



elasticsearch.

SEARCH MADE EASY FOR (WEB) DEVELOPERS

Alexander Reelsen
alexander@reelsen.net
@spinscale



AGENDA

- What is so important about search?
- Scalability, Sharding & Replication
- Configuration, Mapping & Analyzers
- Querying, Facetting, Percolation
- Modules, Plugins, Rivers & Tools
- Production setup & living in the trenches



ABOUT ME - ALEXANDER REELSEN

- Studied information systems
- 10 years linux system engineering, converted to software engineering
- Web framework enthusiast, fed up with complex java environment for simple webapps
- Other interests: Scaling web architectures, Web 2.0 (nosql, search)
- Author of [Play framework cookbook](#)
- Working at [Lusini GmbH](#), building a b2b ecommerce platform
- Streetball/Basketball



WHAT IS SO IMPORTANT ABOUT SEARCH?

- No search, no google, no bing, no twitter, no amazon, no ebay, ...
- Functional requirements: Relevance (finds the *right* stuff)
- Non-functional requirements: Scalability, performance, concurrent updates
- Solutions: [Google commerce search](#), [Sphinx](#), [SearchBlox](#), [Solr](#), [elasticsearch](#), [IndexTank](#), [Sensei DB](#)



SEARCH MUST SEARCH FOR IDS

Lusini.de Alles für professionelle Gastgeber
Gastronomie | Hotel | Gewerbe | Catering

JBEDMCYH94VB

Warenkorb (0 Artikel)

Küche & Kühlraum | Theke & Gastraum | Buffet & Bankett | Berufsbekleidung | Zimmer & Rezeption | Bad & WC | Garten & Terrasse | Mehr

Startseite > Suche > "JBEDMCYH94VB"

"JBEDMCYH94VB" bei Lusini (1 Artikel)

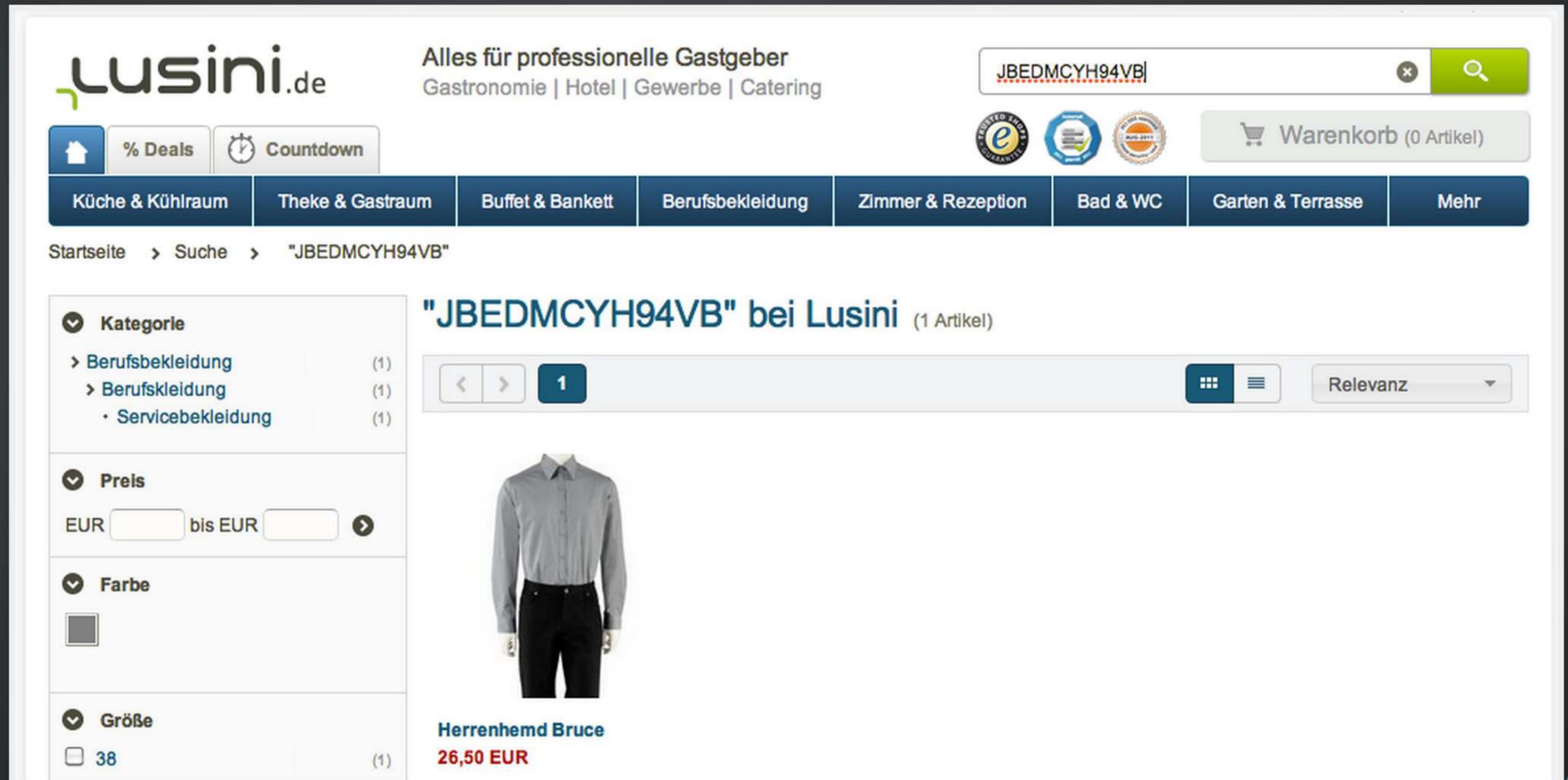
Kategorie
Berufsbekleidung (1)
Berufskleidung (1)
Servicebekleidung (1)

Preis
EUR bis EUR

Farbe

Größe
38 (1)

Herrenhemd Bruce
26,50 EUR



SEARCH MUST SEARCH FOR COLORS

lusini.de Alles für professionelle Gastgeber
Gastronomie | Hotel | Gewerbe | Catering

teelichthalter rot

Warenkorb (0 Artikel)

Startseite > Suche > "teelichthalter rot"

"teelichthalter rot" bei Lusini (17 Artikel)

Relevanz ▾

 <p>Teelichthalter 14,99 EUR</p>	 <p>Teelichthalter 5,99 EUR</p>	 <p>Teelichthalter 3,99 EUR</p>	 <p>Teelichthalter 2,79 EUR</p>
--	---	---	---

Kategorie

- > Theke & Gastraum (14)
- > Raum-Deko (14)
 - > Dekomaterial (3)
 - > Beleuchtung (7)
 - Sonstige Raumdekoration (4)
 - > Zimmer & Rezeption (3)

Preis

EUR bis EUR

Farbe



SEARCH MUST SEARCH FOR BRANDS

Lusini.de Alles für professionelle Gastgeber
Gastronomie | Hotel | Gewerbe | Catering

ritzenhoff glas

Startseite > Suche > "ritzenhoff glas"

