



TECHNICAL UNIVERSITY OF LIBEREC  
Faculty of Mechatronics, Informatics  
and Interdisciplinary Studies ■

# ELASTICSEARCH

*Lukáš Matějů*  
29.4.2024 | DPB



# ČÁST I.: OPAKOVÁNÍ



# OPAKOVÁNÍ

- analýza
  - často také jako textová analýza
    - reálná aplikace jen na textová pole / hodnoty
  - vstupní data dokumentu v poli `_source`
    - ta ale nejsou přímo používána pro prohledávání
      - např. dlouhý popis produktu nelze efektivně prohledávat bez předzpracování
- při indexaci jsou textová pole analyzována
  - o analýzu se stará analyzátor skládající se ze tří komponent
    - znakový filtr (character filter)
    - tokenizér (tokenizer)
    - token filtr (token filter)
  - v základu standard analyzátor
    - žádný znakový filtr
    - tokenizér (Unicode segmentation)
    - lowercase token filtr
- výsledky analýzy uloženy v polích efektivních pro prohledávání



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# OPAKOVÁNÍ

- invertovaný index
  - mapování mezi výrazy a dokumenty, které je obsahují
  - výrazy jsou řazeny abecedně
  - dokumenty jsou odkazovány pomocí `_id`



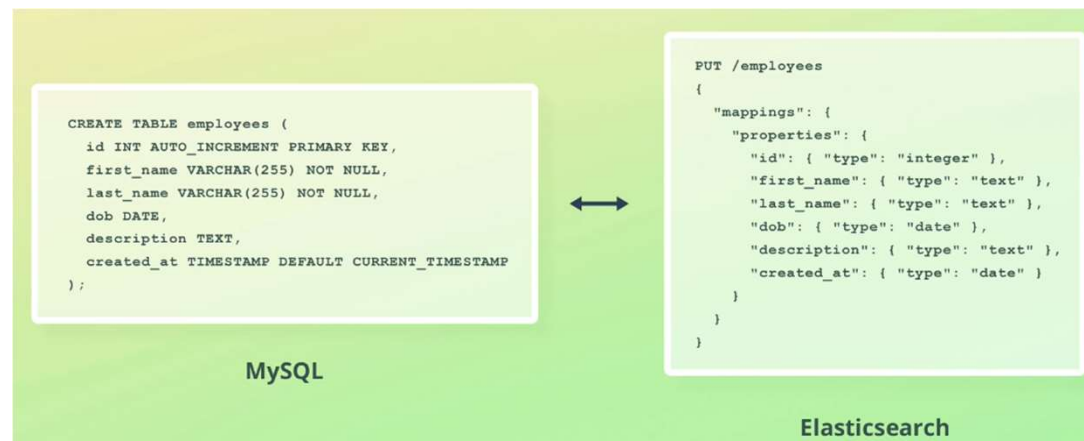
<https://www.udemy.com/course/elasticsearch-complete-guide/>





# OPAKOVÁNÍ

- mapování
  - definuje strukturu dokumentů, a jak jsou indexovány a ukládány
    - pole a jejich datové typy
  - zjednodušeně se dá přirovnat k relaci v relačních databázích
  - explicitní a dynamické
  - povoluje chybějící pole dokumentů



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# OPAKOVÁNÍ

- prohledávání – dotazovací typy
  - dotazy na úrovni termínů (term-level queries)

```
1 GET /products/_search
2 {
3   "query": {
4     "term": {
5       "name": "lobster"
6     }
7   }
8 }
```

```
10 "hits": {
11   "total": {
12     "value": 5,
13     "relation": "eq"
14   },
15   "max_score": 6.035804,
16   "hits": [
17     {
18       "_index": "products",
19       "_type": "doc",
20       "_id": "19",
21       "_score": 6.035804,
22       "_source": {
23         "name": "Lobster - Live",
```

```
1 GET /products/_search
2 {
3   "query": {
4     "term": {
5       "name": "Lobster"
6     }
7   }
8 }
```

```
1 {
2   "took": 1,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 0,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": [ ]
17  }
18 }
```

- fulltextové dotazy (full-text queries)

```
1 GET /products/_search
2 {
3   "query": {
4     "match": {
5       "name": "Lobster"
6     }
7   }
8 }
```

```
10 "hits": {
11   "total": {
12     "value": 5,
13     "relation": "eq"
14   },
15   "max_score": 6.035804,
16   "hits": [
17     {
18       "_index": "products",
19       "_type": "doc",
20       "_id": "19",
21       "_score": 6.035804,
22       "_source": {
23         "name": "Lobster - Live",
```



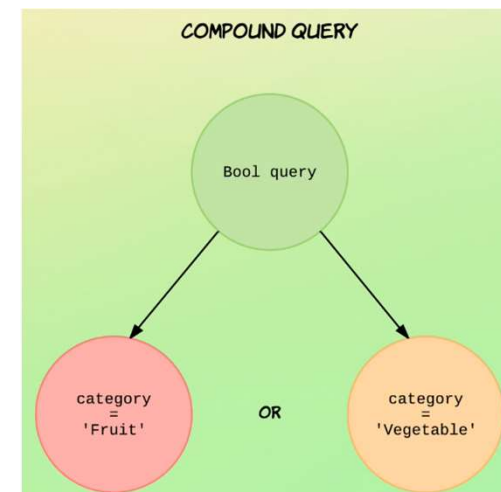


# OPAKOVÁNÍ

- prohledávání – dotazy s Booleovskou logikou
  - leaf queries
    - koncové dotazy
    - nejjednodušší dotazy provádějící jen jednu operaci
    - vyhledávají hodnotu v daném poli
    - např. term nebo match query
  - compound queries
    - tvoří složitější dotazy zaobalením leaf nebo dalších compound queries
    - např. dvě leaf query spojené přes bool query
    - využívají booleovskou logiku



<https://www.udemy.com/course/elasticsearch-complete-guide/>

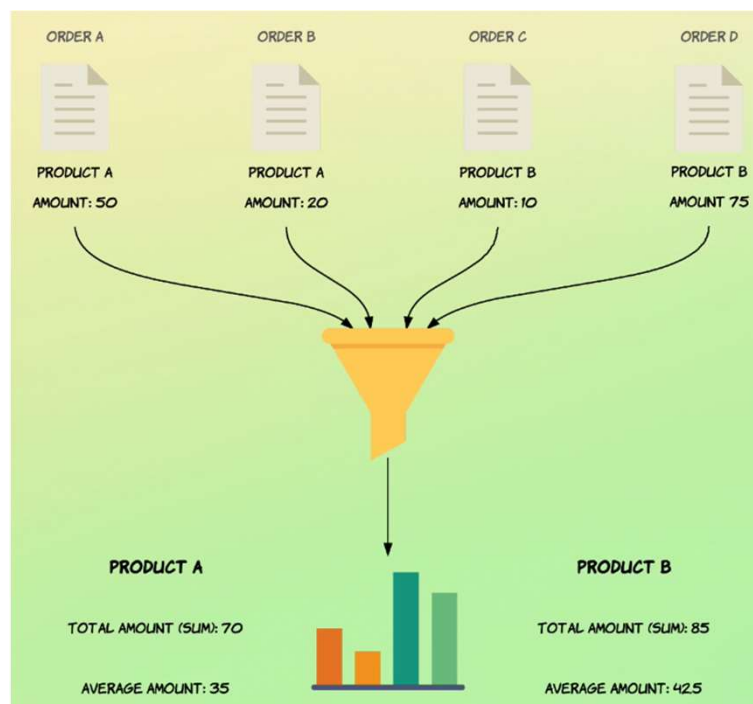


# ČÁST II.: ELASTICSEARCH AGREGACE



# AGREGACE

- způsob jak seskupit data a jak z nich získat statistiky a závěry



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# AGREGACE

- mapování a data pro následující ukázky
  - lines
    - specifikuje jednotlivé objednávky
    - nested objekt
    - id, amount, quantity
  - total\_amount
    - celkové množství
  - status
    - keyword
    - status objednávky
  - sales\_channel
    - kde bylo zboží prodáno
  - salesman
    - prodejce
    - id, jméno

```
1 {  
2   "_index" : "order",  
3   "_type" : "_doc",  
4   "_id" : "1",  
5   "_version" : 1,  
6   "_seq_no" : 0,  
7   "_primary_term" : 1,  
8   "found" : true,  
9   "_source" : {  
10    "purchased_at" : "2016-07-10T16:52:43Z",  
11    "lines" : [  
12      {  
13        "product_id" : 6,  
14        "amount" : 71.32,  
15        "quantity" : 1  
16      },  
17      {  
18        "product_id" : 3,  
19        "amount" : 58.96,  
20        "quantity" : 3  
21      },  
22      {  
23        "product_id" : 1,  
24        "amount" : 29.8,  
25        "quantity" : 3  
26      }  
27    ],  
28    "total_amount" : 160.08,  
29    "salesman" : {  
30      "id" : 11,  
31      "name" : "Matthus Mitkov"  
32    },  
33    "sales_channel" : "store",  
34    "status" : "processed"  
35  }  
36 }
```

```
1 PUT /order  
2 {  
3   "mappings": {  
4     "properties": {  
5       "purchased_at": {  
6         "type": "date"  
7       },  
8       "lines": {  
9         "type": "nested",  
10        "properties": {  
11          "product_id": {  
12            "type": "integer"  
13          },  
14          "amount": {  
15            "type": "double"  
16          },  
17          "quantity": {  
18            "type": "short"  
19          }  
20        },  
21      },  
22      "total_amount": {  
23        "type": "double"  
24      },  
25      "status": {  
26        "type": "keyword"  
27      },  
28      "sales_channel": {  
29        "type": "keyword"  
30      },  
31      "salesman": {  
32        "type": "object",  
33        "properties": {  
34          "id": {  
35            "type": "integer"  
36          },  
37          "name": {  
38            "type": "text"  
39          }  
40        }  
41      }  
42    }  
43  }  
44 }
```



# AGREGACE

- agregace metrik (metric aggregation)
- nejzákladnější agregace
- dvě kategorie
  - single-value numeric metric aggregation
    - výstupem je jediná hodnota (např. průměr nebo suma)
  - multi-value numeric metric aggregation
    - výstupem je více hodnot
- pole aggs
  - zajímá nás statistika všech prodejů přes všechny objednávky
    - celkový objem zboží na objednávku
    - minimum, maximum, průměr a suma celkového objemu

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "total_sales": {
6       "sum": {
7         "field": "total_amount"
8       }
9     },
10    "avg_sale": {
11      "avg": {
12        "field": "total_amount"
13      }
14    },
15    "min_sale": {
16      "min": {
17        "field": "total_amount"
18      }
19    },
20    "max_sale": {
21      "max": {
22        "field": "total_amount"
23      }
24    }
25  }
26 }
```

```
1 {
2   "took": 1,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 1000,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": [ ]
17  },
18  "aggregations": {
19    "max_sale": {
20      "value": 281.77
21    },
22    "avg_sale": {
23      "value": 109.20961
24    },
25    "min_sale": {
26      "value": 10.27
27    },
28    "total_sales": {
29      "value": 109209.61
30    }
31  }
32 }
```



# AGREGACE

- agregace metrik (metric aggregation)
  - cardinality
    - počet unikátních hodnot pro dané pole
  - value\_count
    - počet hodnot pro dané pole
  - stats
    - vybrané základní agregace v jednom balíčku

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "total_salesmen": {
6       "cardinality": {
7         "field": "salesman.id"
8       }
9     }
10  }
11 }
```

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "values_count": {
6       "value_count": {
7         "field": "total_amount"
8       }
9     }
10  }
11 }
```

