie in reader:
    print ( "{ \"create\" : { \"_index\": \"movies\", \"_type\": \"movie\", \"_id\" : \"\" , movie['movieId']
    title = re.sub(" \\.*\)$", "", re.sub('\"', '', movie['title']))
    year = movie['title'][-5:-1]
    if (not year.isdigit()):
        year = "2016"
    genres = movie['genres'].split('|')
    print ( "{ \"id\": \"\", movie['movieId'], "\", \"title\": \"\", title, "\", \"year\":", year, ", \"genre
    for genre in genres[:-1]:
        print( "\"\", genre, "\",", end='', sep='')
    print( "\"\", genres[-1], "\",", end = '', sep='')
    print ( " ] }")
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# python example:

```
import csv
import re

csvfile = open('ml-latest-small/movies.csv', 'r')

reader = csv.DictReader( csvfile )
for movie in reader:
    print ("{ \"create\" : { \"_index\": \"movies\", \"_type\": \"movie\", \"_id\" : \"\" , movie['movieId'], \"\" } }\", sep='')
    title = re.sub(" \\.*)$", "", re.sub("'",'', movie['title']))
    year = movie['title'][-5:-1]
    if (not year.isdigit()):
        year = "2016"
    genres = movie['genres'].split('|')
    print ("{ \"id\": \"\", movie['movieId'], "\", \"title\": \"\", title, "\", \"year\":", year, ", \"genre\":[", end='', sep='')
    for genre in genres[:-1]:
        print("\", genre, "\",", end='', sep='')
    print("\", genres[-1], \"\", end = '', sep='')
    print ("] }")
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

┌ importing  
via client api's └



**INNOVATION**  
**SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# | client libraries

free elasticsearch client libraries are available for pretty much any language.

- **java** has a client maintained by elastic.co
- **python** has an elasticsearch package
- elasticsearch-**ruby**
- several choices for **scala**
- elasticsearch.pm module for **perl**

**You don't have to wrangle JSON.**

```
es = elasticsearch.Elasticsearch()
```

```
es.indices.delete(index="ratings", ignore=404)  
deque(helpers.parallel_bulk(es, readRatings(), index="ratings", doc_t  
es.indices.refresh()
```



# python full script

- function to read Movies
  - open csv
  - build dictionary
  - match movieid with title
- function to read Ratings
  - open csv
  - build dictionary
  - match movieid with title
- create ES instance,
- delete ratings table
- load ratings into "ratings" index.

```
import csv
from collections import deque
import elasticsearch
from elasticsearch import helpers
```

```
def readMovies():
    csvfile = open('ml-latest-small/movies.csv', 'r')
    reader = csv.DictReader( csvfile )
    titleLookup = {}

    for movie in reader:
        titleLookup[movie['movieId']] = movie['title']
    return titleLookup
```

```
def readRatings():
    csvfile = open('ml-latest-small/ratings.csv', 'r')

    titleLookup = readMovies()

    reader = csv.DictReader( csvfile )
    for line in reader:
        rating = {}
        rating['user_id'] = int(line['userId'])
        rating['movie_id'] = int(line['movieId'])
        rating['title'] = titleLookup[line['movieId']]
        rating['rating'] = float(line['rating'])
        rating['timestamp'] = int(line['timestamp'])
        yield rating
```

```
es = elasticsearch.Elasticsearch()

es.indices.delete(index="ratings", ignore=404)
deque(helpers.parallel_bulk(es, readRatings(), index="ratings", doc_type="rating"), maxlen=0)
es.indices.refresh()
```



LEARN. EMPOWER. INNOVATE.

RE  
ATE.

# lab: python scripts

- Lab 12: Python scripts to import data
- We will also do the exercise on the next slide.

# python - tags

## exercise

write a script to import the tags.csv data from ml-latest-small into a new “tags” index.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# one solution

- tag = line ['tag'] instead of float line (no longer a number)

```
import csv
from collections import deque
import elasticsearch
from elasticsearch import helpers

def readMovies():
    csvfile = open('ml-latest-small/movies.csv', 'r')
    reader = csv.DictReader( csvfile )
    titleLookup = {}

    for movie in reader:
        titleLookup[movie['movieId']] = movie['title']
    return titleLookup
```

```
def readTags():
    csvfile = open('ml-latest-small/tags.csv', 'r')

    titleLookup = readMovies()

    reader = csv.DictReader( csvfile )
    for line in reader:
        tag = {}
        tag['user_id'] = int(line['userId'])
        tag['movie_id'] = int(line['movieId'])
        tag['title'] = titleLookup[line['movieId']]
        tag['tag'] = line['tag']
        tag['timestamp'] = int(line['timestamp'])
        yield tag
```

```
es = elasticsearch.Elasticsearch()