"ritzenhoff glas" bei Lusini (625 Artikel)

... Relevanz ▾

 6709 623628 625635	 4791 117318	 4791 117318	 6680 617184 617191
Glasteller Balzano - 8,36 EUR	Glasschüssel Point - 6 5,84 EUR	Glasschüssel Point - 6 15,92 EUR	Glasleuchter Fiona - 2,48 EUR

Kategorie

- > Küche & Kühlraum (18)
- > Theke & Gastraum (471)
 - > Geschirr (117)
 - > Gläser, Becher & Karaffen (255)
 - > Raum-Deko (99)
- > Buffet & Bankett (83)
- > Zimmer & Rezeption (52)
- > Essen & Trinken (1)

Preis

EUR bis EUR

Farbe



SEARCH MUST ADVICE



The screenshot shows the top navigation bar of the website Lusini.de. On the left is the logo 'Lusini.de'. To its right is the tagline 'Alles für professionelle Gastgeber' and the categories 'Gastronomie | Hotel | Gewerbe | Catering'. Below the logo are three buttons: a home icon, '% Deals', and 'Countdown'. A horizontal menu contains five categories: 'Küche & Kühlraum', 'Theke & Gastraum', 'Buffet & Bankett', 'Berufsbekleidung', and 'Zimmer & Rezeption'. A search bar on the right contains the text 'glasvase' and has a magnifying glass icon. A dropdown menu below the search bar lists suggestions: 'glasvase', 'glasvase', 'glaswindlicht', 'glasplatte', and 'glasteelichthalter'. To the right of the search bar is a button labeled 'Artikel)' and a 'Mehr' button. At the bottom left of the header, there is a breadcrumb trail: 'Startseite > Küche & Kühlraum'. A teal four-pointed star icon is located in the bottom right corner of the overall image.

Lusini.de Alles für professionelle Gastgeber
Gastronomie | Hotel | Gewerbe | Catering

Home % Deals Countdown

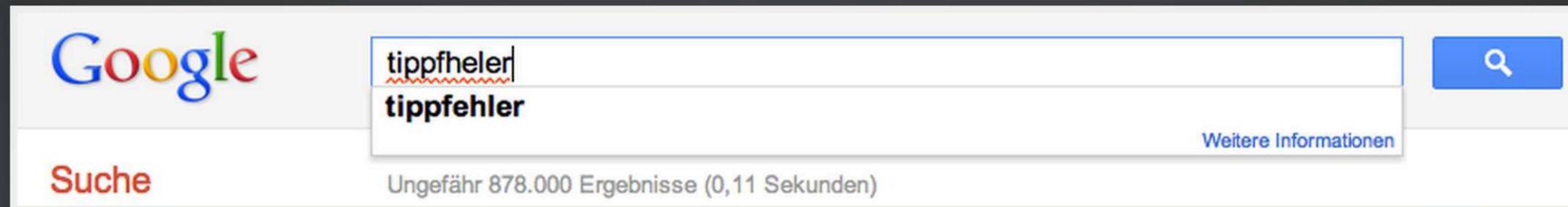
Küche & Kühlraum Theke & Gastraum Buffet & Bankett Berufsbekleidung Zimmer & Rezeption

Startseite > Küche & Kühlraum

glasvase
glasvase
glaswindlicht
glasplatte
glasteelichthalter

Artikel) Mehr

SEARCH MUST BE INTELLIGENT



The image shows a Google search bar with the text "tippfheler" entered. A red squiggly line is under the "h" in "tippfheler", and a dropdown menu shows the corrected text "tippfehler". To the right of the search bar is a blue search button with a magnifying glass icon. Below the search bar, the word "Suche" is written in red. To the right of "Suche", the text "Ungefähr 878.000 Ergebnisse (0,11 Sekunden)" is displayed. Further to the right, there is a link that says "Weitere Informationen".

Google

tippfheler
tippfehler

Weitere Informationen

Suche

Ungefähr 878.000 Ergebnisse (0,11 Sekunden)



SEARCH MUST AGGREGATE

lusini.de Alles für professionelle Gastgeber
Gastronomie | Hotel | Gewerbe | Catering

herren hemd |

% Deals Countdown

Warenkorb (0 Artikel)

Küche & Kühlraum | Theke & Gastraum | Buffet & Bankett | Berufsbekleidung | Zimmer & Rezeption | Bad & WC | Garten & Terrasse | Mehr

Startseite > Suche > "herren hemd"

Kategorie

- > Berufsbekleidung (77)
- > Berufskleidung (77)
 - > Servicebekleidung (77)
 - Hemden & Blusen (77)

Preis

EUR bis EUR

Farbe

Größe

- 44 (46)
- 46 (44)
- 45 (42)
- 43 (42)
- 41 (42)
- 40 (42)
- 38 (42)
- 47 (41)
- 42 (41)
- 39 (41)

Mehr...

Material

- Mischgewebe (33)
- Baumwolle (31)
- Nicht angegeben (16)

Ausführung

- Herrenhemd (2)
- Nicht angegeben (75)

Serie

- Herrenhemd Rico (8)
- Herrenhemd Keno (7)
- Herrenhemd Marcello (6)
- Herrenhemd (6)
- Herrenhemd Tori (4)
- Herrenhemd Daniel (4)
- Herrenhemd Veith (3)
- Herrenhemd Tyler (3)
- Herrenhemd Marc (3)
- Herrenhemd Lius (3)

Mehr...



WHY AN OWN SEARCH ENGINE?

- Because you can - telling this your CTO doesn't work.
- Your data, your search - noone spying...
- Customize your search - Rank your own style
- Customize your data - Extend your search?!
- Best support and in-sourced know-how - Lower TCO
- No blackbox - Lower TCO



ELASTICSEARCH IN TEN SECONDS

- Java, based on Apache Lucene
- Scales out, replicates, shards, fail-over
- Schema-free
- Document-based
- Every interaction can be done via HTTP & JSON
- References: Mozilla, StumbleUpon, Sony, Infochimps, Assistly, Klout