```
18 GET /order/_search
19 {
20   "size": 0,
21   "aggs": {
22     "amount_stats": {
23       "stats": {
24         "field": "total_amount"
25       }
26     }
27   }
28 }
```

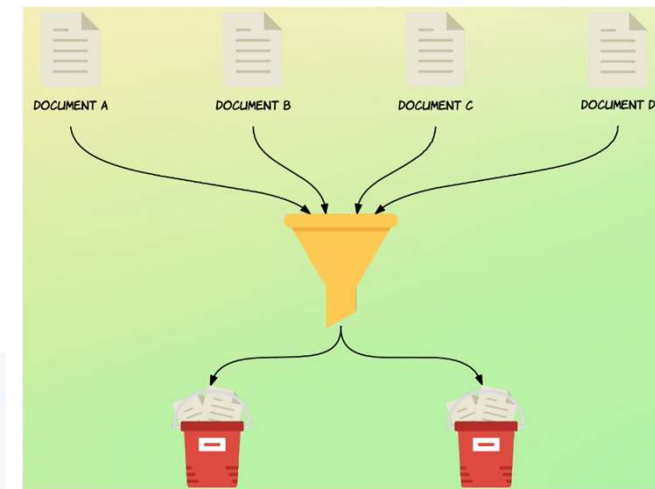
```
18 "aggregations": {
19   "amount_stats": {
20     "count": 1000,
21     "min": 10.27,
22     "max": 281.77,
23     "avg": 109.20961,
24     "sum": 109209.61
25   }
26 }
27 }
28 }
```

# AGREGACE

- bucket agregace
  - vytváří skupiny (buckets) dokumentů
  - každý bucket má kritéria, která rozhodují, jestli do nich dokument spadá
  - agregace výrazů (terms)
    - vytváří bucket pro každý unikátní výraz
    - agregace podle statusu objednávky

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "status_terms": {
6       "terms": {
7         "field": "status",
8         "missing": "N/A",
9         "order": {
10          "_key": "asc"
11        }
12      }
13    }
14  }
15 }
```

```
18 "aggregations": {
19   "status_terms": {
20     "doc_count_error_upper_bound": 0,
21     "sum_other_doc_count": 0,
22     "buckets": [
23       {
24         "key": "cancelled",
25         "doc_count": 196
26       },
27       {
28         "key": "completed",
29         "doc_count": 204
30       },
31       {
32         "key": "confirmed",
33         "doc_count": 192
34       },
35       {
36         "key": "pending",
37         "doc_count": 199
38       },
39       {
40         "key": "processed",
41         "doc_count": 209
42       }
43     ]
44   }
45 }
46 }
```



<https://www.udemy.com/course/elasticsearch-complete-guide/>



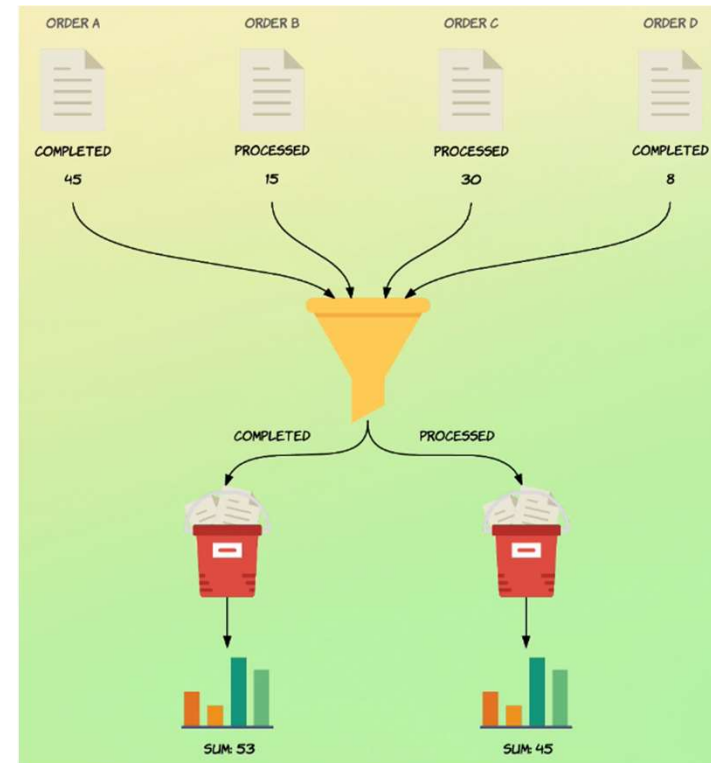
# AGREGACE

- bucket agregace
  - vnořené (nested) agregace
    - také jako subagregace
    - vnoření dalšího aggs pole
      - další agregace na získaných buckets
    - statistika pro každý status

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "status_terms": {
6       "terms": {
7         "field": "status"
8       }
9     },
10    "aggs": {
11      "status_stats": {
12        "stats": {
13          "field": "total_amount"
14        }
15      }
16    }
17  }
18 }
```

```
{
  "key": "processed",
  "doc_count": 209,
  "status_stats": {
    "count": 209,
    "min": 10.27,
    "max": 281.77,
    "avg": 109.30703349282295,
    "sum": 22845.17
  }
},
```

```
{
  "key": "completed",
  "doc_count": 204,
  "status_stats": {
    "count": 204,
    "min": 10.93,
    "max": 260.59,
    "avg": 113.54058823529411,
    "sum": 23162.28
  }
},
```



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# AGREGACE

- bucket agregace
  - filter agregace
    - filter
    - filtrování dokumentů

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "low_value": {
6       "filter": {
7         "range": {
8           "total_amount": {
9             "lt": 50
10          }
11        }
12      },
13      "aggs": {
14        "avg_amount": {
15          "avg": {
16            "field": "total_amount"
17          }
18        }
19      }
20    }
21  }
22 }
```

```
1 {
2   "took" : 6,
3   "timed_out" : false,
4   "_shards" : {
5     "total" : 1,
6     "successful" : 1,
7     "skipped" : 0,
8     "failed" : 0
9   },
10  "hits" : {
11    "total" : {
12      "value" : 1000,
13      "relation" : "eq"
14    },
15    "max_score" : null,
16    "hits" : [ ]
17  },
18  "aggregations" : {
19    "low_value" : {
20      "doc_count" : 164,
21      "avg_amount" : {
22        "value" : 32.59371951219512
23      }
24    }
25  }
26 }
```





# AGREGACE

- bucket agregace
  - filters agregace
    - filters
    - definování pravidel určujících dokumenty patřící do vybraného bucket
    - např. bucket pro trička bude obsahovat jen trička

```
1 GET /recipes/_search
2 {
3   "size": 0,
4   "aggs": {
5     "my_filter": {
6       "filters": {
7         "pasta": {
8           "match": {
9             "title": "pasta"
10          }
11        },
12        "spaghetti": {
13          "match": {
14            "title": "spaghetti"
15          }
16        }
17      }
18    },
19    "aggs": {
20      "avg_rating": {
21        "avg": {
22          "field": "ratings"
23        }
24      }
25    }
26  }
27 }
28 }
29 }
```

```
1 {
2   "took": 4,
3   "timed_out": false,
4   "_shards": {
5     "total": 1,
6     "successful": 1,
7     "skipped": 0,
8     "failed": 0
9   },
10  "hits": {
11    "total": {
12      "value": 21,
13      "relation": "eq"
14    },
15    "max_score": null,
16    "hits": [ ]
17  },
18  "aggregations": {
19    "my_filter": {
20      "buckets": {
21        "pasta": {
22          "doc_count": 9,
23          "avg_rating": {
24            "value": 3.4125
25          }
26        },
27        "spaghetti": {
28          "doc_count": 4,
29          "avg_rating": {
30            "value": 2.3684210526315788
31          }
32        }
33      }
34    }
35  }
36 }
```





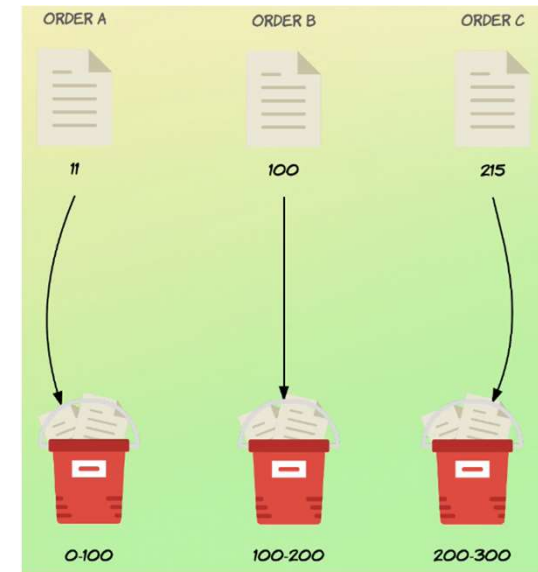
# AGREGACE

- bucket agregace
  - agregace rozsahů (range)
    - definice rozsahů, kde každý rozsah vytvoří bucket
    - dvě varianty
      - range – hodnoty
      - date\_range – datумы

```
18 GET /order/_search
19 {
20   "size": 0,
21   "aggs": {
22     "amount_distribution": {
23       "range": {
24         "field": "total_amount",
25         "ranges": [
26           {
27             "to": 50
28           },
29           {
30             "from": 50,
31             "to": 100
32           },
33           {
34             "from": 100
35           }
36         ]
37       }
38     }
39   }
40 }
```

```
18 "aggregations" : {
19   "amount_distribution" : {
20     "buckets" : [
21       {
22         "key" : "*-50.0",
23         "to" : 50.0,
24         "doc_count" : 164
25       },
26       {
27         "key" : "50.0-100.0",
28         "from" : 50.0,
29         "to" : 100.0,
30         "doc_count" : 347
31       },
32       {
33         "key" : "100.0-*",
34         "from" : 100.0,
35         "doc_count" : 489
36       }
37     ]
38   }
39 }
40 }
```

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "purchased_ranges": {
6       "date_range": {
7         "field": "purchased_at",
8         "ranges": [
9           {
10            "from": "2016-01-01",
11            "to": "2016-01-01|+6M"
12          },
13          {
14            "from": "2016-01-01|+6M",
15            "to": "2016-01-01|+1y"
16          }
17        ]
18      }
19    }
20  }
21 }
```



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# AGREGACE

- bucket agregace
  - agregace podle chybějících polí
    - missing
    - pole s hodnotou null nebo úplně chybící

```
1 POST /order/_doc/1001
2 {
3   "total_amount": 100
4 }
5
6 POST /order/_doc/1002
7 {
8   "total_amount": 200,
9   "status": null
10 }
```

```
1 GET /order/_search
2 {
3   "size": 0,
4   "aggs": {
5     "orders_without_status": {
6       "missing": {
7         "field": "status"
8       }
9     }
10 }
11 }
```

```
1 {
2   "took" : 1,
3   "timed_out" : false,
4   "_shards" : {
5     "total" : 1,
6     "successful" : 1,
7     "skipped" : 0,
8     "failed" : 0
9   },
10  "hits" : {
11    "total" : {
12      "value" : 1002,
13      "relation" : "eq"
14    },
15    "max_score" : null,
16    "hits" : [ ]
17  },
18  "aggregations" : {
19    "orders_without_status" : {
20      "doc_count" : 2
21    }
22  }
23 }
```

# AGREGACE

- bucket agregace
  - agregace vnořených objektů
    - nested

```
1 GET /departments/_search
2 {
3   "size": 0,
4   "aggs": {
5     "employees": {
6       "nested": {
7         "path": "employees"
8       },
9       "aggs": {
10        "minimum_age": {
11          "min": {
12            "field": "employees.age"
13          }
14        }
15      }
16    }
17  }
18 }
```