es.indices.delete(index="tags", ignore=404)
deque(helpers.parallel_bulk(es, readTags(), index="tags", doc_type="tag"), maxlen=0)
es.indices.refresh()
```



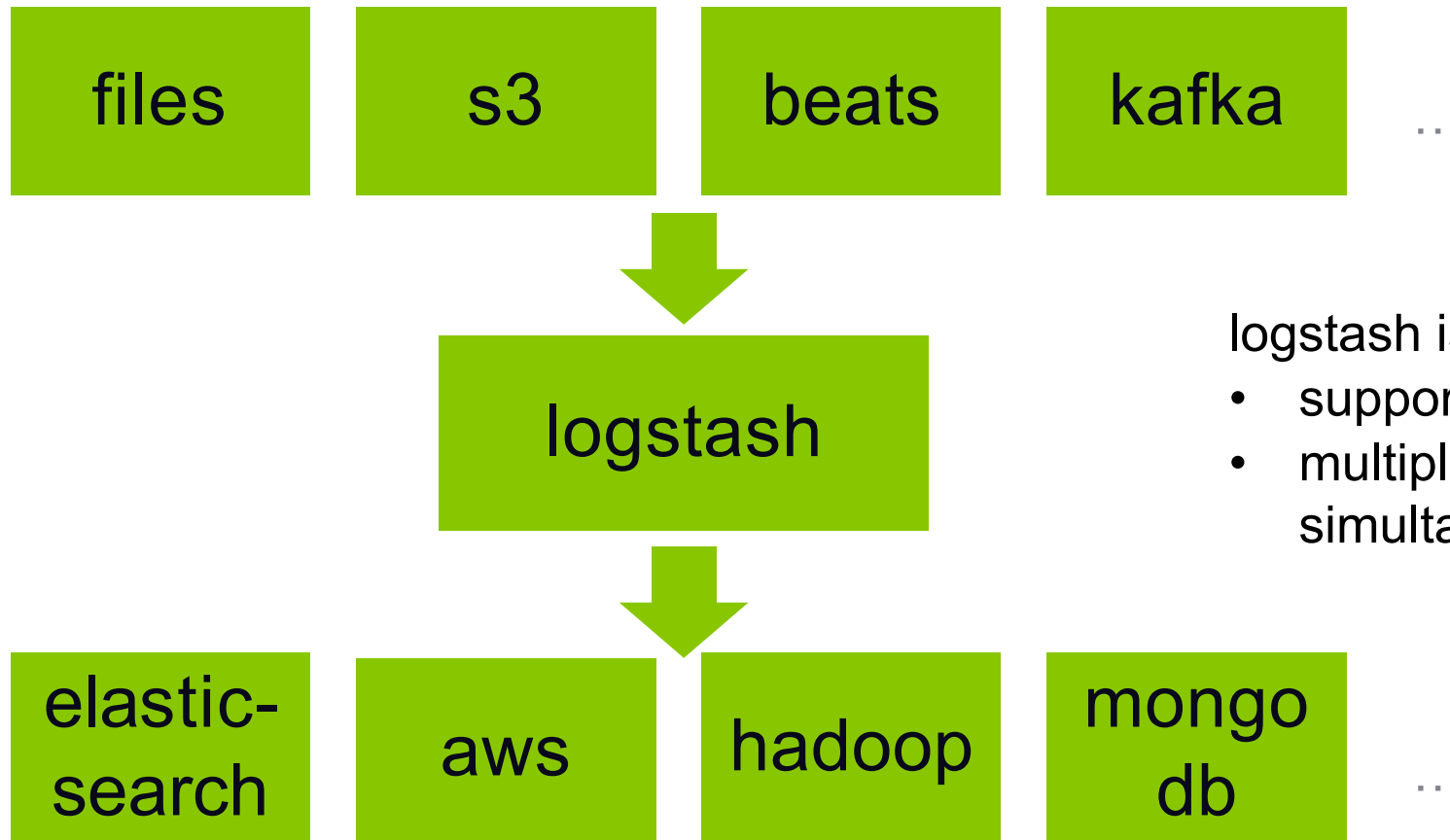
LEARN. EMPower. INNOVATE.

┌ introducing  
logstash ─┐



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# logstash



logstash is very powerful

- supports many systems
- multiple in/out simultaneously.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# Powerful features

- logstash **parses**, **transforms**, and **filters** data as it passes through.
- it can **derive structure** from unstructured data
- it can **anonymize** personal data or exclude it entirely
- it can do **geo-location** lookups
- it can scale across many nodes
- it guarantees at-least-once delivery
- it absorbs throughput from load spikes

See <https://www.elastic.co/guide/en/logstash/current/filter-plugins.html> for the huge list of filter plugins.

- Can serve as buffer between source and destination

# logstash - input sources

elastic beats – cloudwatch – couchdb – drupal – elasticsearch – windows event log – shell output – local files – ganglia – gelf – gemfire – random generator – github – google pubsub – graphite – heartbeats – heroku – http – imap – irc – jdbc – jmx – kafka – lumberjack – meetup – command pipes – puppet – rabbitmq – rackspace cloud queue – redis – relp – rss – s3 – salesforce – snmp – sqlite – sqs – stdin – stomp – syslog – tcp – twitter – udp – unix sockets – varnish log – websocket – wmi – xmpp – zenoss – zeromq



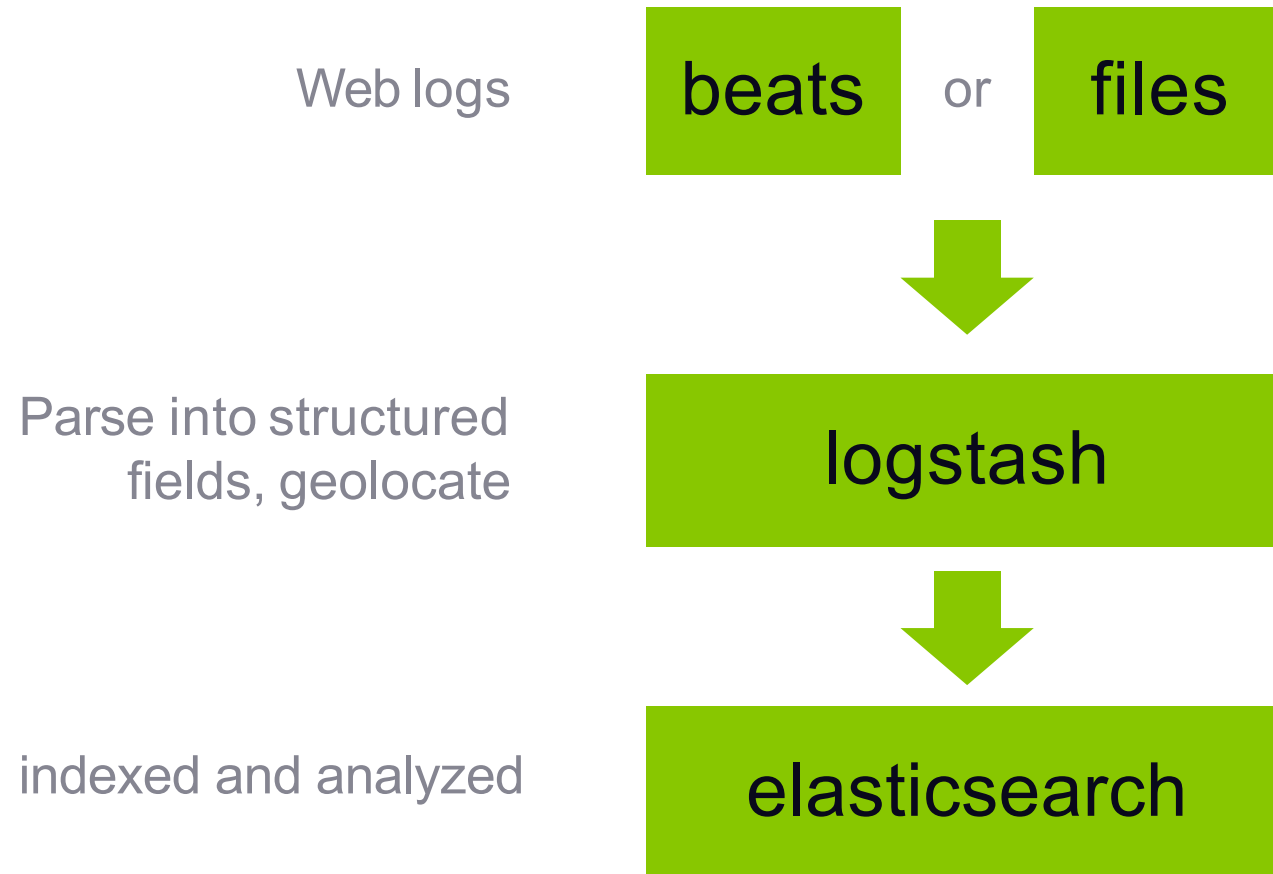
# logstash - output

boundary – circonus – cloudwatch – csv – datadoghq –  
elasticsearch – email – exec – local file – ganglia – gelf –  
bigquery – google cloud storage – graphite – graphtastic –  
hipchat – http – influxdb – irc – jira – juggernaut – kafka –  
librato – loggly – lumberjack – metriccatcher – mongodb –  
nagios – new relic insights – opentsdb – pagerduty – pipe  
to stdin – rabbitmq – rackspace cloud queue – redis –  
redmine – riak – riemann – s3 – sns – solr – sqs – statsd  
– stdout – stomp – syslog – tcp – udp – webhdfs –  
websocket – xmpp – zabbix - zeromq



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# typical usage



┌ installing  
logstash ─┐

# installing logstash