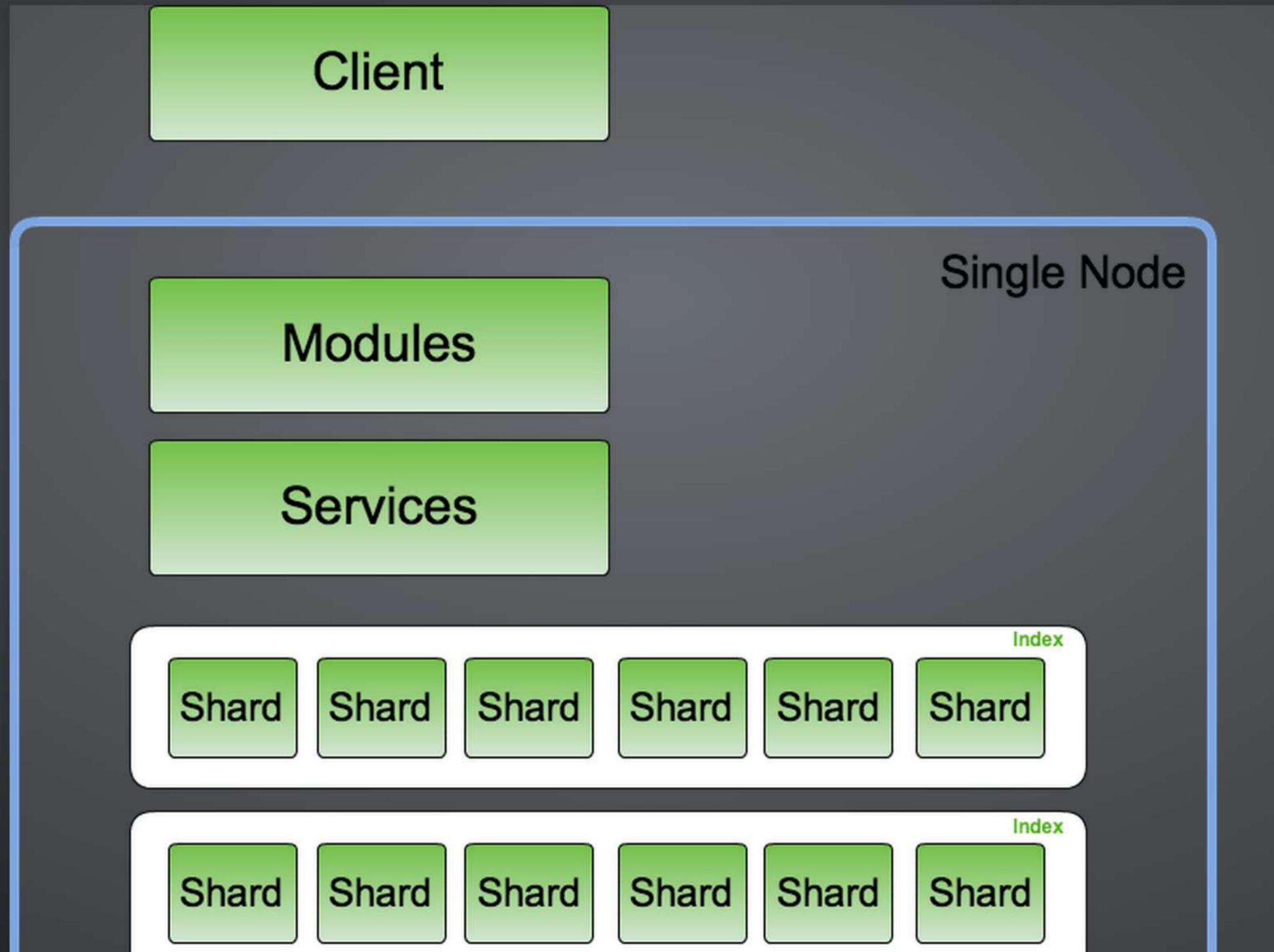


STANDING ON THE SHOULDERS OF GIANTS

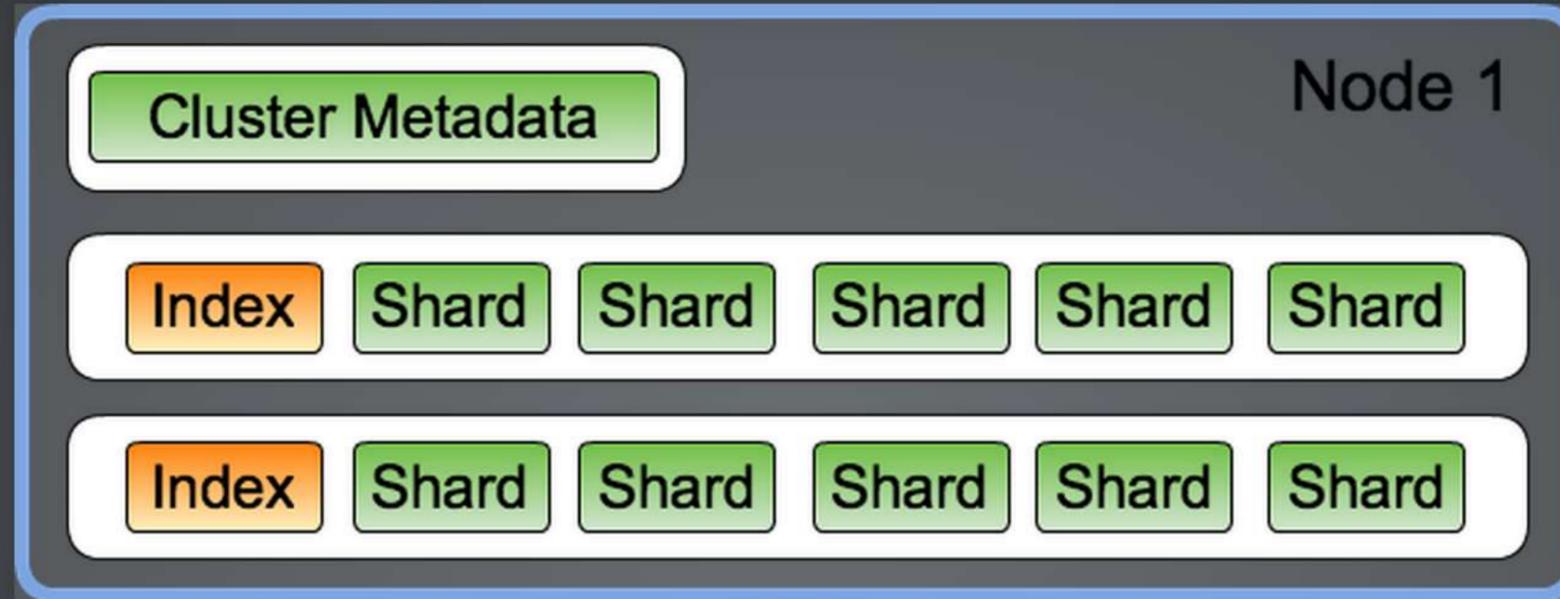
- Lucene, JBoss Netty, Jackson, log4j
- Google Guice, Google Guava, MVEL, Groovy
- Jodatime, JLine, snakeyaml
- hamcrest, testng
- sigar via JNA