```
1 {
2   "took" : 3,
3   "timed_out" : false,
4   "_shards" : {
5     "total" : 1,
6     "successful" : 1,
7     "skipped" : 0,
8     "failed" : 0
9   },
10  "hits" : {
11    "total" : {
12      "value" : 2,
13      "relation" : "eq"
14    },
15    "max_score" : null,
16    "hits" : [ ]
17  },
18  "aggregations" : {
19    "employees" : {
20      "doc_count" : 15,
21      "minimum_age" : {
22        "value" : 19.0
23      }
24    }
25  }
26 }
```

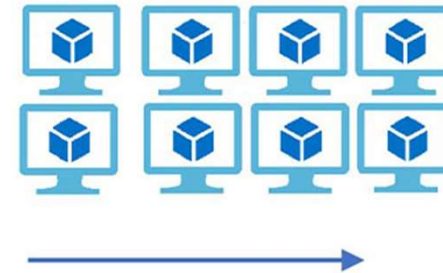
# ČÁST III.: ELASTICSEARCH SHARDING A REPLIKACE

# SHARDING

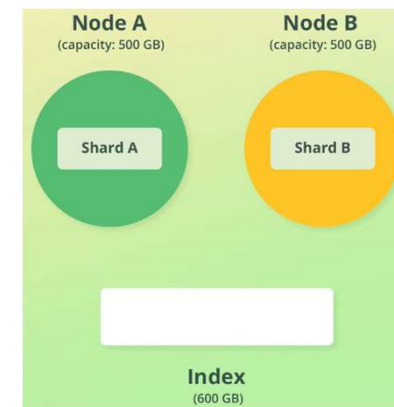
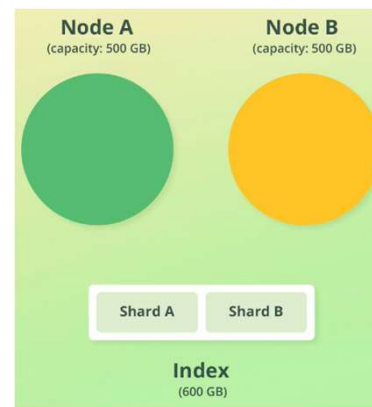
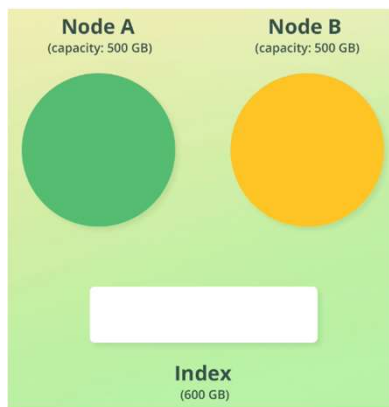
- rozdělení indexů na menší části
  - každá část je nazývána shard
  - prováděno na úrovni indexů
    - ne uzlů nebo clusterů
- horizontální škálování
  - škálování datového úložiště

## Horizontal Scaling

( Add more instances )



<https://www.webairy.com/horizontal-and-vertical-scaling/>



<https://www.udemy.com/course/elasticsearch-complete-guide/>



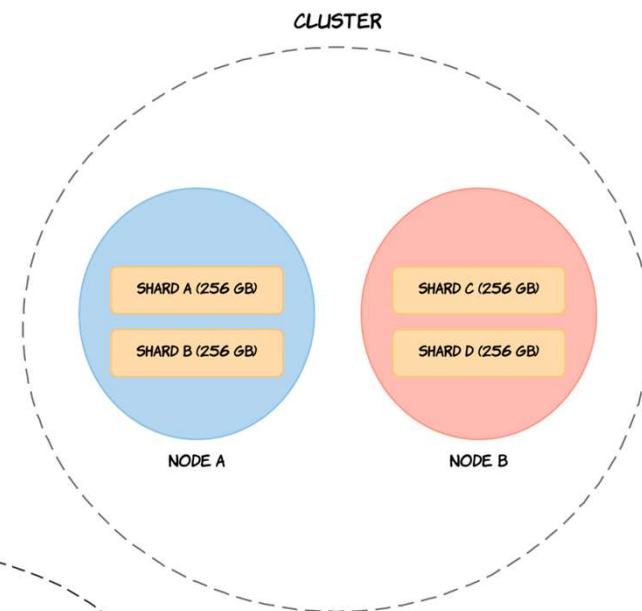
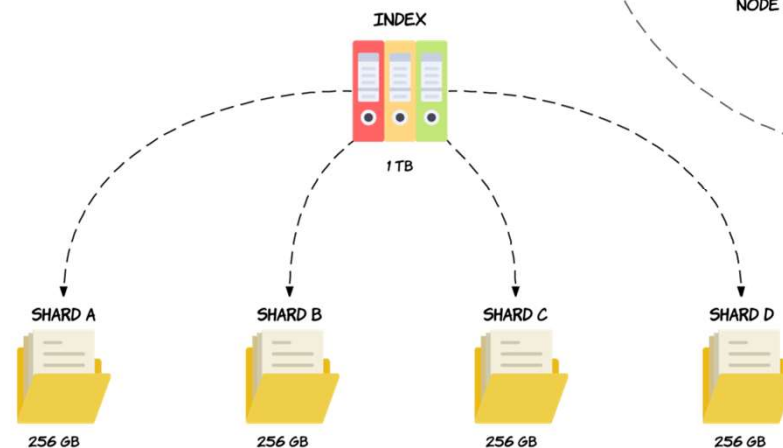
# SHARDING

- každý shard
  - slouží víceméně jako samostatný index
    - Apache Lucene index
    - každý index je složen z alespoň jednoho Lucene indexu
  - nemá předdefinovanou velikost
    - roste s vkládanými dokumenty
  - může obsahovat až 2 miliardy dokumentů
- v základu každý index má jeden shard
  - možné konfigurovat
    - zvýšení počtu shardů pomocí Split API
    - snížení Shrink API
  - vhodné rozmyslet dopředu
    - neexistuje optimální počet shardů
    - závisí na konkrétní aplikaci



# SHARDING

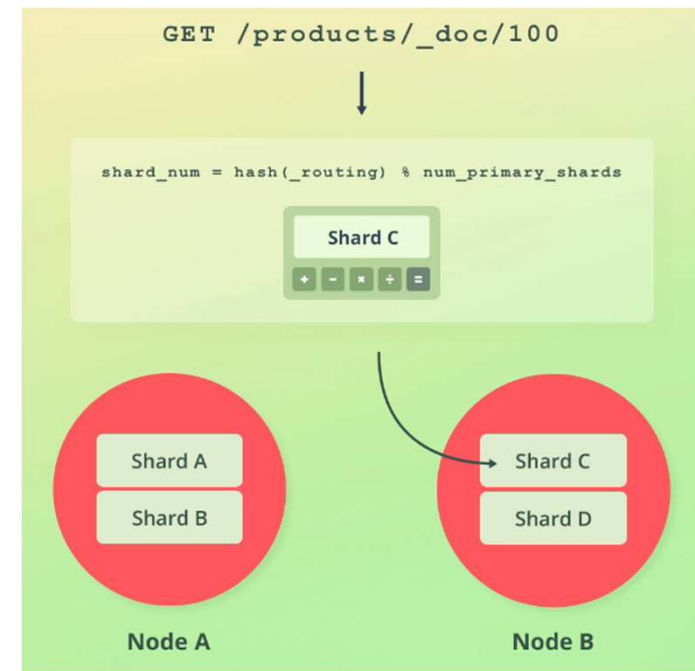
- proč shardovat?
  - ukládání více dokumentů
    - miliardy dokumentů v jednom indexu
  - rozdělení velkých indexů
    - rozdělení na části a uložení v uzlech
  - zlepšení výkonu
    - paralelní zpracování dotazů v shardech



<https://codingexplained.com/coding/elasticsearch/understanding-sharding-in-elasticsearch>

# ROUTING

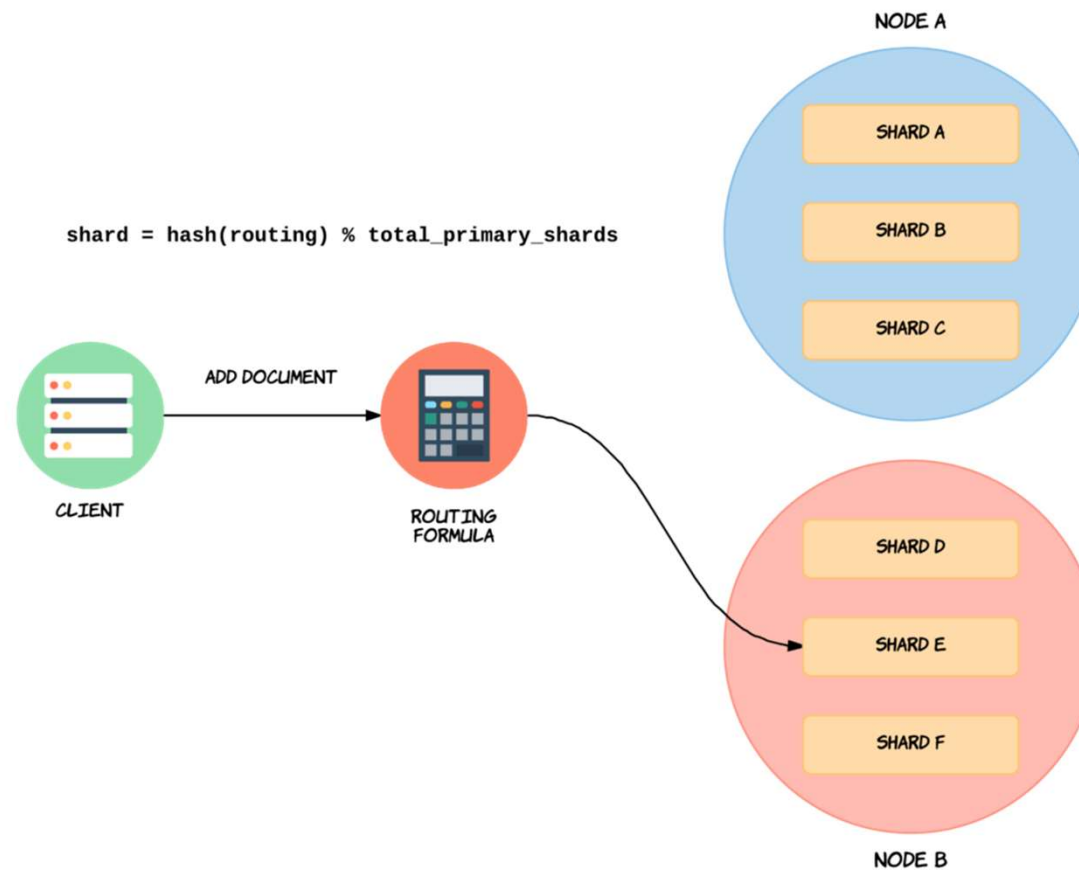
- proces přiřazující dokumentu shard
  - určuje, kde bude dokument uložen
    - na jakém konkrétním shardu
  - udává, jak jsou dokumenty po zaindexování hledány
    - pro operace čtení, update nebo mazání
- Elasticsearch využívá jednoduchý vzorec
$$\text{shard\_num} = \text{hash}(\text{\_routing}) \% \text{num\_primary\_shards}$$
  - `\_routing` – odpovídá `\_id` dokumentu
    - možnost vlastního nastavení
  - `num\_primary\_shards` – celkový počet shardů
- zajišťuje rovnoměrné rozmístění dokumentů mezi shardy



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# ROUTING



<https://codingexplained.com/coding/elasticsearch/understanding-sharding-in-elasticsearch>



# REPLIKACE

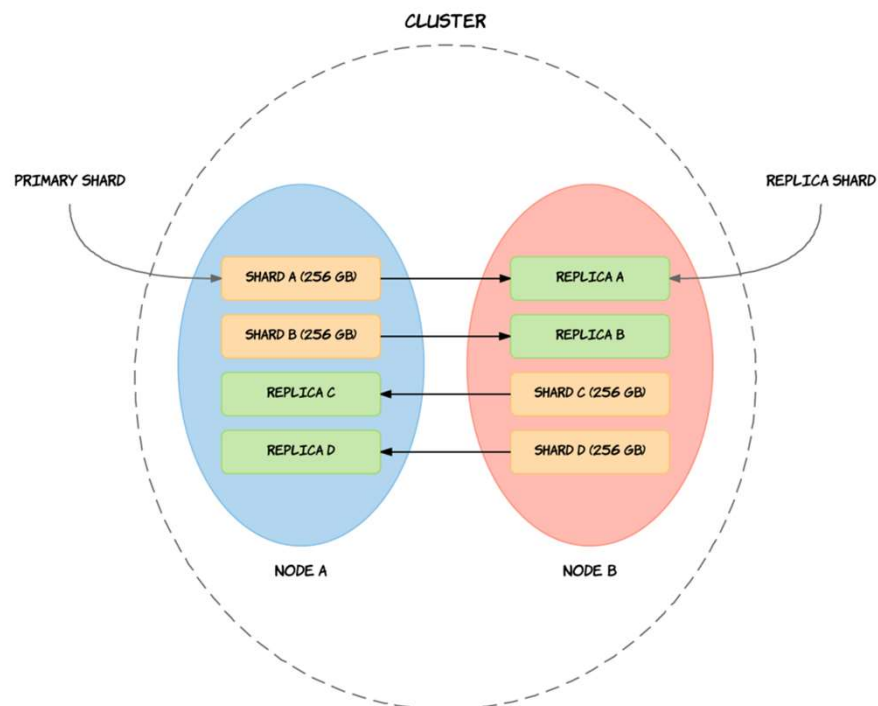
- automatická distribuce změn v originále do jeho kopií
- replikace v Elasticsearch
  - ochrana proti chybám (fault tolerance)
    - selhání hardware
  - nativní podpora, zapnuta automaticky
    - často obtížné nastavit v jiných databázích