```
sudo apt-get update  
sudo apt-get install logstash
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# configuring logstash

input = log location  
start\_position  
ignore\_older

Set timestamp  
format being used in  
access\_log

Send to  
Elasticsearch &  
standard out

```
sudo vi /etc/logstash/conf.d/logstash.conf
```

```
input {  
  file {  
    path => "/home/<user>/access_log"  start_position => "beginning"  ignore_older => 0  
  }  
}  
  
filter {  
  grok {  
    match => { "message" => "%{COMBINEDAPACHELOG}" }  
  }  
  date {  
    match => [ "timestamp", "dd/MMM/yyyy:HH:mm:ss Z" ]  
  }  
}  
  
output {  
  elasticsearch {  
    hosts => ["localhost:9200"]  
  }  
  stdout {  
    codec => rubydebug  
  }  
}
```



I N

LEARN. EMPOWER. INNOVATE.

running  
logstash

```
cd /usr/share/logstash/
```

```
sudo bin/logstash -f /etc/logstash/conf.d/logstash.conf
```

┌ logstash  
with mysql └



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# jdbc driver

get a mysql connector from <https://dev.mysql.com/downloads/connector/j/>

wget <https://dev.mysql.com/get/Downloads/Connector-J/mysql-connector-java-5.1.42.zip>

unzip mysql-connector-java-5.1.42.zip



# configure logstash

```
input {  
  jdbc {  
    jdbc_connection_string => "jdbc:mysql://localhost:3306/movielens"  
    jdbc_user => "root"  
    jdbc_password => "password"  
    jdbc_driver_library => "/home/<user>/mysql-connector-java-5.1.46/mysql-connector-java-5.1.46-bin.jar"  
    jdbc_driver_class => "com.mysql.jdbc.Driver"  
    statement => "SELECT * FROM movies"  
  }  
}
```

```
output {  
  stdout { codec => json_lines }  
  elasticsearch {  
    "hosts" => "localhost:9200"  
    "index" => "movielens-sql"  
    "document_type" => "data"  
  }  
}
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# lab: logstash & mysql

- Lab 13: Install Logstash and integrate with MySQL

┌ logstash  
with s3 └



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# what is s3

amazon web services' **simple storage service**

cloud-based distributed storage system



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

| integration is  
easy

```
input {  
  s3{  
    bucket => "learning-es"  
    access_key_id => "AKIAIS****C26Y***Q"  
    secret_access_key => "d*****FEN0XcCuNC4iTbSLbibA*****eyn****"  
  }  
}
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

┌ logstash  
with kafka └



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# what is kafka

- apache kafka
- open-source stream processing platform
- high throughput, low latency
- publish/subscribe
- process streams
- store streams

has a lot in common with logstash, really.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# integration is easy

```
input {  
  kafka{  
    bootstrap_servers => "127.0.0.1:9200"  
    topics = ["kafka-logs"]  
  }  
}
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

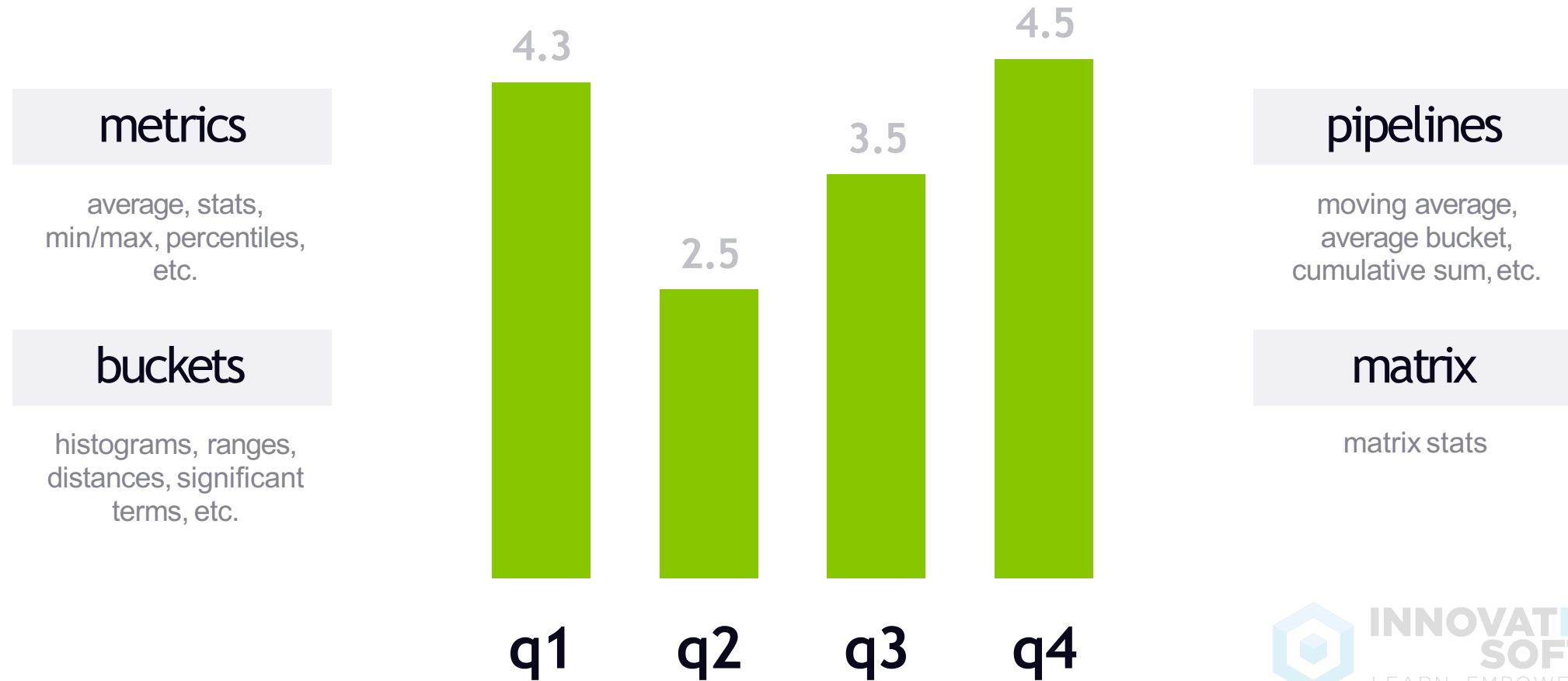


「aggregations」



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# not just for search anymore



aggregations  
are amazing

elasticsearch aggregations can  
sometimes take the place of hadoop /  
spark / etc – and return results instantly!



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

it gets better

you can even nest aggregations  
together!

Pair aggregations with search queries

# let's learn by example

bucket by rating value:

```
curl -XGET  
'127.0.0.1:9200/ratings/rating/_search?size=0&pretty' -d '  
{  
  "aggs": {  
    "ratings": {  
      "terms": {  
        "field": "rating"  
      }  
    }  
  }  
}'
```

Sum up all documents that have a "rating" value.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# let's learn by example

- query narrows down movies to only 5.0 rating
- aggregates sum up all those values.

count only 5-star ratings:

```
curl -XGET  
'127.0.0.1:9200/ratings/rating/_search?size=0&pretty' -d '  
{  
  "query": {  
    "match": {  
      "rating": 5.0  
    }  
  },  
  "aggs" : {  
    "ratings": {  
      "terms": {  
        "field" : "rating"  
      }  
    }  
  }  
}'
```



INNOVATION  
SOFTWARE

LEARN. EMPOWER. INNOVATE.

# let's learn by example

average rating for Star Wars:

- Find Star Wars Episode IV
- Return average rating for that specific movie.