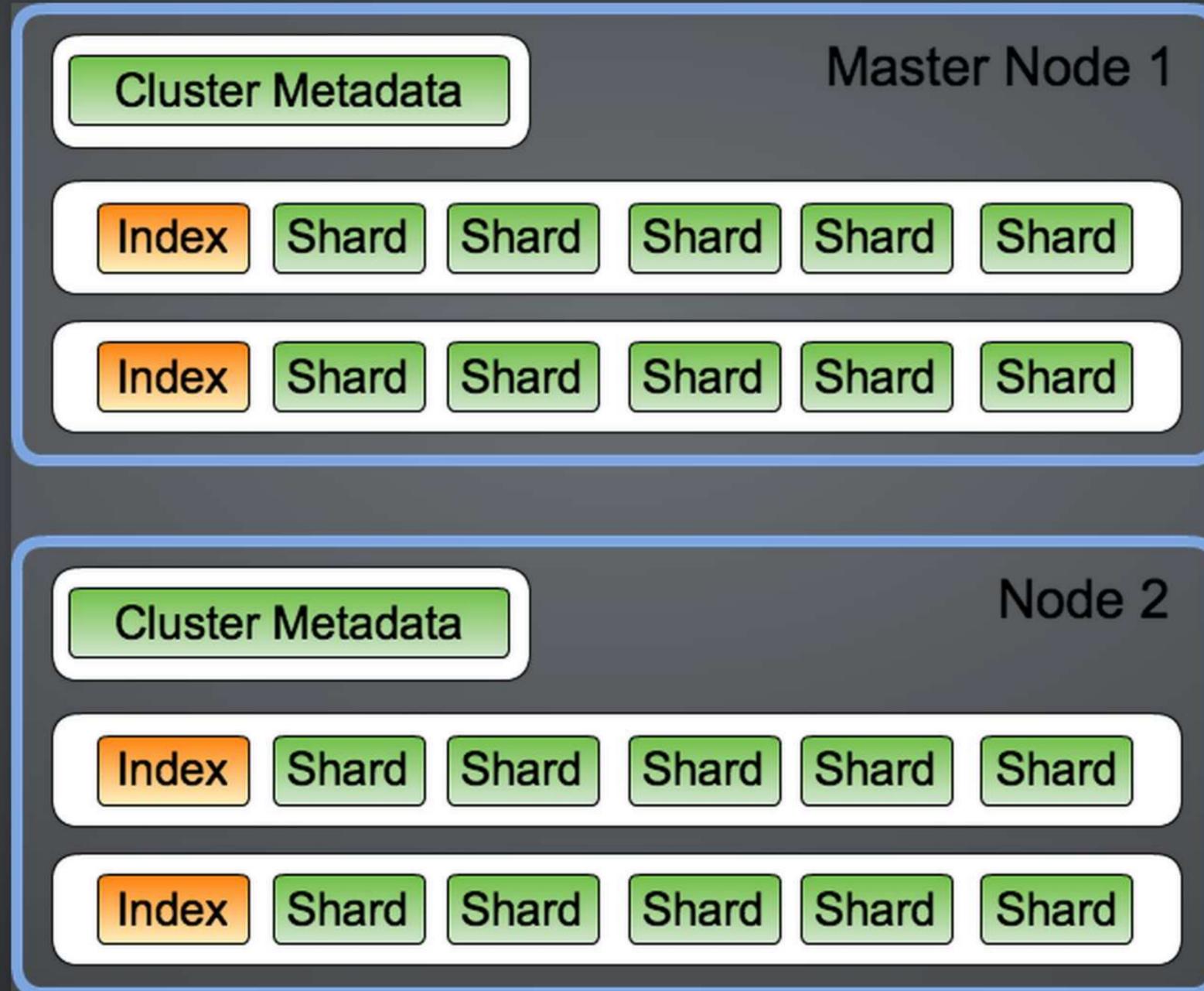
ELASTICSEARCH ARCHITECTURE



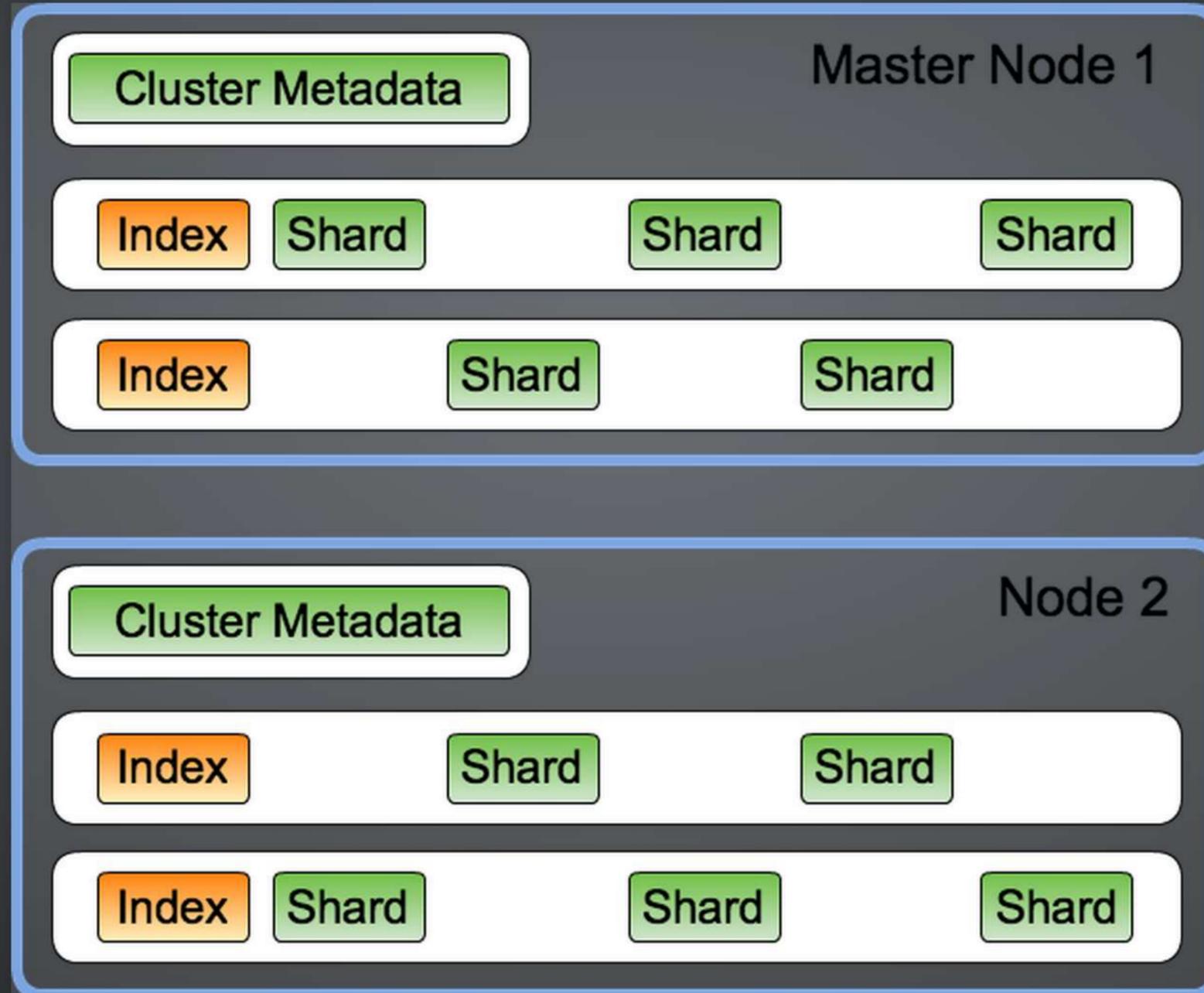
SINGLE NODE SETUP



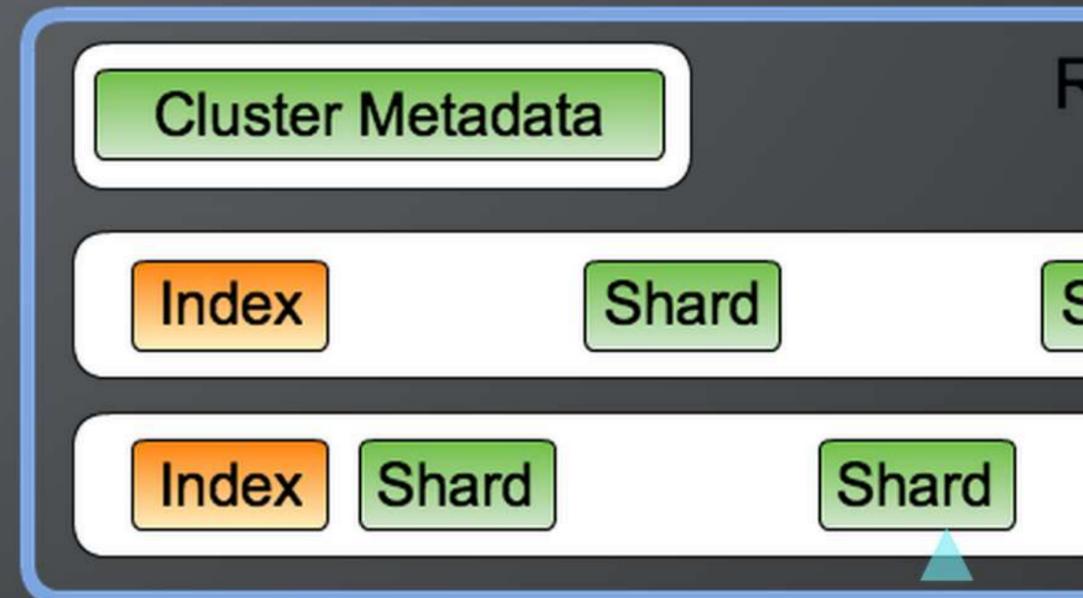