  - definována na úrovni indexů
  - vytváří úplné kopie shardů
    - repliky shardů (replica shards)
    - původní shard se nazývá primární shard (primary shard)
    - primary shard a replica shards tvoří replication group
    - mohou plně obsluhovat dotazy
  - při vytváření indexu možnost nastavení počtu replik (1 v základu)



<https://www.udemy.com/course/elasticsearch-complete-guide/>

# REPLIKACE

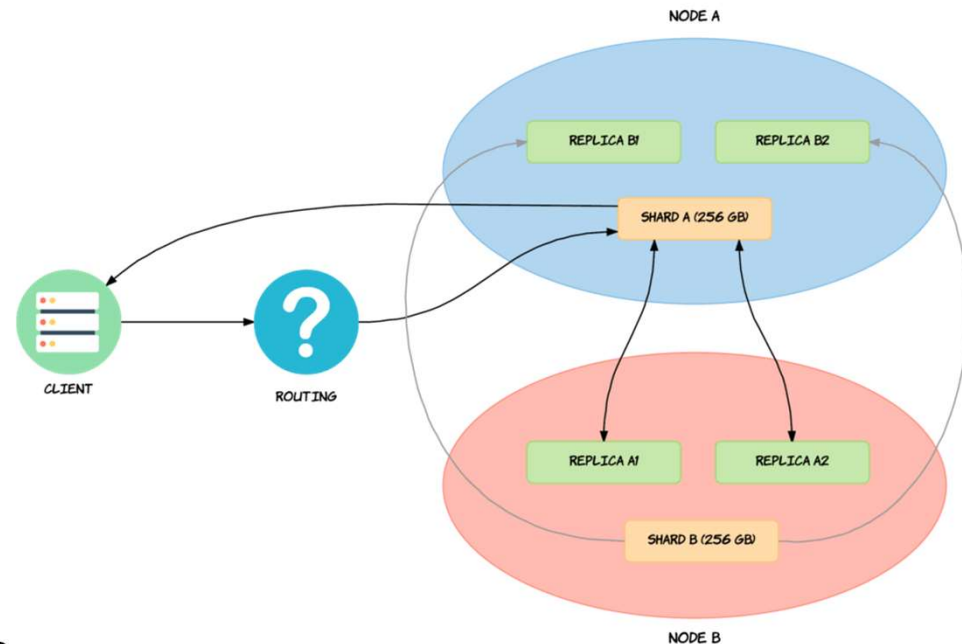
- replikace v Elasticsearch



<https://codingexplained.com/coding/elasticsearch/understanding-replication-in-elasticsearch>

# REPLIKACE

- replikace v Elasticsearch
  - pro synchronizaci využíván model primary-backup
  - primární shard je vstupem pro indexovací operace
    - vkládání, update, mazání, ...
  - primární shard operaci validuje a následně lokálně aplikuje na svá data
  - po dokončení je operace předána všem replikám
  - operace je paralelně provedena na všech replikách
  - po potvrzení od všech replik informuje primární shard klienta

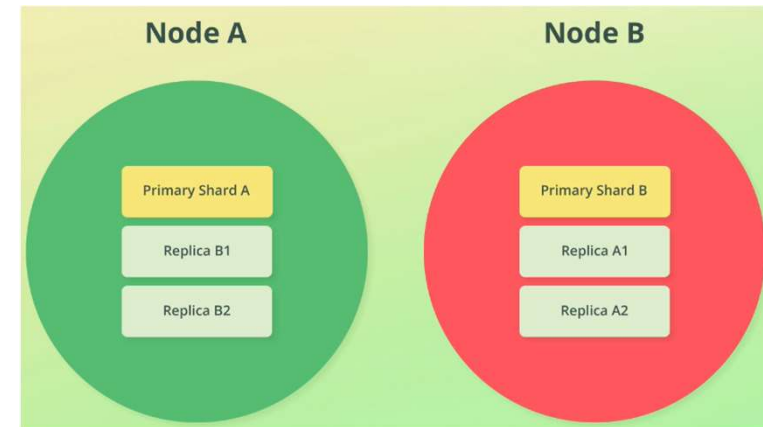


<https://codingexplained.com/coding/elasticsearch/understanding-replication-in-elasticsearch>

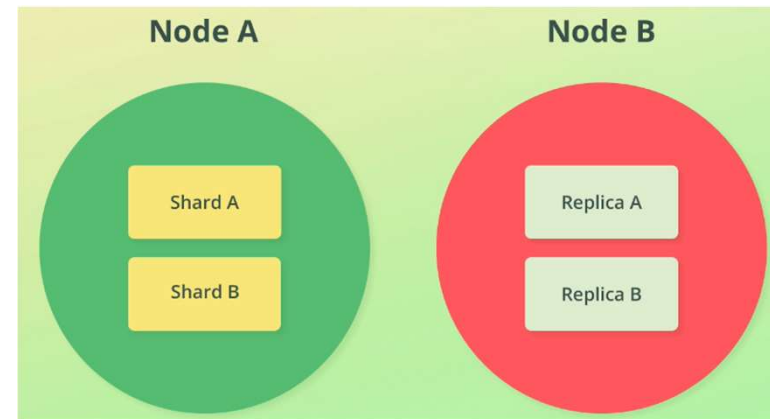
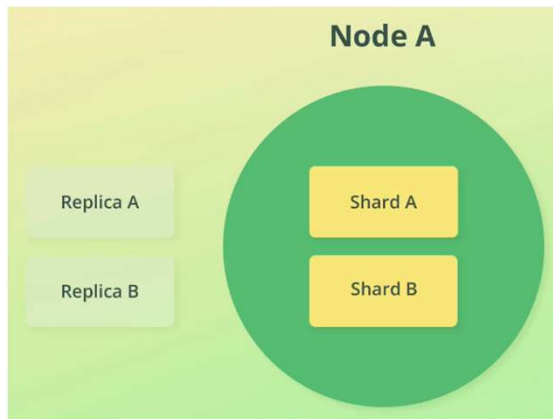


# REPLIKACE

- replikace v Elasticsearch
  - jak zabránit ztrátě dat?
    - repliky nejsou na stejném uzlu jako primární shard
      - jsou potřeba alespoň dva uzly
      - Elasticsearch přidává repliky jen pro clustery s alespoň dvěma uzly
  - zvyšuje dostupnost indexu



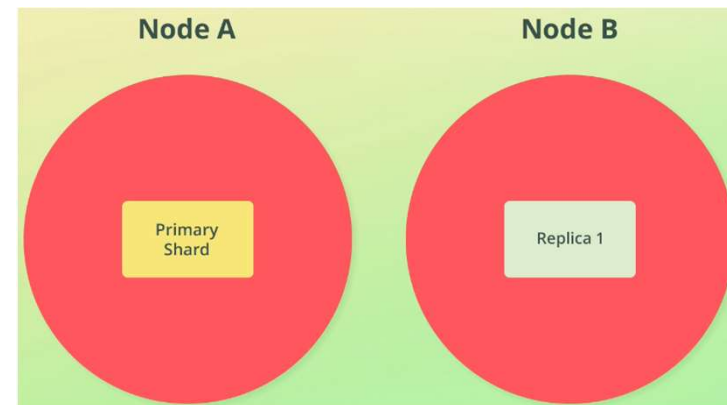
<https://www.udemy.com/course/elasticsearch-complete-guide/>



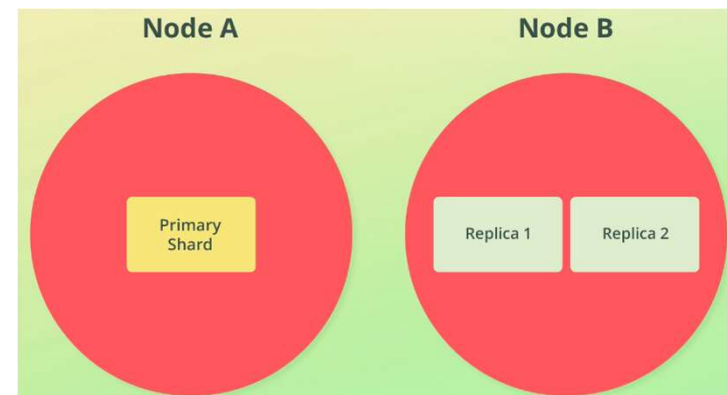


# REPLIKACE

- replikace v Elasticsearch
  - ideální počet replik?
    - závisí na aplikaci
    - je možnost data obnovit z jiného zdroje?
    - je problémem krátká nedostupnost?
      - jedna replika v případě, že ztráta dat není katastrofická
      - alespoň dvě pro klíčové systémy
        - např. aplikace pro nemocnice, ...
  - může také zvyšovat propustnost dotazování
    - repliky plně obsluhují dotazy
    - obsluhují je paralelně
    - dochází ke zvýšení zpracovaných dotazů ve stejný čas



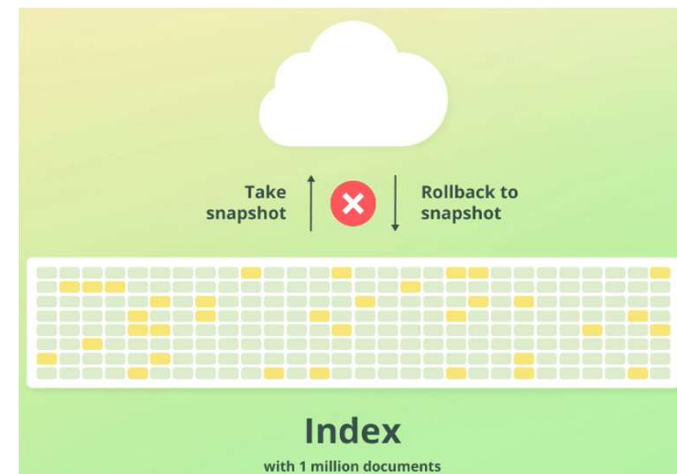
<https://www.udemy.com/course/elasticsearch-complete-guide/>





# SNAPSHOTY

- snímek databáze určený pro zálohu stavu databáze
- snímky v Elasticsearch
  - umožňují obnovu databáze do stavu, ve kterém byl snímek pořízen
  - mohou být definované nad indexy ale i nad celým clusterem
  - export dat do souboru
- snapshoty vs. replikace
  - snímky slouží pro zálohu stavu
    - nezajišťují dostupnost aktuálních dat
  - replikace pro vysokou dostupnost
    - jen pro aktuální data
    - poskytuje obnovu při výpadku uzlu
    - neumožňuje obnovu historických dat



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# PŘIDÁNÍ UZLU (PRO VÝVOJ)

- kompletně změněno od verze 8.0
  - dříve jen další instance Elasticsearch
  - nyní založené na tzv. enrollment token
  - nejprve je potřeba enrollment token vygenerovat na libovolném uzlu