```
curl -XGET
'127.0.0.1:9200/ratings/rating/_search?size=0&pretty' -d '{
  "query": {
    "match_phrase": {
      "title": "Star Wars Episode IV"
    }
  },
  "aggs" : {
    "avg_rating": {
      "avg": {
        "field" : "rating"
      }
    }
  }
}'
```



**ATION**  
**SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

┌ histograms ┐

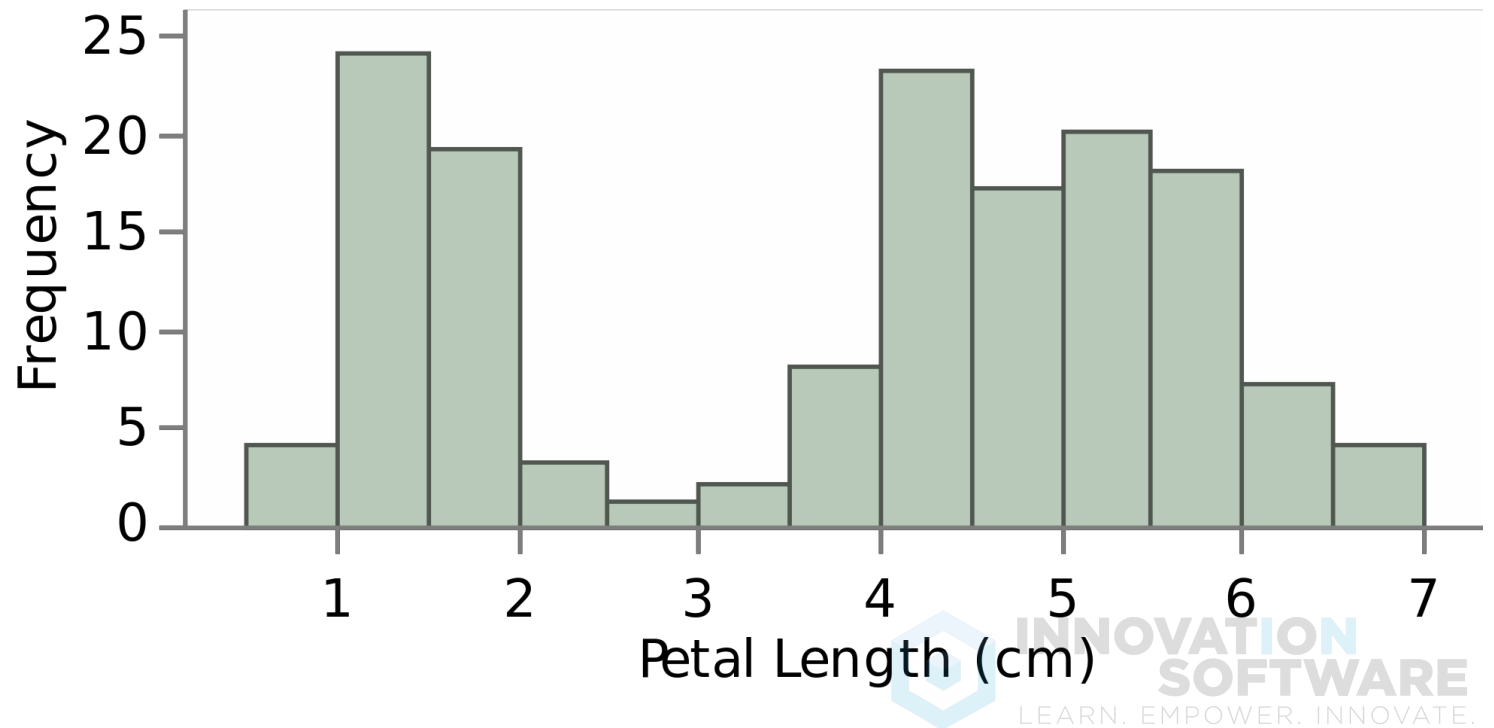


**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.



# what is a histogram

display totals of documents  
bucketed by  
some **interval range**



# rating intervals

Group ratings together by whole number.

```
curl -XGET
'127.0.0.1:9200/ratings/rating/_search?size=0&pretty' -d '{
  {
    "aggs" : {
      "whole_ratings": {
        "histogram": {
          "field": "rating",
          "interval": 1.0
        }
      }
    }
  }
}'
```



**INNOVARE**  
LEARN. EMPOWER. INNOVATE.

# sort by decade

```
curl -XGET  
'127.0.0.1:9200/movies/movie/_search?size=0&pretty' -d '  
{  
  "aggs" : {  
    "release": {  
      "histogram": {  
        "field": "year",  
        "interval": 10  
      }  
    }  
  }  
}
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

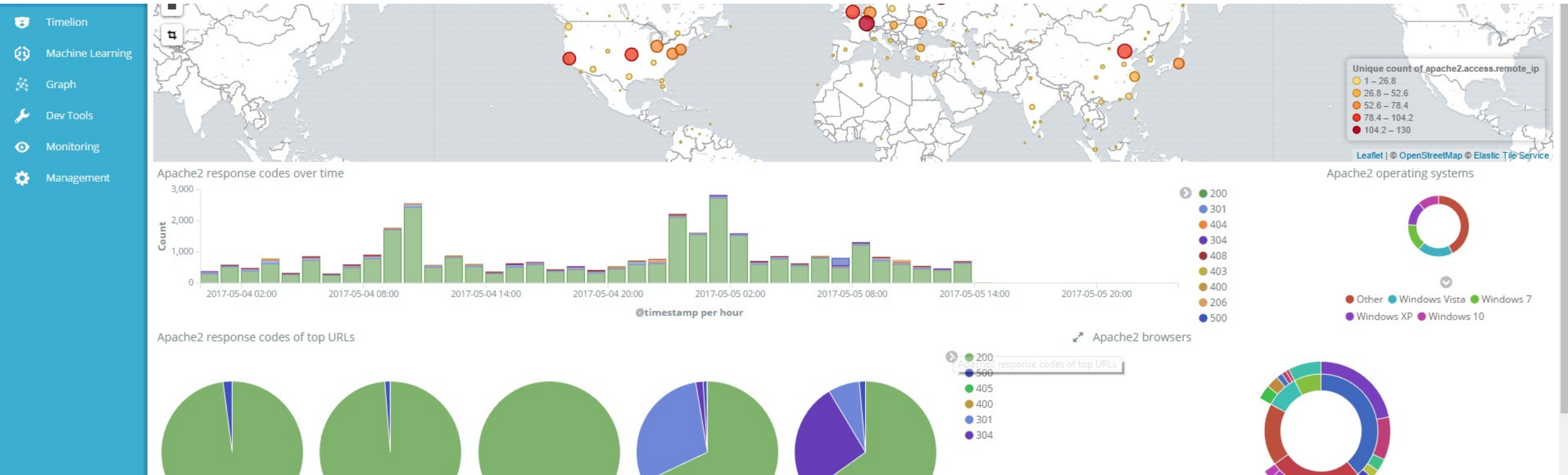


using  
kibana



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# what is kibana



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# installing kibana

```
sudo apt-get install kibana  
sudo vi /etc/kibana/kibana.yml  
change server.host to 0.0.0.0
```

```
sudo /bin/systemctl daemon-reload  
sudo /bin/systemctl enable kibana.service  
sudo /bin/systemctl start kibana.service
```