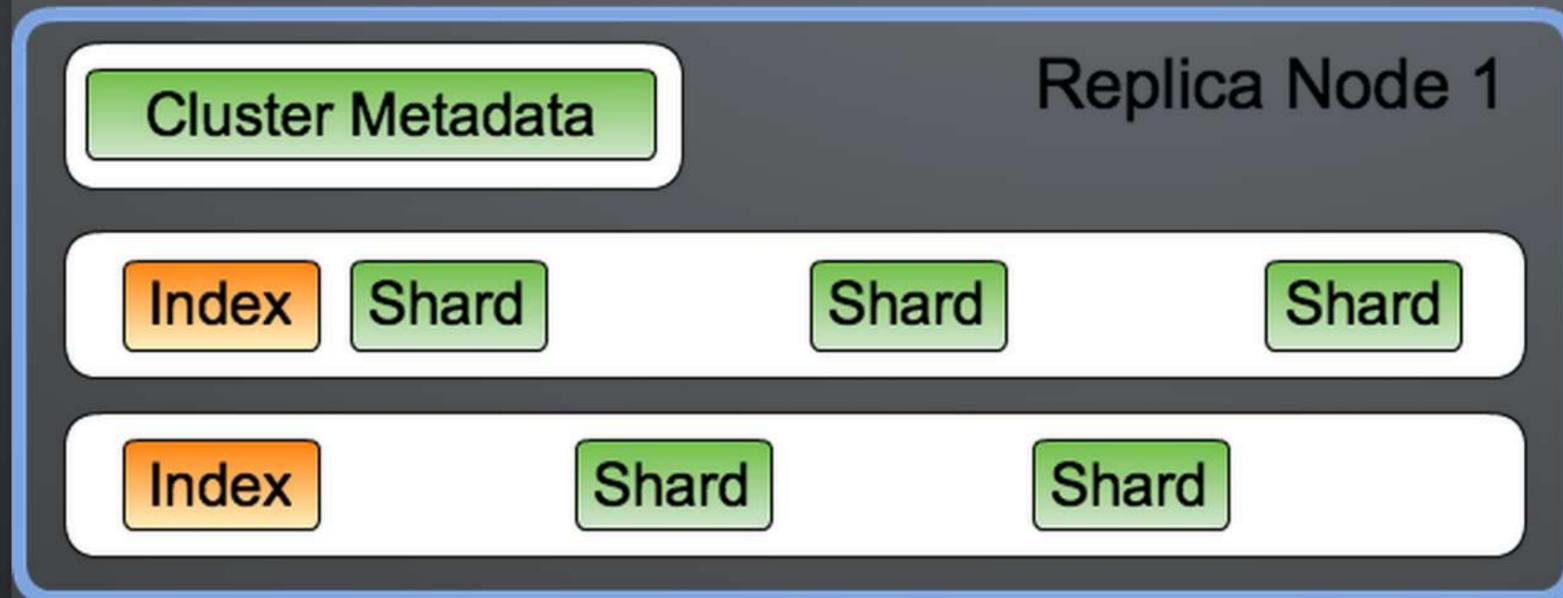
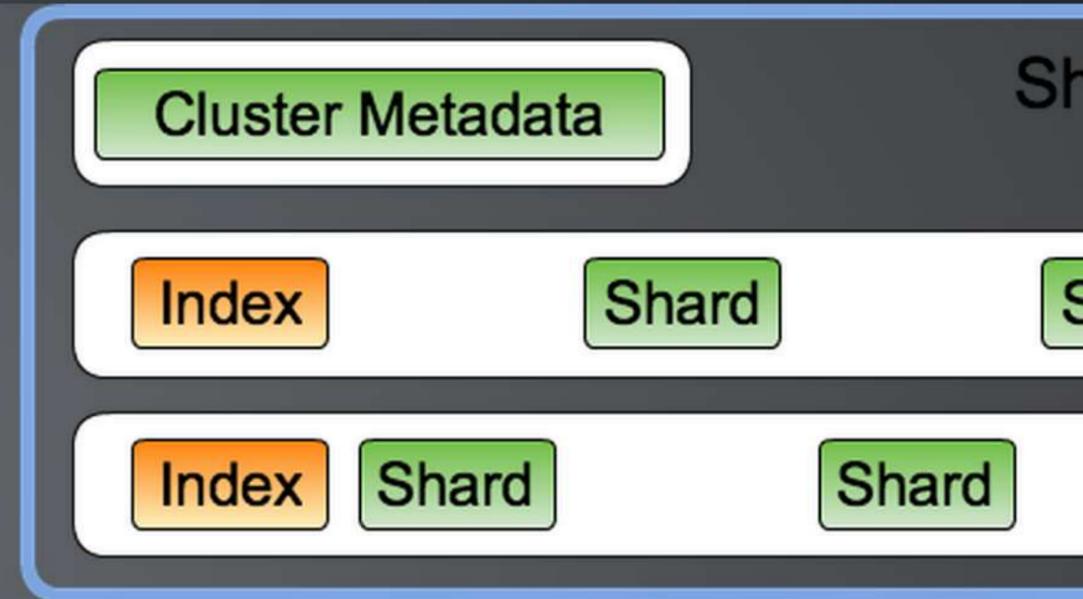
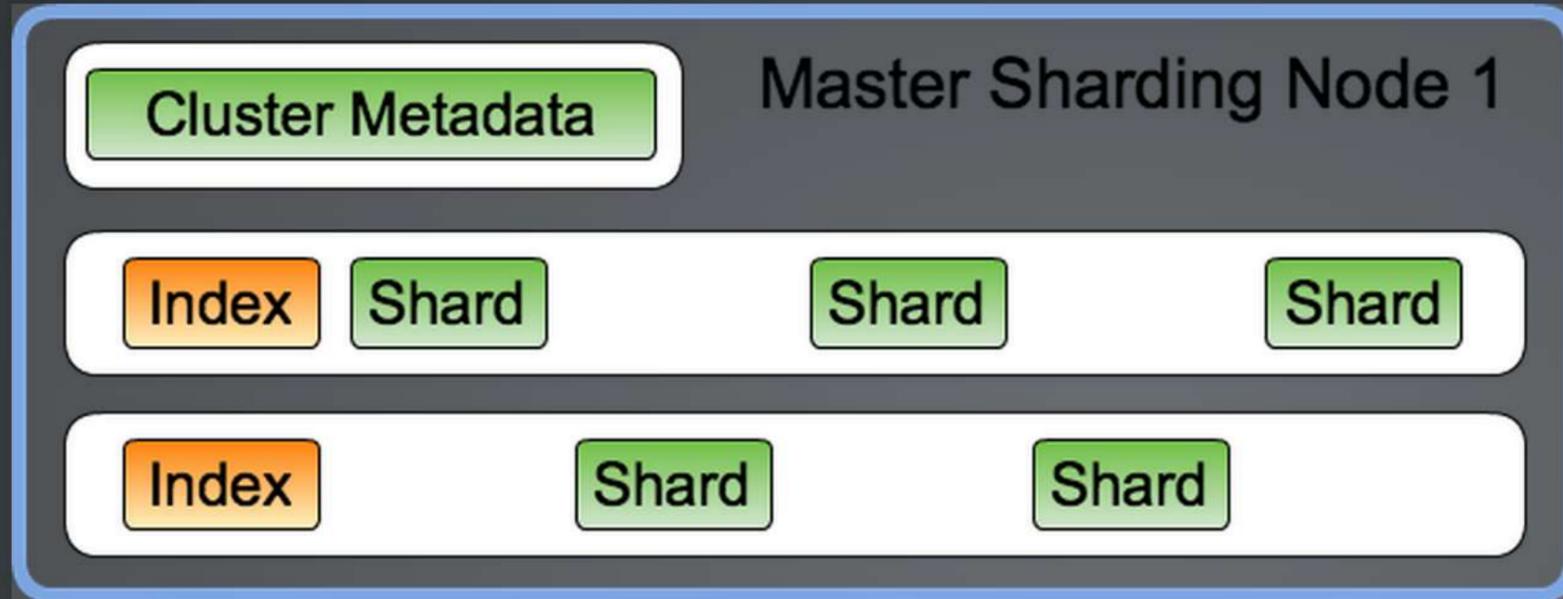
REPLICATION