```
bin/elasticsearch-create-enrollment-token -s node
```

- následně je enrollment token předán při spuštění uzlu

```
bin/elasticsearch --enrollment-token <enrollment-token>
```

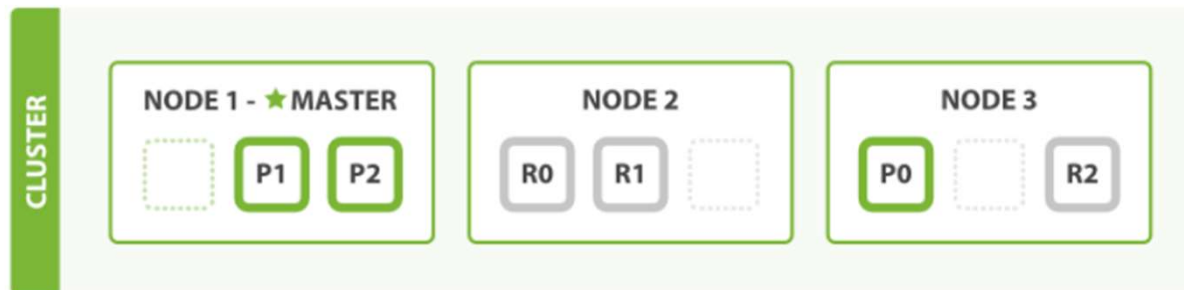
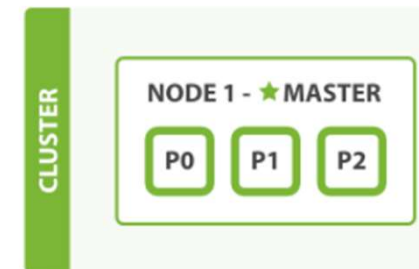
- certifikáty a klíče jsou automaticky vygenerovány do config/certs
- enrollment token může být předán i dalším uzlům pro další rozšíření clusteru
- jen uzly na stejném stroji se mohou připojit bez dalšího nastavení

<https://www.elastic.co/guide/en/elasticsearch/reference/current/configuring-stack-security.html#stack-enroll-nodes>



## PŘIDÁNÍ UZLU (PRO VÝVOJ)

- pro praktické nasazení v clusteru potřeba i [další nastavení](#)
  - [transport.host](#) musí být nastaven pro připojení uzlů z různých zdrojů
- nelze použít pro Elastic Cloud
  - [přidání/odebrání](#) uzlů v interface Elastic Cloud
  - další uzly jsou samozřejmě placené
  - přidělení replica shards automaticky



<https://www.elastic.co/guide/en/elasticsearch/reference/current/add-elasticsearch-nodes.html>



# PŘIDÁNÍ UZLU (PRO VÝVOJ)

```
[2020-11-18T14:10:41,619][INFO ][o.e.c.r.a.AllocationService] [node-1] updating number_of_replicas to [1] for indices [ilm-history-3-000001, .apm-custom-link, .apm-agent-configuration, .kibana_1, .kibana_task_manager_1, .kibana-event-log-7.10.0-000001]
[2020-11-18T14:10:41,621][INFO ][o.e.c.s.MasterService] [node-1] node-join[{node-2}{lmhdqVKZQuuMO6qsIELH8Q}{vSeL76bNTrC1I75EX4OumA}{127.0.0.1}{127.0.0.1:9301}{cdhilmrstw}{ml.machine_memory=34315870208, ml.max_open_jobs=20, xpack.installed=true, transform.node=true} join existing leader], term: 797, version: 153, delta: added [{node-2}{lmhdqVKZQuuMO6qsIELH8Q}{vSeL76bNTrC1I75EX4OumA}{127.0.0.1}{127.0.0.1:9301}{cdhilmrstw}{ml.machine_memory=34315870208, ml.max_open_jobs=20, xpack.installed=true, transform.node=true}]
[2020-11-18T14:10:42,449][INFO ][o.e.c.s.ClusterApplierService] [node-1] added [{node-2}{lmhdqVKZQuuMO6qsIELH8Q}{vSeL76bNTrC1I75EX4OumA}{127.0.0.1}{127.0.0.1:9301}{cdhilmrstw}{ml.machine_memory=34315870208, ml.max_open_jobs=20, xpack.installed=true, transform.node=true}], term: 797, version: 153, reason: Publication{term=797, version=153}
[2020-11-18T14:10:57,570][INFO ][o.e.c.r.a.AllocationService] [node-1] Cluster health status changed from [YELLOW] to [GREEN] (reason: [shards started [{.kibana-event-log-7.10.0-000001}[0]]]).
```