kibana is now available on port 5601

┌ playing with  
kibana ┐



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.



let's analyze the works  
of **william shakespeare...**

because we can.





# lab: aggs & kibana

- Lab 14: Using aggregates and histograms
- Lab 15: Installing and using Kibana

# | elasticsearch

## ┌ exercise └

find the longest shakespeare plays -  
create a vertical bar chart that  
aggregates the count of documents by  
play name in descending order.



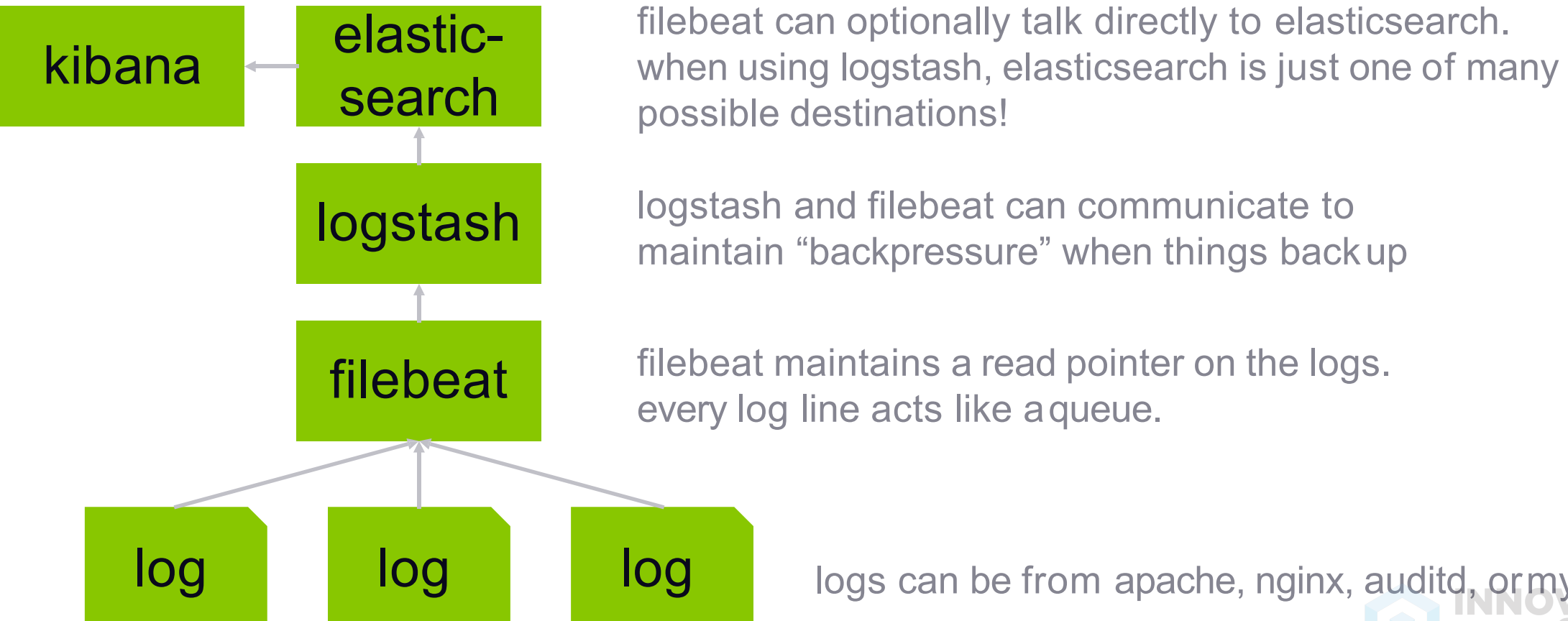
**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

using  
filebeat



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# filebeat: lightweight shipper



# this is called the elastic stack

prior to beats, you'd hear about the “ELK stack” –  
elasticsearch, logstash, kibana.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# filebeat vs logstash

- Wrong question: Use them together.
- it won't let you overload your pipeline.
- you get more flexibility on scaling your cluster.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

┌ installing  
filebeat └

# lab: kibana & filebeats

- Lab 16: Install filebeats and configure with Kibana



# elasticsearch operations



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

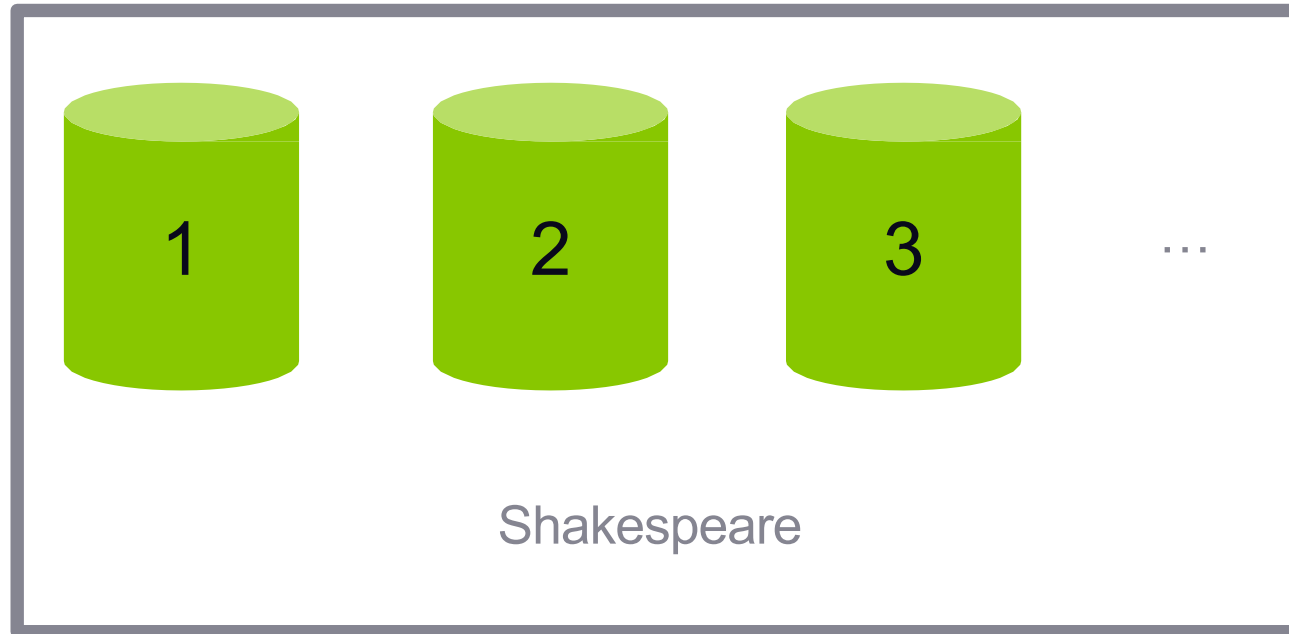
「choosing your  
shards」



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# an index is split into shards

Documents are **hashed** to a particular **shard**.