SHARDING



REPLICATION & SHARDING



INSTALLATION - TAKES TWO MINUTES

- Download from github or elasticsearch.org zip file
- Unzip
`unzip elasticsearch-0.19.1.zip`
- Run
`elasticsearch/bin/elasticsearch -f`
- Check by connecting with your browser to <http://localhost:9200>



CONFIGURATION

- `config/elasticsearch.yml` or `config/elasticsearch.json`
- Application-wide settings (zen discovery, available analyzers)
- index default configurations (number of shards)
- Seperate logging file: `config/logging.yml` (simplified log4)



CONFIGURATION

```
discovery:
  zen:
    multicast.enabled: false

http:
  max_content_length: 100000

index:
  number_of_shards: 1

analysis:
  analyzer:
    default:
      type: standard

  lowercase_analyzer:
    type: custom
    tokenizer: standard
    filter: [standard, lowercase]
```



DATA REPRESENTATION IN JSON

```
{  
  "id": "1",  
  "name": "MacBook Air",  
  "price": 1099,  
  "descr": "Some lengthy never-read description",  
  "attributes": {  
    "color": "silver",  
    "display": 13.3,  
    "ram": 4  
  }  
}
```



INDEX YOUR PRODUCT

```
curl -X PUT localhost:9200/products/product/1 -d '{
  "id": "1",
  "name" : "MacBook Air",
  "price": 1099,
  "descr" : "Some lengthy never-read description",
  "attributes" : {
    "color" : "silver",
    "display" : 13.3,
    "ram" : 4
  }
}'
```

<http://localhost:9200/products/product/1>



JSON AS QUERY LANGUAGE

http://host:9200/products/product/_search

```
{ "query" : { "term" : { "name": "MacBook Air" } } }
```

```
{ "query" : { "prefix" : { "name": "Mac" } } }
```

```
{ "query" : { "range" : { "price" : { "from" : 1000, "to": 2000 } } } }
```

```
{ "from": 0, "size": 10, "query" : { "term" : { "name": "MacBook Air" } } }
```

```
{ "sort" : { "name" : { "order": "asc" } }, "query" : { "term" : { "name": "MacBook Air" } } }
```



JSON AS QUERY LANGUAGE

`http://host:9200/products/product/_search`

```
{ "took":206,"timed_out":false,
  "_shards":{"total":1,"successful":1,"failed":0},
  "hits":{"total":1,"max_score":2.098612,
    "hits":[ {
      "_index":"products1","_type":"product","_id":"1",
      "_score":2.098612, "_source" : {
        "id": "1",
        "name" : "MacBook Air",
        "price": 1099,
        "descr" : "Some lengthy never-read description",
        "attributes" : {
          "color" : "silver",
          "display" : 13.3,
          "ram" : 4
        }
      }
    }
  ]
}
```



CONFIGURATION - MAPPING

- On indexing the JSON document is parsed and all data types are extracted
- Mapping fields to datatypes is done automatically on first indexing
- Can be configured on a per-type basis
- Strings can have their own analyzer
- Sample types: float, long, boolean, date (+formatting), object
- One field can have multiple fields analyzed differently (lowercase, query)



SAMPLE MAPPING

```
{
  "product": {
    "properties": {
      "ProductId": { "type": "string", "index": "not_analyzed" },

      "ProductEnabled": { "type": "boolean" },
      "PiecesIncluded": { "type": "long" },
      "LastModified": { "type": "date", "format": "yyyy-MM-dd HH:mm:ss.SSS" },

      "AvailableInventory": { "type": "float" },
      "Price": { "type": "float" },

      "LongDescription": { "type": "string", "include_in_all" : true },

      "ProductName" : {
        "type" : "multi_field",
        "include_in_all" : true,
        "fields" : {
          "ProductName": { "type": "string", "index": "not_analyzed" },
          "lowercase": { "type": "string", "analyzer": "lowercase_analyzer" },
          "suggest" : { "type": "string", "analyzer": "suggest_analyzer" }
        }
      }
    }
  }
}
```