```
1 PUT /test
2
3 GET /_cluster/health
4
5 GET /_cat/shards?v
```

```
1 {
2   "cluster_name" : "elasticsearch",
3   "status" : "green",
4   "timed_out" : false,
5   "number_of_nodes" : 2,
6   "number_of_data_nodes" : 2,
7   "active_primary_shards" : 7,
8   "active_shards" : 14,
9   "relocating_shards" : 0,
10  "initializing_shards" : 0,
11  "unassigned_shards" : 0,
12  "delayed_unassigned_shards" : 0,
13  "number_of_pending_tasks" : 0,
14  "number_of_in_flight_fetch" : 0,
15  "task_max_waiting_in_queue_millis" : 0,
16  "active_shards_percent_as_number" : 100.0
17 }
```

index	shard	prirep	state	docs	store	ip	node
.kibana_task_manager_1	0	r	STARTED	5	122.9kb	127.0.0.1	node-2
.kibana_task_manager_1	0	p	STARTED	5	132.2kb	127.0.0.1	node-1
ilm-history-3-000001	0	r	STARTED			127.0.0.1	node-2
ilm-history-3-000001	0	p	STARTED			127.0.0.1	node-1
.apm-agent-configuration	0	r	STARTED	0	208b	127.0.0.1	node-2
.apm-agent-configuration	0	p	STARTED	0	208b	127.0.0.1	node-1
.kibana-event-log-7.10.0-000001	0	r	STARTED	3	16.4kb	127.0.0.1	node-2
.kibana-event-log-7.10.0-000001	0	p	STARTED	3	16.4kb	127.0.0.1	node-1
test	0	r	STARTED	0	208b	127.0.0.1	node-2
test	0	p	STARTED	0	208b	127.0.0.1	node-1
.apm-custom-link	0	r	STARTED	0	208b	127.0.0.1	node-2
.apm-custom-link	0	p	STARTED	0	208b	127.0.0.1	node-1
.kibana_1	0	r	STARTED	26	10.4mb	127.0.0.1	node-2
.kibana_1	0	p	STARTED	26	10.4mb	127.0.0.1	node-1





# ROLE UZLŮ

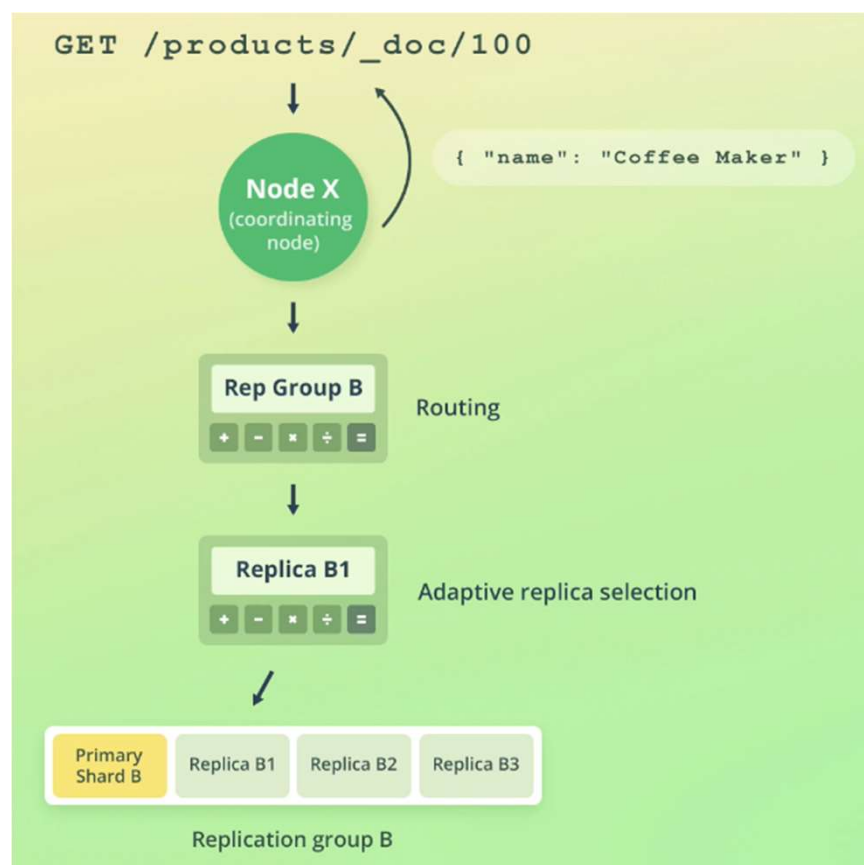
- možno definovat pro každý uzel až několik rolí
  - podle aplikace samotného uzlu
- uzly podle rolí
  - master
    - jeden, často dedikovaný
    - v případě více uzlů s touto rolí jen jeden je zvolen masterem
    - zodpovědnost za vytváření a mazání indexů, přiřazování shardů uzlům, ...
  - data
    - ukládání dat
    - dotazy nad daty
  - ingest
    - ingest pipeline – série kroků prováděná při indexaci dokumentů
    - přidání dokumentů do indexů
      - má právo dokumenty měnit

# ROLE UZLŮ

- uzly podle rolí
  - machine learning
    - možnost provádět úlohy strojového učení
  - coordination
    - distribuce dotazů a agregace výsledků
  - voting-only
    - účast na volbě nového master uzlu
    - nemůže být master

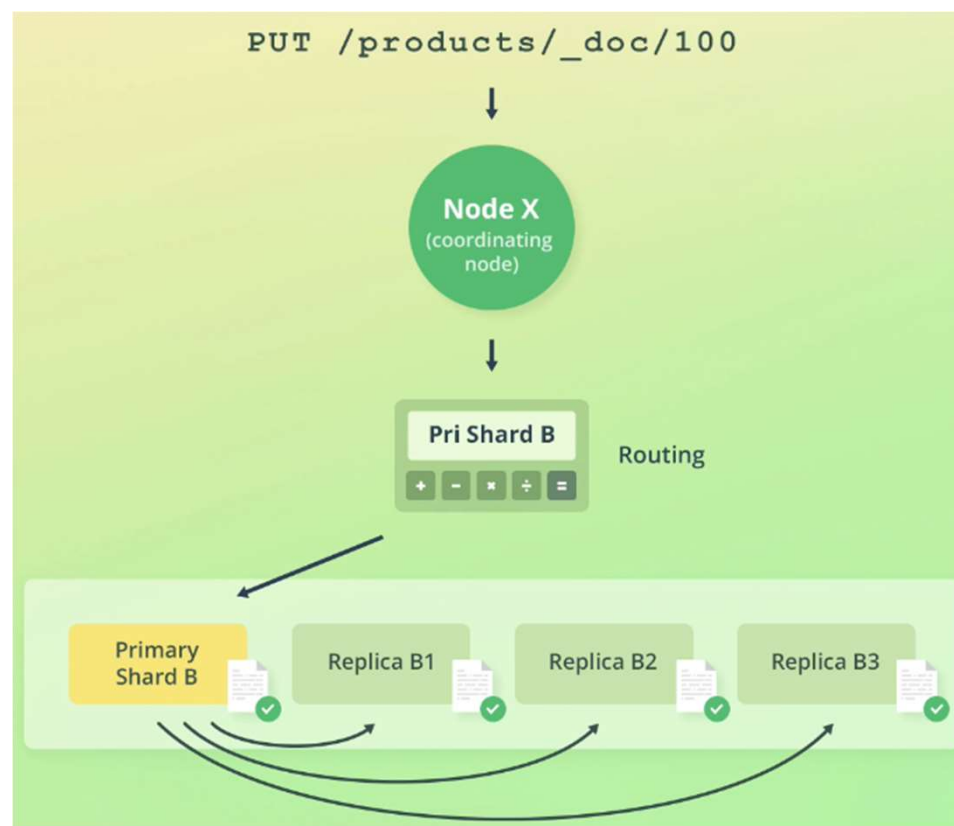
1	GET /_cat/nodes?v		1	ip	heap.percent	ram.percent	cpu	load_1m	load_5m	load_15m	node.role	master	name
2			2	127.0.0.1	21	43	8				cdhilmrstw	-	node-2
3			3	127.0.0.1	32	43	8				cdhilmrstw	*	node-1

# ČTENÍ DOKUMENTU PODLE ID



<https://www.udemy.com/course/elasticsearch-complete-guide/>

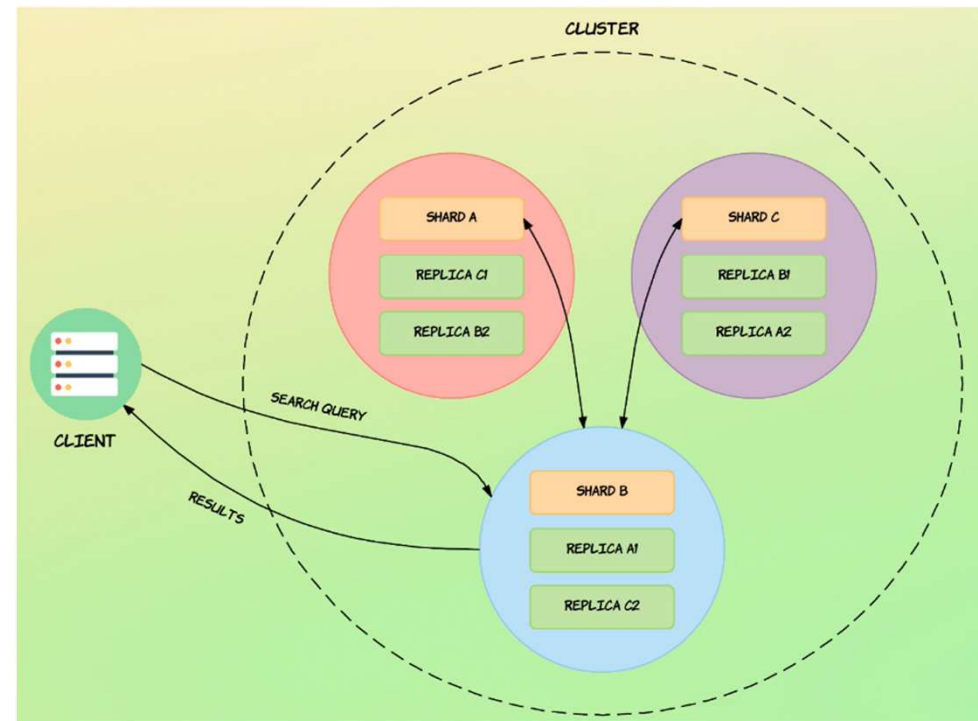
# ZÁPIS DOKUMENTU PODLE ID



<https://www.udemy.com/course/elasticsearch-complete-guide/>

# VYHLEDÁVÁNÍ

- princip
  - dotaz dorazí od klienta na jeden z uzlů (modrý)
  - ten se stává koordinačním uzlem
  - broadcastem přeposílá dotaz ostatním shardům
    - primárním i replikám
  - shardy odpovídají
  - koordinační uzel sloučí výsledky a vrací je klientovi
- liší se výrazně od čtení podle `_id`



<https://www.udemy.com/course/elasticsearch-complete-guide/>

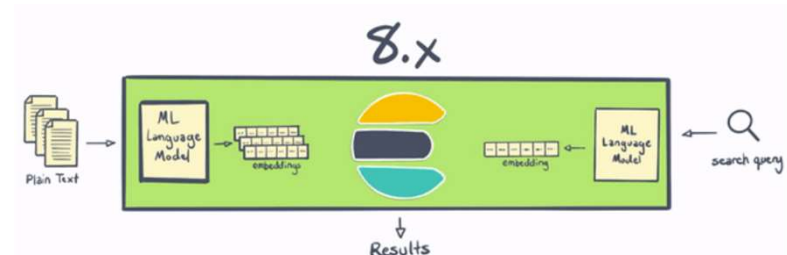
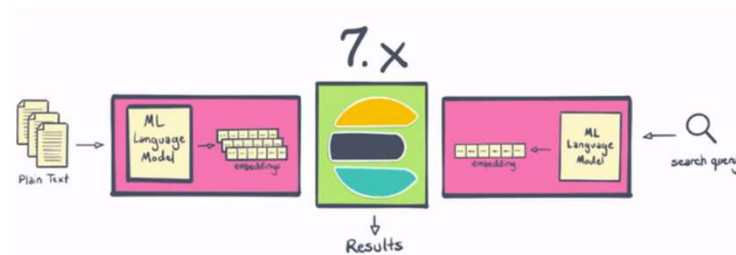
# ČÁST IV.: ELASTICSEARCH 8.0





# CO JE NOVÉHO?

- nová verze [8.0](#) vydaná 02/2022 přináší [různá vylepšení](#)