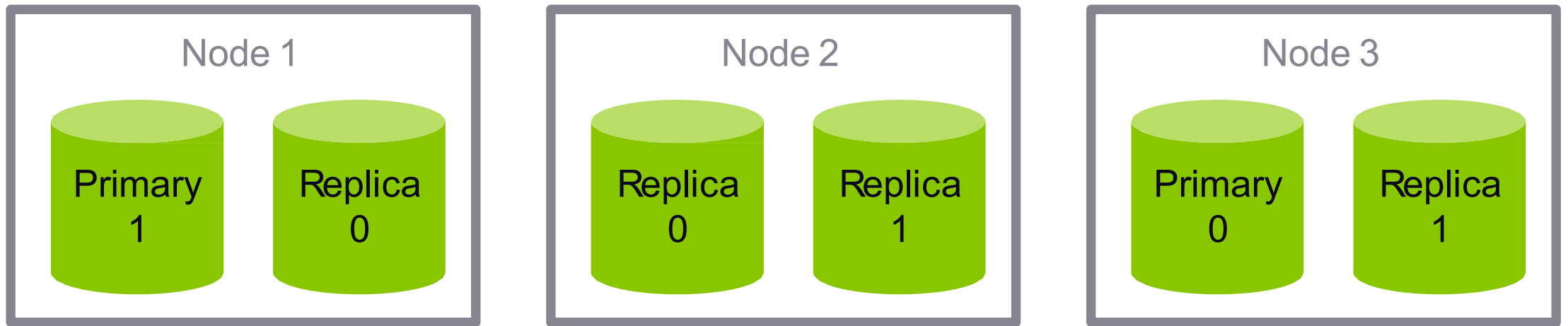
Each shard may be on a different **node** in a **cluster**.  
Every shard is a self-contained Lucene index of its own.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# primary and replica shards

This **index** has two **primary shards** and two **replicas**.  
Your application should round-robin requests amongst nodes.



**Write** requests are routed to the primary shard, then replicated  
**Read** requests are routed to the primary or any replica



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# how many shards do i need?

- you can't add more shards later without re-indexing
- but shards aren't free – you can't just make 1,000 of them and stick them on one node at first.
- you want to overallocate, but not too much
- consider scaling out in phases, so you have time to re-index before you hit the next phase



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# shard planning

- the “right” number of shards depends on your data and your application. there’s no secret formula.
- start with a single server using the same hardware you use in production, with one shard and no replication.
- fill it with real documents and hit it with real queries.
- push it until it breaks – now you know the capacity of a single shard.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# remember replica shards can be added

- read-heavy applications can add more replica shards without re-indexing.
- note this only helps if you put the new replicas on extra hardware!
- Shards have a limit of 2B documents



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

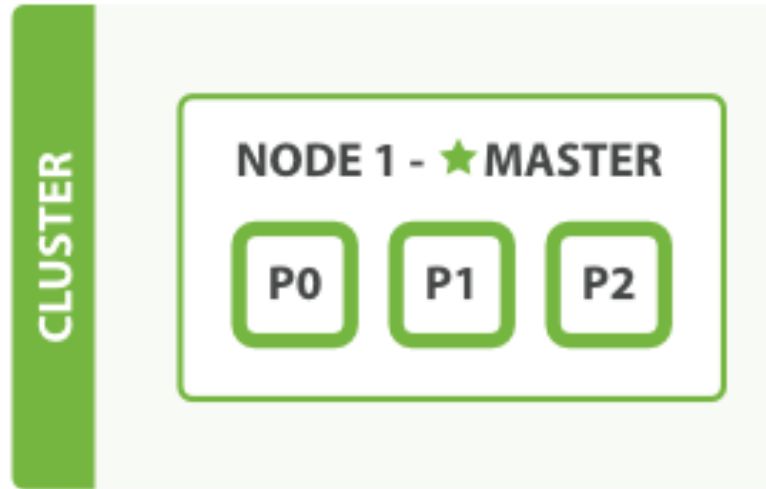
「multi-node cluster」



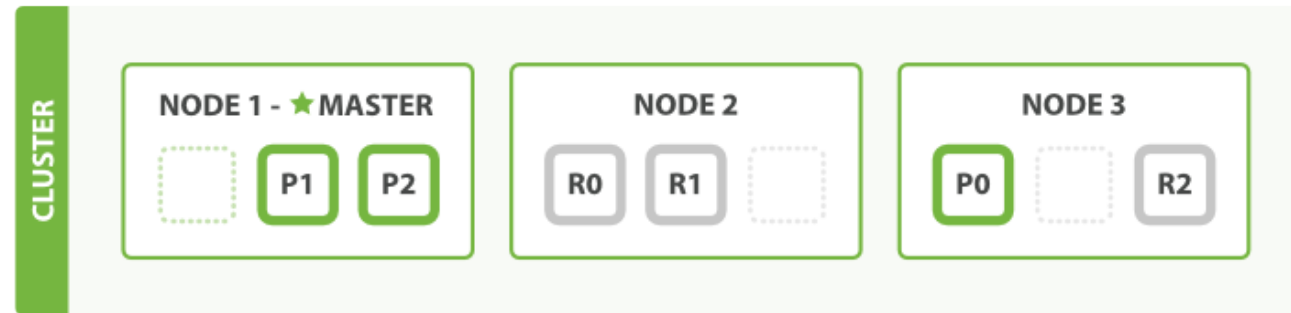
**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.