CONFIGURATION - ANALYZERS

- An analyzer consists of a Tokenizer and an arbitrary amount of filters
- Example:

```
suggest_analyzer:  
  type: custom  
  tokenizer: whitespace  
  filter: [standard, lowercase, shingle]
```

- Stripping html code:

```
char_filter: html_strip
```



JAVA API - CREATING A CLIENT

```
Settings settings = ImmutableSettings.settingsBuilder().  
    put("cluster.name", clusterName).build();  
  
InetSocketAddress addr =  
    new InetSocketAddress(host, port)  
  
Client client = new TransportClient(settings).  
    addTransportAddress(addr);
```



STARTING AN EMBEDDED SERVER

```
File config = new File("elasticsearch-local.yml");
String config = FileUtils.readFileToString(config);

Builder settingsBuilder = ImmutableSettings.settingsBuilder().
    loadFromSource(config);

Node node = NodeBuilder.nodeBuilder().
    settings(settingsBuilder).node();

Client client = node.client();
```



EXECUTING A QUERY

```
CountRequestBuilder countRequestBuilder =  
    new CountRequestBuilder(client)  
        .setQuery(QueryBuilders.termQuery("foo", "bar"))  
        .setIndices("products")  
        .setTypes("product");  
  
CountResponse response =  
    countRequestBuilder.execute().actionGet();  
long count = response.count();
```



SEARCH API OVERVIEW

- Index, Delete, Delete-By-Query, Get, Multiget, Bulk
- Search/Count queries (term query, prefix query, id, fuzzy...)
- Geo-based queries, TTL
- More like this, Highlighting
- Facetting, Percolation, Scripting



SEARCH - FACETTING

- Facetting adds aggregated information to a standard search query
- Term: Group results by a term
- Range: Group by price or date ranges
- Histogram: Group results in equally sized buckets, also as date histogram
- Statistical: Include statistical data like min, max, sum, avg & some more
- Geo distance: Group results around a coordinate



FACET QUERY

```
SearchRequestBuilder searchRequestBuilder = new SearchRequestBuilder(client)
    .setIndices("products")
    .setTypes("product");

searchRequestBuilder.setQuery(QueryBuilders.prefixQuery("ProductName.suggest", "macbook"));

searchRequestBuilder.addFacet(FacetBuilders.termsFacet("categoryFacet").field("CategoryId"));

SearchResponse searchResponse = searchRequestBuilder.execute().actionGet();

TermsFacet facet = searchResponse.getFacets().facet(TermsFacet.class, "categoryFacet");
List entries = facet.entries();
String term = entries.get(0).term();
int count = entries.get(0).count();
```



SEARCH - SCRIPTING

THIS IS WHERE YOUR OWN INTEGRATION BEATS ALL OTHERS

- Score down all your products without an image
- Dont include them in your results
- Score up products by an attribute like its product quality or stock
- Apply math operations on fields to change score



SEARCH API - PERCOLATION

IMPLEMENT A PRICE AGENT FOR FREE!

```
curl -X PUT localhost:9200/_percolator/products/pricecheck -d '{
  "query" : {
    "bool" : {
      "must" : { "term" : { "name" : "MacBook Air" } },
      "must" : { "range" : { "price" : { "from" : 200, "to" : 999 } } }
    }
  }
}'
{"ok":true,"_index":"_percolator","_type":"products","_id":"pricecheck","_version":1}

curl -X PUT 'localhost:9200/products/product/1?percolate=*' -d '{ "price": 1000, "name" : "MacBook Air"
}'
{"ok":true,"_index":"products","_type":"product","_id":"1","_version":1,"matches":[]}

curl -X PUT 'localhost:9200/products/product/2?percolate=*' -d '{ "price": 999, "name" : "MacBook Air"
}'
{"ok":true,"_index":"products","_type":"product","_id":"2","_version":1,"matches":["pricecheck"]}
```



INDICES API

- Aliases, Analyze
- Create, Delete, Exists, Open, Close, Optimize, Refresh, Flush, Settings
- Get, Put, Delete Mapping
- Get, update settings
- Snapshot
- Index templates (mappings + settings)
- Stats, Status
- Segments, Clear cache



CLUSTER API

- Health, State, Settings
- Nodes Info, Nodes Stats, Nodes Shutdown



MODULES

- REST, Thrift, Memcached, ZeroMQ
- JMX
- Scripting (MVEL, javascript, groovy, python, native)
- Discovery: EC2, Zen
- Cluster, Indices, Network, Transport



PLUGINS

- Analysis: Smart Chinese, ICU, IK, Mmseg, Hunspell
- Transport: Memcached, Thrift, ZeroMQ, Servlet
- Scripting: javascript, groovy, python
- Site plugins: BigDesk, Elasticsearch Head
- Misc: Mapper attachments, Hadoop, AWS cloud, Mock Solr, Suggester, PartialUpdate



RIVERS

- Interface to import data into elasticsearch
- CouchDB, Wikipedia, Twitter, RabbitMQ
- RSS, MongoDB
- Hint: When writing your own river, make sure you are implementing streaming