  - vyšší rychlost a menší paměťové nároky
  - nativní podpora pro zpracování přirozeného jazyka (NLP)
    - analýza sentimentu, rozpoznávání entit, klasifikace textu, ...
    - možnost přímého použití PyTorch modelů
  - vyšší relevance výsledků díky nativnímu vektorovému vyhledávání
  - vyšší rychlost a škálovatelnost vyhledávání
    - algoritmus approximate nearest neighbor search



<https://www.elastic.co/blog/whats-new-elastic-8-0-0>



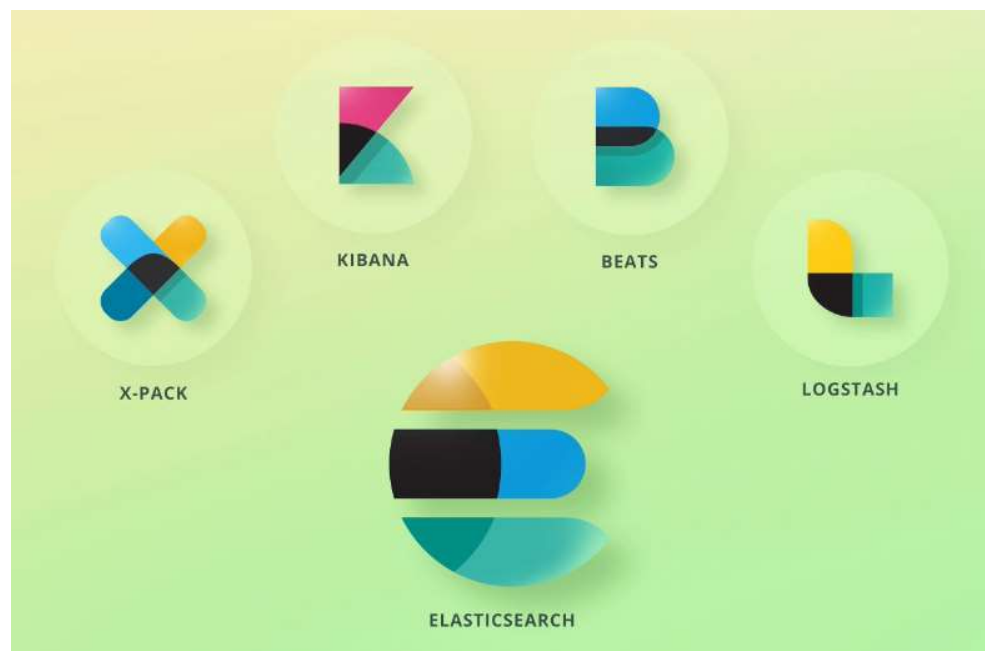
# CO JE NOVÉHO?

- nová verze [8.0](#) vydaná 02/2022 přináší [různá vylepšení](#)
  - vylepšená integrace AWS
  - vylepšená bezpečnost
    - autentizace, autorizace, šifrovaná komunikace (mezi uzly i mezi Elasticsearch a Kibana)
    - zapnuté a nastavené v základu
- informace ke [zpětné kompatibilitě](#) REST API
  - verze 8.0 přináší určité změny
    - ideální implementovat
  - ale nabízí také zpětnou kompatibilitu díky 7.x API compatibility headers
    - možné použít před přechodem na nové API
- výrazné problémy tedy nejsou

# ČÁST V.: ELASTIC STACK

# ELASTIC STACK

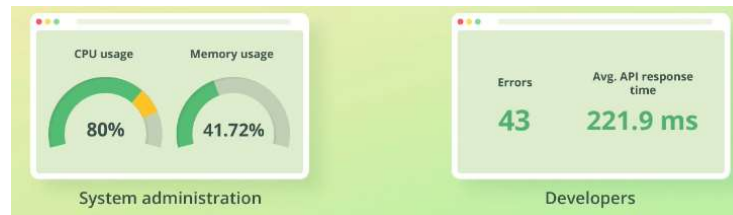
- technologie vyvinuté a spravované Elastic BV
  - silná synergie mezi technologiemi
    - často používané společně s Elasticsearch



<https://www.udemy.com/course/elasticsearch-complete-guide/>

# KIBANA

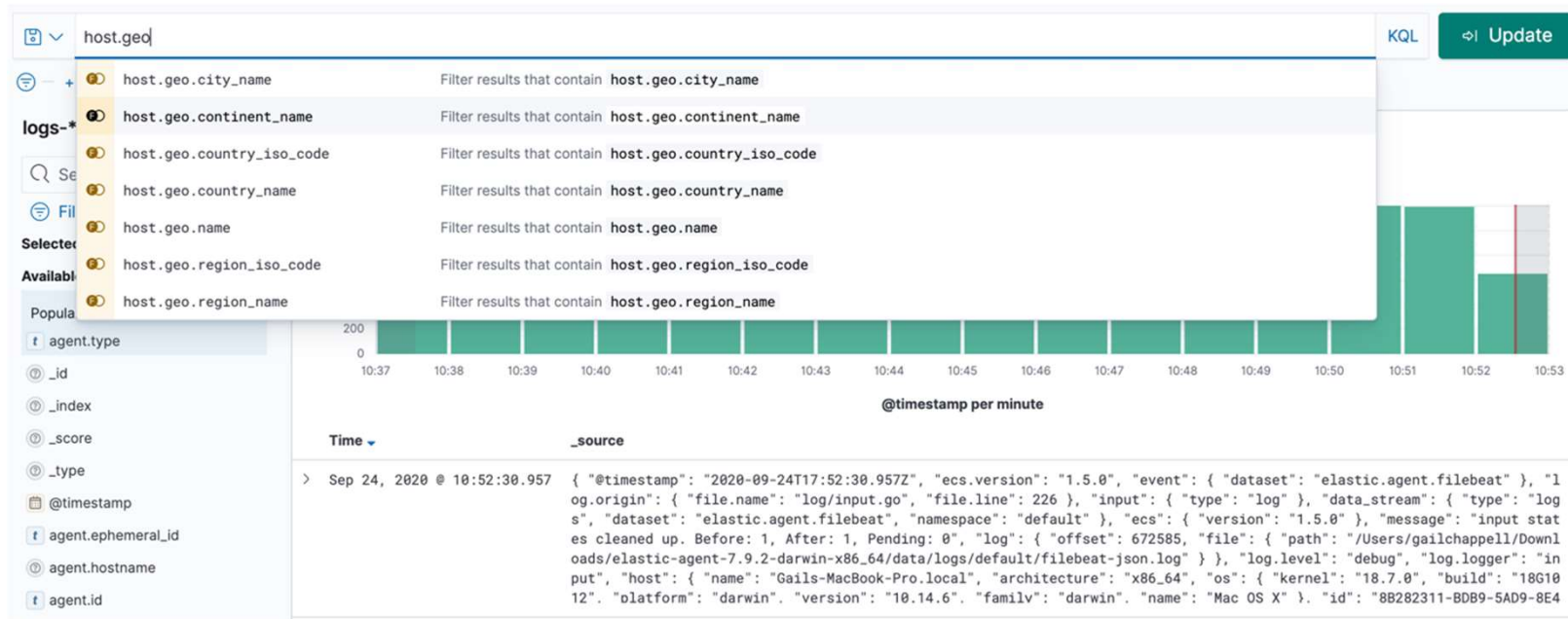
- platforma pro analýzu a vizualizaci dat v Elasticsearch
  - webový interface pro data uložená v Elasticsearch
  - vestavěná real-time vizualizace
    - grafy, mapy, ...
    - dashboards
  - nastavení strojového učení
    - predikce, detekce anomálií
  - může spravovat určité části Elasticsearch a Logstash
    - autentifikace, autorizace



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# KIBANA

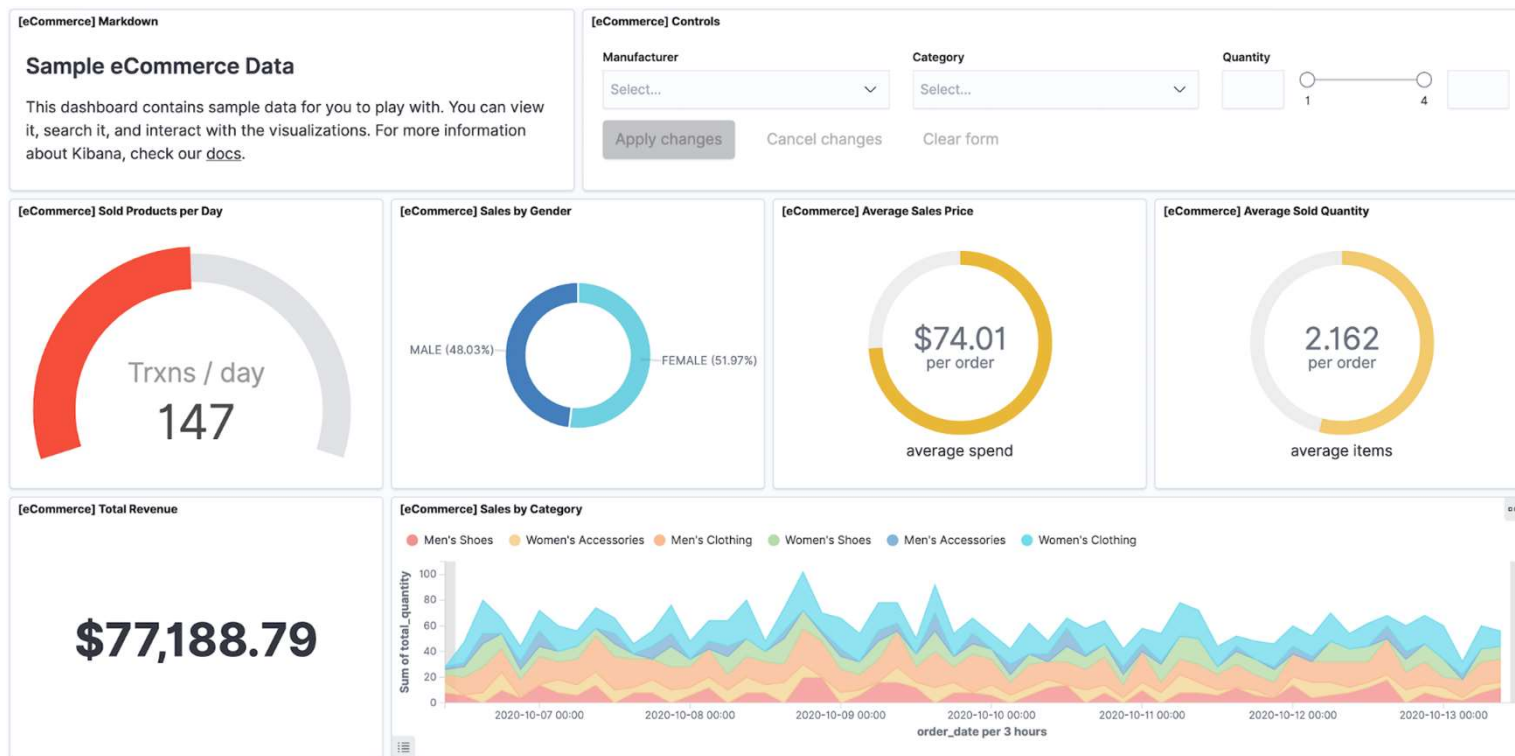


<https://www.elastic.co/guide/en/kibana/current/kuery-query.html>



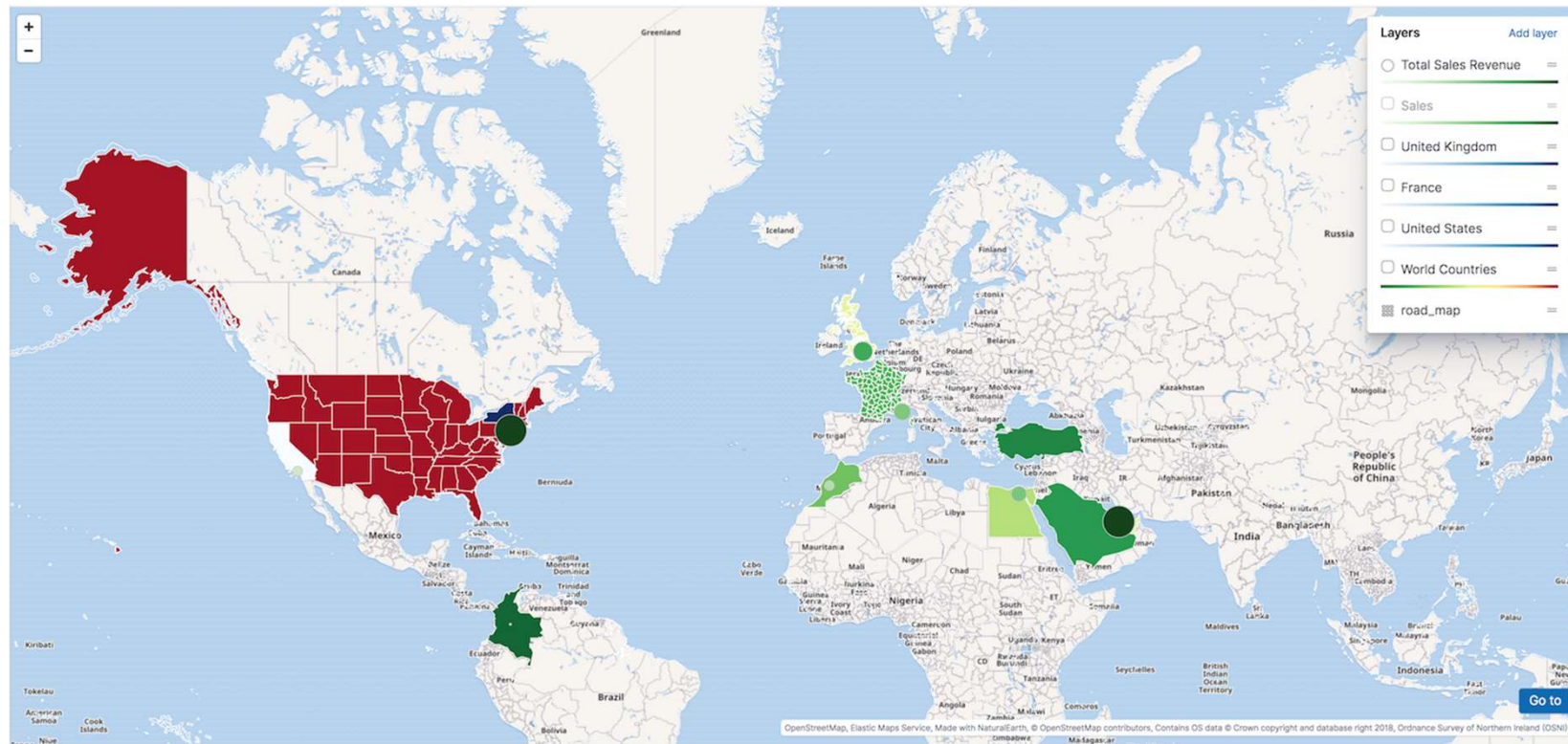


# KIBANA



<https://www.elastic.co/guide/en/kibana/current/dashboard.html>

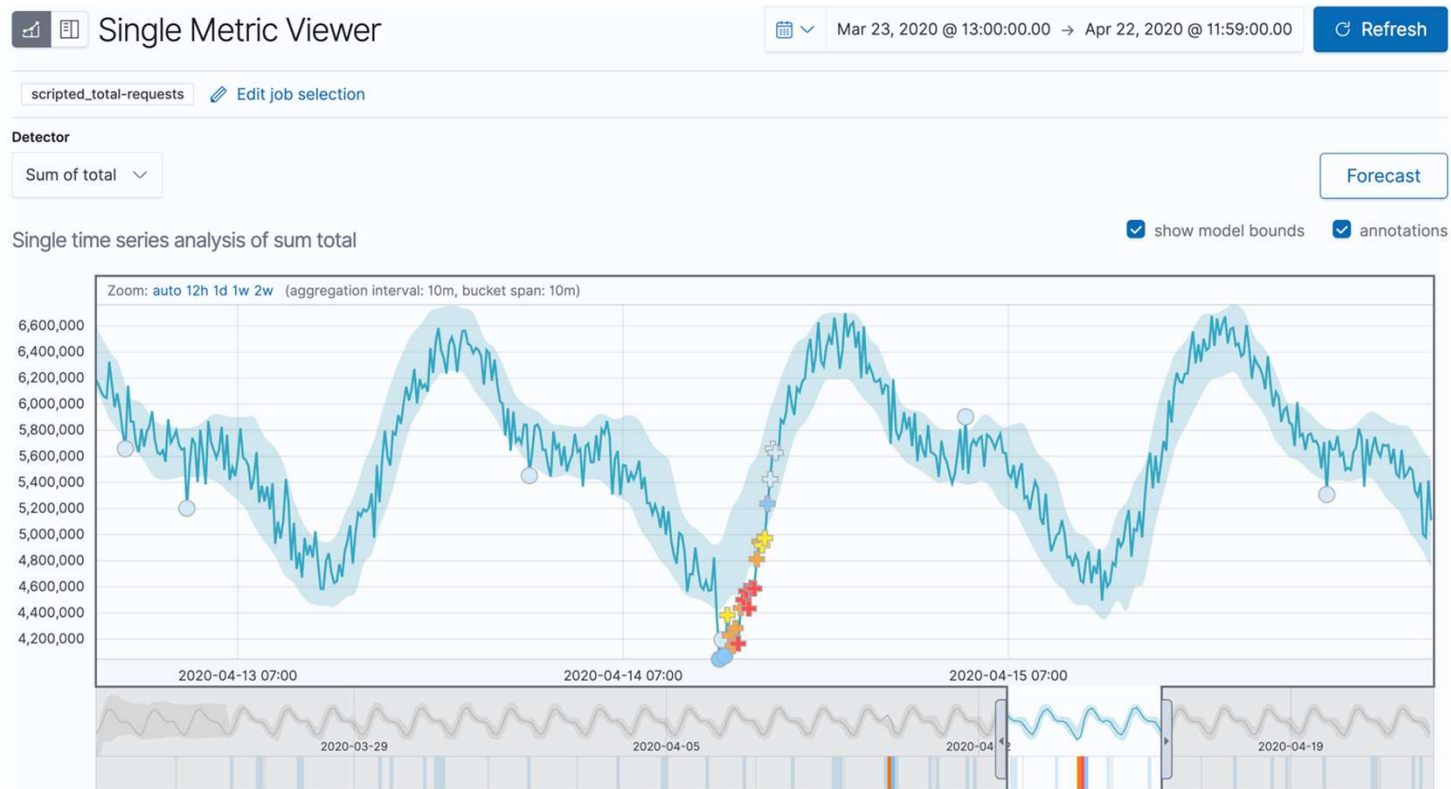
# KIBANA



<https://www.elastic.co/guide/en/kibana/current/maps.html>



# KIBANA



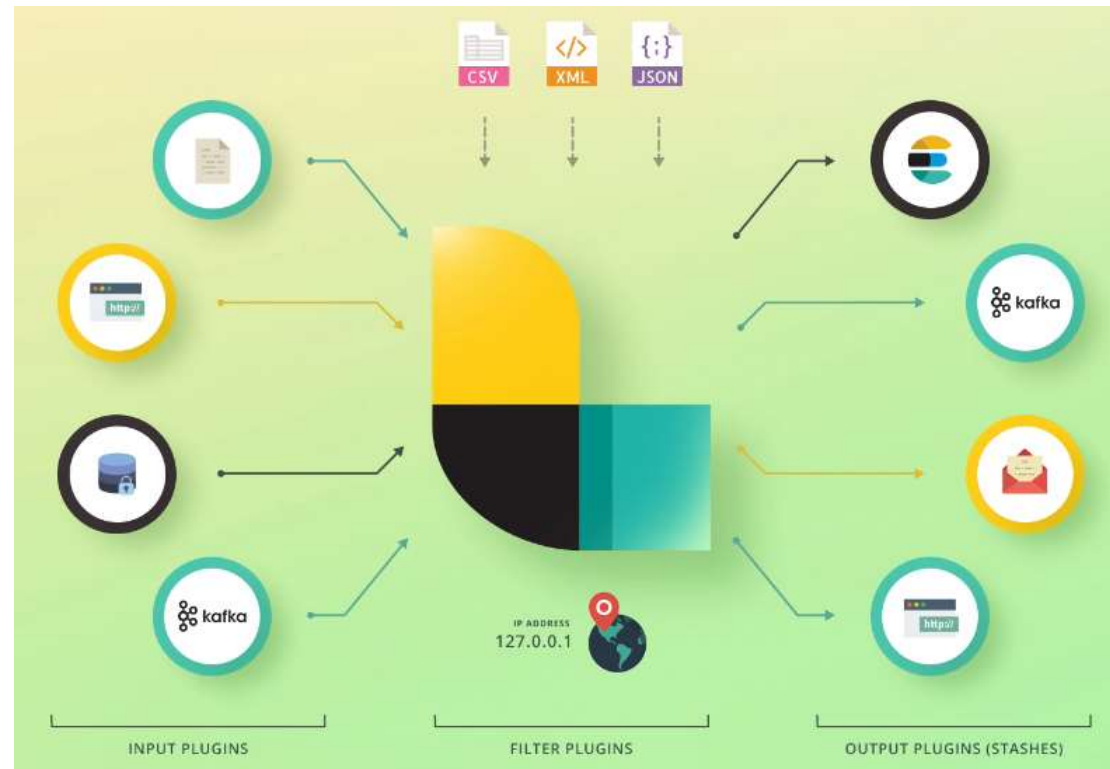
<https://www.elastic.co/guide/en/kibana/current/xpack-ml-anomalies.html>

# LOGSTASH

- původně nástroj na zpracování logů a předávání do Elasticsearch
- dnes roura (pipeline) na zpracování dat
  - obecnější nástroj
  - hlavní princip
    - data vstupují do Logstash a jsou považována za obecné události
      - mohou být cokoliv (serverové logy, objednávky, zákazníci, zprávy z chatu, ...)
    - Logstash události zpracovává a předává dál
      - může být kamkoliv (Elasticsearch, e-mail, web, ...)
- roura se skládá ze tří částí
  - vstupy, filtry, výstupy
  - každá část využívá pluginy pro práci s událostmi
    - vstupní pluginy – soubory, data přes HTTP, relační databáze, ...
    - filtrační pluginy – parsování XML a JSON, obohacení dat (IP adresa, geo-informace), ...
    - výstupní pluginy – stashes, místa, kam jsou události posílány dál
- horizontální škálování



# LOGSTASH



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# LOGSTASH

- roura na zpracování dat
  - definovaný vlastní značkovací formát (podobnost s JSON)
  - podporuje ale i podmíněné výrazy
    - dynamická roura