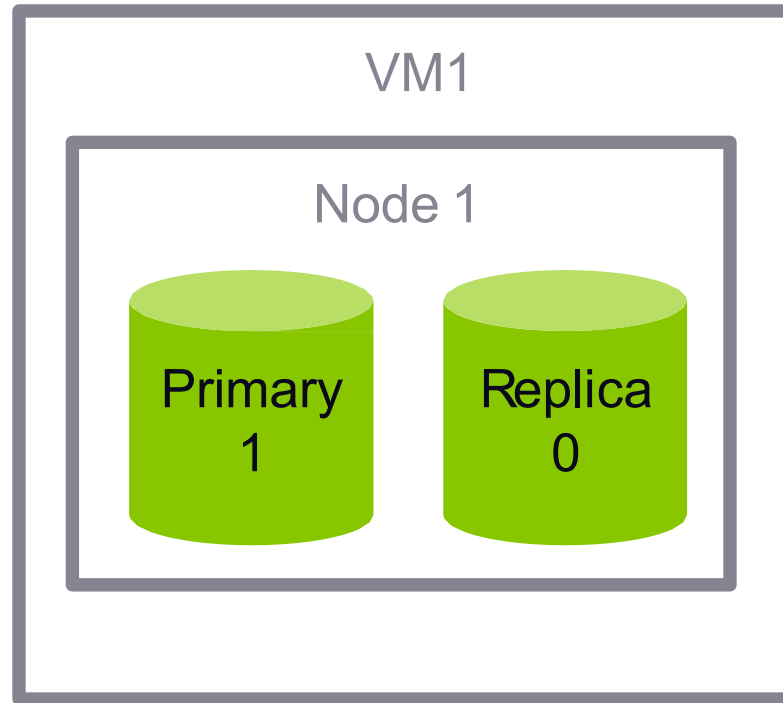
# elasticsearch multi-node cluster



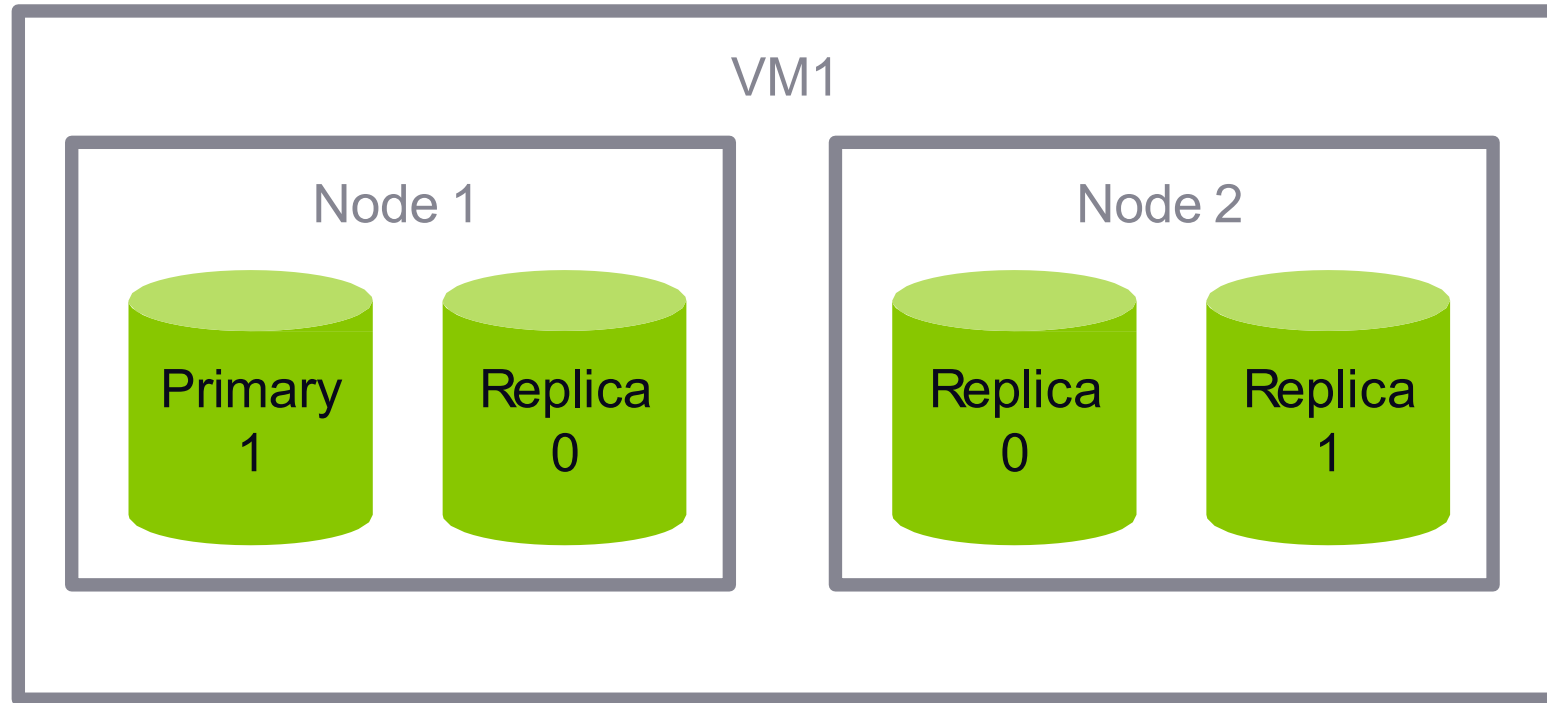
# elasticsearch multi-node cluster



# elasticsearch multi-node cluster



# elasticsearch multi-node cluster



# elasticsearch multi-node cluster

- Set cluster.name
- Set node.name
- Add
  - node.max\_local\_storage\_nodes: 2
- Copy config directory for new node
- Edit new node yml and set
  - cluster.name
  - network.port (something unique)
- Create log directory
- Create new init script

# lab: multi-node cluster

- Lab 17: Multi-node cluster

「adding an index」



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# creating a new index

- Remember this is an important step for production clusters!

```
PUT /new_index
{
  "settings": {
    "number_of_shards": 10,
    "number_of_replicas": 1
  }
}
```

You can use *index templates* to automatically apply mappings, analyzers, aliases, etc.



# scaling strategies

- Adding a new shard is not only option for scaling.
- multiple-indices:
  - make a new index to hold new data
  - search both indices
  - use *index aliases* to make this easy to do
  - Much easier than re-indexing

# scaling strategies

- multiple-indices:
  - with time-based data, you can have one index per time frame
  - common strategy for log data where you usually just want current data, but don't want to delete old data either
  - again you can use index aliases, ie "logs\_current", "last\_3\_months", to point to specific indices as they rotate



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# alias rotation example

POST /\_aliases

```
{
  "actions": [
    { "add": { "alias": "logs_current", "index": "logs_2017_06" }},
    { "remove": { "alias": "logs_current", "index": "logs_2017_05" }},
    { "add": { "alias": "logs_last_3_months", "index": "logs_2017_06" }},
    { "remove": { "alias": "logs_last_3_months", "index": "logs_2017_03" }}
  ]
}
```

optionally....

DELETE /logs\_2017\_03

「choosing your  
hardware」



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

RAM is likely your bottleneck

64GB per machine is the sweet spot  
(32GB to elasticsearch, 32GB to the  
OS/ disk cache for lucene)

under 8GB not recommended



# other hardware considerations

- fast disks are better – SSD's if possible
- use RAID0 – your cluster is already redundant
- cpu not that important
- need a fast network
- don't use NAS
- use medium to large configurations; too big is bad, and too many small boxes is bad too.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

┌ heap sizing ┐



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# your heap size is wrong

the default heap size is only 1GB!

half or less of your physical memory should be allocated to elasticsearch

- the other half can be used by lucene for caching
- if you're not aggregating on analyzed string fields, consider using less than half for elasticsearch
- smaller heaps result in faster garbage collection and more memory for caching

```
export ES_HEAP_SIZE=10g
```

or

```
ES_JAVA_OPTS="-Xms10g -Xmx10g" ./bin/elasticsearch
```

don't cross 32GB! pointers blow up then.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.