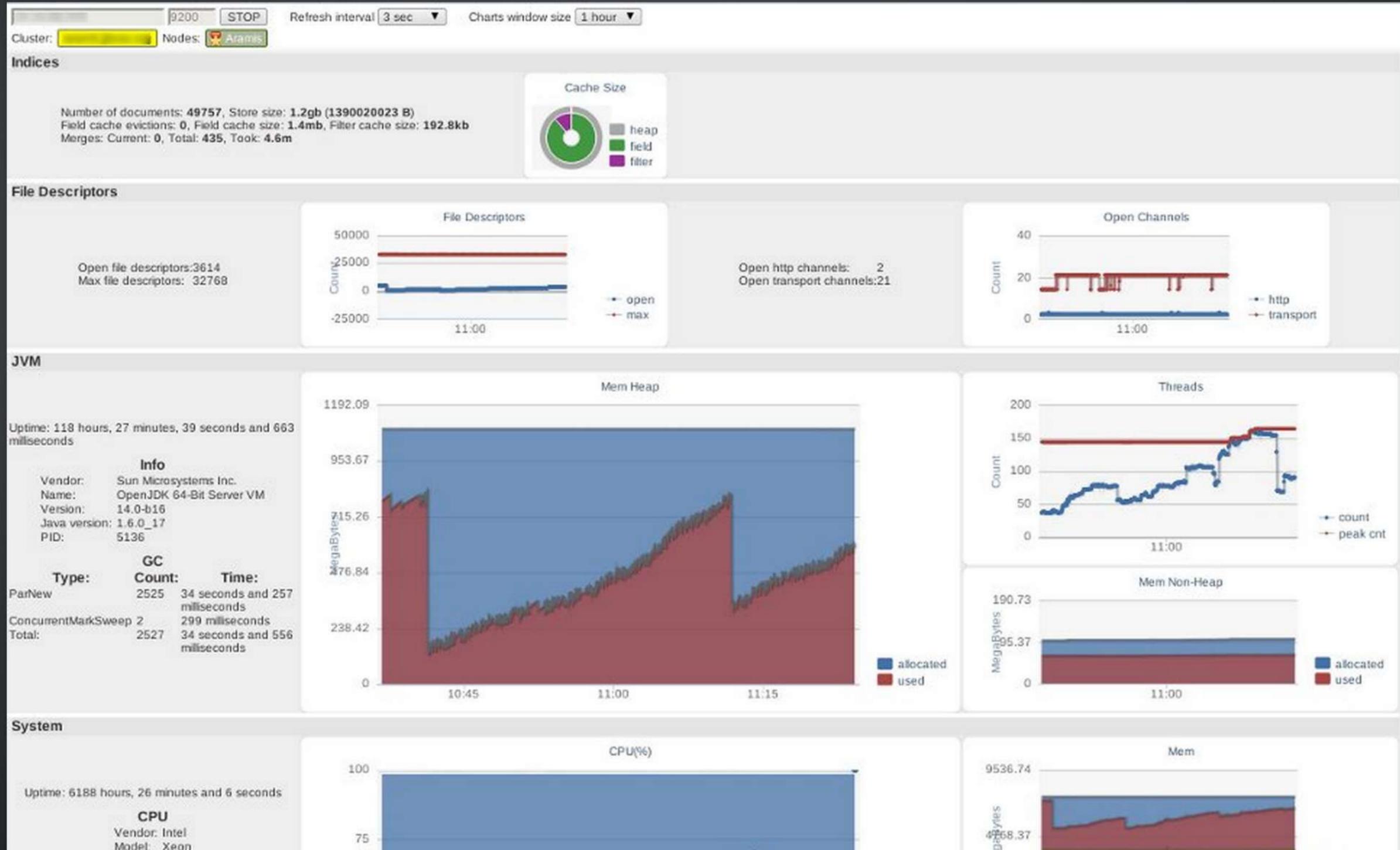


TOOLS

- BigDesk, Elasticsearch Head
- Chef, puppet
- RPMs and debian packages
- daikon CLI



BIGDESK SCREENSHOT



ELASTICSEARCH-HEAD SCREENSHOT

The screenshot displays the ElasticSearch Head interface. At the top, the title 'ElasticSearch' is followed by the URL 'http://192.168.7.8:9200/'. A 'Connect' button is visible, along with the cluster name 'Rick' and its health status 'cluster health: yellow (6, 18)'. Below this, a navigation bar includes 'Overview', 'Browser', 'Structured Query', and 'Any Request'. On the right side of the navigation bar, there are tabs for 'Info', 'Status', 'Nodes Stats', 'Cluster Nodes', 'Cluster State', and 'Cluster Health'. The main content area is titled 'Cluster Overview' and features a 'New Index' button. The interface is divided into several sections: a list of nodes on the left, and three index details on the right: 'cu_docs', 'bnvil', and 'cu_msg'. The 'cu_docs' index has a size of 180Gb (540Gb) and 995131 docs. The 'bnvil' index has a size of 80kb (480kb) and 90 docs. The 'cu_msg' index has a size of 313Gb (1.56Tb) and 10047450 docs. The 'anvil' index is also visible with a 'close' button. The node list includes Leon, Pris, Rick, Rachel, Zhora, Roy, and Unassigned. Each node has an 'Info' and 'Actions' button. The 'Actions' menu for the 'cu_msg' index is open, showing options like 'Refresh', 'Flush', 'Gateway Snapshot', 'Test Analyser', 'Close', and 'Delete...'. A modal window is open in the bottom right corner, displaying the configuration for the 'Leon' node, including its name, transport and http addresses, and OS details like refresh interval, CPU vendor, model, and frequency.

ElasticSearch <http://192.168.7.8:9200/> [Connect](#) **Rick** **cluster health: yellow (6, 18)**

Overview | **Browser** | Structured Query | Any Request | Info | Status | Nodes Stats | Cluster Nodes | Cluster State | Cluster Health

Cluster Overview [New Index](#)

Node	cu_docs	bnvil	cu_msg	anvil
Leon 3Wqr1xaCRu-b0uEzDkmrDg inet[/192.168.7.8:9202] Info Actions	0 1	0 1	0 1	Info Actions
Pris L8qx7ilfSI-kcKq_6bMbWw inet[/192.168.7.8:9204] Info Actions	0 1	0 1	0 1	Info Actions
Rick Vnpra1FNTGirwRfZsZ2RxQ inet[/192.168.7.8:9200] Info Actions	1 2	0 1	0 1 2 3 4	Info Actions
Rachel 87KsIv0FTVSkkqwENaja6A inet[/192.168.7.8:9203] Info Actions	1 2	0 1	0 1 2 3	Info Actions
Zhora b6NxRTxsR_WUQI5cXPKHbw inet[/192.168.7.8:9205] Info Actions	0 2	0 1	0 1 2 3 4	Info Actions
Roy _8RI2wYVT7Svn_v5F97jJA inet[/192.168.7.8:9201] Info Actions	0 2	0 1	0 1 2 3 4	Info Actions
Unassigned		0 0		

cu_docs
size: 180Gb (540Gb)
docs: 995131 (995131)
[Info](#) [Actions](#)

bnvil
size: 80kb (480kb)
docs: 90 (90)
[Info](#) [Actions](#)

cu_msg
size: 313Gb (1.56Tb)
docs: 10047450 (10140915)
[Info](#) [Actions](#)

anvil
index: close
[Info](#) [Actions](#)

```
{
  name: "Leon",
  transport_address: "inet[/192.168.7.8:9302]",
  attributes: {},
  http_address: "inet[/192.168.7.8:9202]",
  os: {
    refresh_interval: 5000,
    cpu: {
      vendor: "Intel",
      model: "Macmini4,1",
      mhz: 2400,
      total_cores: 2
    }
  }
}
```

LANGUAGE SUPPORT & SOFTWARE

- java, groovy, python, perl, [ruby](#), erlang, .net, [clojure](#)
- Integrations: grails, django, rails, catalyst, flume, terrastore, hadoop, symfony2, drupal, couchdb, play framework, node.js
- Software: [Graylog2](#)
- Elasticsearch as SaaS: [bonsai.io](#)



RUNNING IN PRODUCTION

- 220k products, one index, one shard (due to result grouping)
- Almost all queries have a big facetting query part (with filters)
- Don't expose your search engine to the internet!
- Write your own river
- Be prepared to upgrade every now and then



THANKS FOR LISTENING!

QUESTIONS?

Slides available at
<http://spinscale.github.com/>

alexander@reelsen.net
@spinscale



DOCUMENTATION & CREDITS

- elasticsearch.org
- [@elasticsearch](#)
- [Elasticsearch google group](#)

- Presentation done with [reveal.js](#)
- Cool zooming done with [zoom.js](#)