```
input {  
  file {  
    path => "/path/to/apache_access.log"  
  }  
}  
  
filter {  
  if [request] in ["/robots.txt", "/favicon.ico"] {  
    drop { }  
  }  
}  
  
output {  
  file {  
    path => "%{type}_%{+yyyy_MM_dd}.log"  
  }  
}
```

# LOGSTASH



<https://www.udemy.com/course/elasticsearch-complete-guide/>

# X-PACK

- balík funkcí do Elasticsearch a Kibana
  - zabezpečení
    - autentifikace a autorizace
    - integrace s poskytovateli autentifikací
    - nastavení uživatelských práv
  - monitoring
    - sledování výkonu Elastic Stack
    - využití procesorů, pamětí, disků, ...
    - upozornění v případě něčeho neobvyklého
  - upozornění
  - reporty
    - export dat (a vizualizace) z Kibany do různých formátů (CSV, PDF, ...)
    - generované na základě požadavků, podle plánu nebo po splnění určených podmínek
    - upravitelný vzhled



# X-PACK

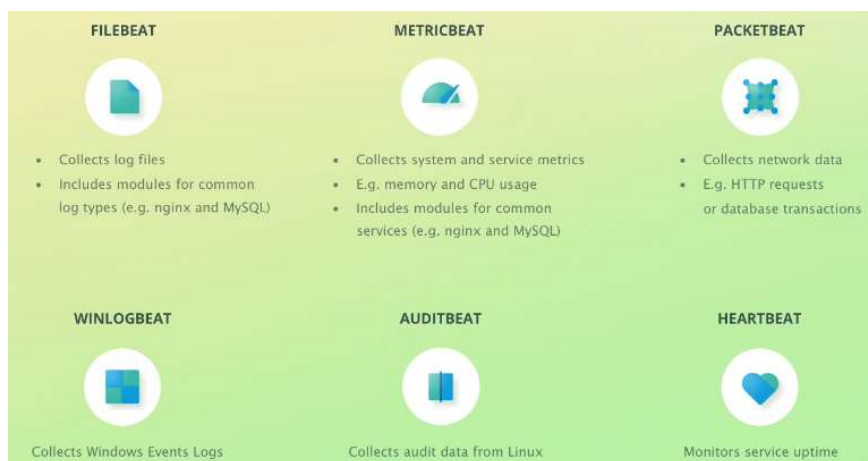
- balík funkcí do Elasticsearch a Kibana
  - strojové učení
    - funkcionalita dodána v X-Pack
    - interface v Kibaně
    - predikce
    - detekce anomálií
  - graph
    - analýza vazeb (vztahů) mezi daty
    - plugin do Kibany pro vizualizaci vazeb
  - Elasticsearch SQL
    - dotazování pomocí SQL
    - dotazy jsou běžně psány v Query DSL
    - překlad SQL do Query DSL na pozadí



<https://www.udemy.com/course/elasticsearch-complete-guide/>

# BEATS

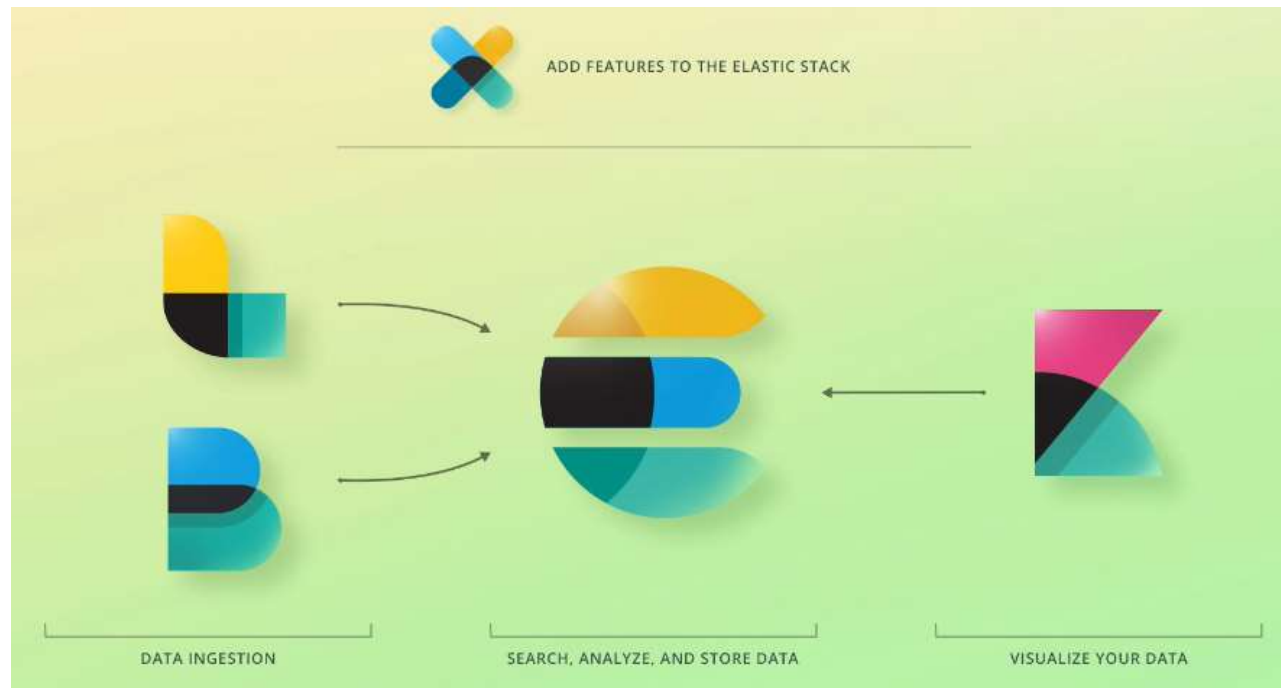
- kolekce data shippers
  - nenároční agenti instalovaní na servery
  - posílají data do Logstash / Elasticsearch
  - různé druhy
    - sběr různých dat
      - filebeat – sběr logů
      - metricbeat – systémové metriky



<https://www.udemy.com/course/elasticsearch-complete-guide/>



# ELASTIC STACK



<https://www.udemy.com/course/elasticsearch-complete-guide/>

## A PŘÍŠTĚ?

- sloupcové databáze
  - Apache Cassandra



Děkuji za pozornost.  
Otázky?