「monitoring with x-  
pack」



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# what is x-pack?

- an elastic stack extension
- security, monitoring, alerting, reporting, graph, and machine learning
- formerly shield / watcher / marvel
- only parts can be had for free – requires a paid license or trial otherwise

# install x-pack

```
cd /usr/share/elasticsearch
sudo bin/elasticsearch-plugin install x-pack

sudo vi /etc/elasticsearch/elasticsearch.yml
(Add xpack.security.enabled:false)

sudo /bin/systemctl stop elasticsearch.service

sudo /bin/systemctl start elasticsearch.service

cd /usr/share/kibana/
sudo -u kibana bin/kibana-plugin install x-pack
sudo /bin/systemctl stop kibana.service

sudo /bin/systemctl start kibana.service
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# lab: Install X-Packs

- Lab 17: Install X-Packs

using  
snapshots



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# back up your indices

store backups to NAS, Amazon S3, HDFS, Azure

smart enough to only store changes since last snapshot



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# create a repository

add it into elasticsearch.yml:

```
path.repo: ["/home/<user>/backups"]
```

```
PUT _snapshot/backup-repo
```

```
{  
  "type": "fs",  
  "settings": {  
    "location": "/home/<user>/backups/backup-repo"  
  }  
}
```



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# using snapshots

snapshot all open indices:

PUT\_snapshot/backup-repo/snapshot-1

get information about a snapshot:

GET\_snapshot/backup-repo/snapshot-1

monitor snapshot progress:

GET\_snapshot/backup-repo/snapshot-1/\_status

restore a snapshot of all indices:

POST/\_all/\_close

POST\_snapshot/backup-repo/snapshot-1/\_restore



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.



rolling  
restarts



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# restarting your cluster



sometimes you have to... OS updates, elasticsearch version updates, etc.

to make this go quickly and smoothly, you want to disable index reallocation while doing this.

# rolling restart procedure

1. stop indexing new data if possible
2. disable shard allocation
3. shut down one node
4. perform your maintenance on it and restart, confirm it joins the cluster.
5. re-enable shard allocation
6. wait for the cluster to return to green status
7. repeat steps 2-6 for all other nodes
8. resume indexing new data

# | cheat sheet

```
PUT _cluster/settings
{
  "transient": {
    "cluster.routing.allocation.enable": "none"
  }
}
```

Disable shard allocation

```
sudo /bin/systemctl stop elasticsearch.service
```

Stop elasticsearch safely

```
PUT _cluster/settings
{
  "transient": {
    "cluster.routing.allocation.enable": "all"
  }
}
```

Enable shard allocation

# cluster basics

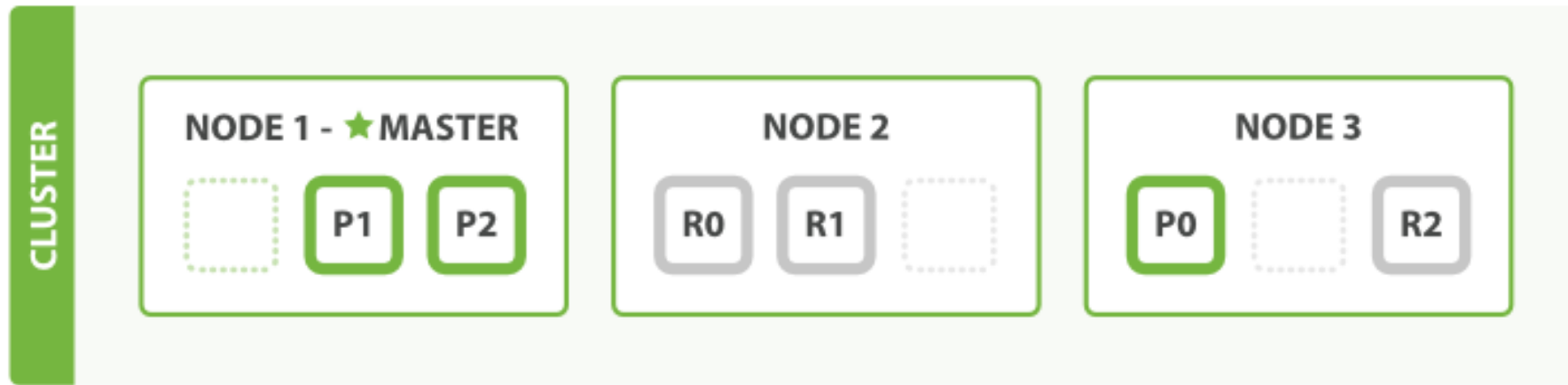
# elasticsearch cluster

- a cluster is a collection of one or more nodes (servers)
- together the cluster holds your entire data and provides federated indexing and search capabilities across all nodes
- clusters are identified by a unique name
- default cluster naming is "elasticsearch"
- nodes join clusters automatically on startup
- nodes join clusters using their configured naming



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# elasticsearch cluster



# putting it all together

- when you start an instance of Elasticsearch, you are starting a node
- an Elasticsearch cluster is a group of nodes that have the same **cluster.name** attribute
- as nodes join or leave a cluster, the cluster automatically reorganizes itself to evenly distribute the data across the available nodes
- nodes can host primary and/or replica depending on configuration
- configuration is contained in the `elasticsearch.yml`, a YAML file
- nodes are both data and master-eligible, by default
- new nodes can be configured for a specific purpose, such as handling data ingestion





# cluster monitoring

- Marvel, an Elastic Stack Feature, is the default option for monitoring Elasticsearch clusters
- monitoring should be done with a separate node or cluster to avoid impacts on monitoring when there are cluster issues
- InfluxDB can be utilized with an Elasticsearch Telegraf agent for monitoring
- three Elastic Stack APIs form the basis for monitoring
  - node stats API for cluster node statistics.
  - cluster health API retrieves status of cluster node health
  - cluster stats API retrieves status from a cluster wide perspective

「 features 」

# Elastic Stack Features

## Security



Protection of your Elasticsearch data in a robust and granular way.

## Alerting



Receive notifications about changes in your data.

## Monitoring



Ensure your clusters are running at peak efficiency and effectiveness

## Reporting



Create and share reports of your Kibana charts and dashboards.

## Graph



Explore meaningful relationships in your data streams

## Machine Learning



Automation of anomaly detection on your Elasticsearch data.



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

┌ query  
performance └



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.

# to optimize performance of our queries

- design your data model based on your queries
- utilize nested documents versus parent-child
- tune the refresh interval, which as we noted earlier is 1s
- use force merge to reduce the number of segments
- understand the global ordinals where you use term aggregation
- consider pre-indexing when using range aggregations
- avoid memory to disk swapping
- always test against production-like data sets
- use bulk for insert queries when possible

「wrapping up」



**INNOVATION  
SOFTWARE**  
LEARN. EMPOWER. INNOVATE.