



# ElasticSearch для .NET

Татьяна Ёлкина  
Старший программист  
Интерфакс



## Обо мне:



**Старший инженер-  
программист  
Интерфакс**

 **@YOTata ytata.512@gmail.com**

 **facebook.com/yotata**



## О чем поговорим:

- Какие задачи решали?
- Какие есть варианты решения подобных задач?
- Что такое Elasticsearch? Когда и зачем создавался?
- Как это работает?
- Про запросы и анализаторы
- Стек ELK. Что это такое?
- NEST. Примеры использования



## Задача:

- Нужен быстрый и качественный поиск по постоянно растущему числу документов
- Нужна возможность изменения конфигурации “на лету”
- Нужно удобство использования ....



## Какие продукты есть сейчас

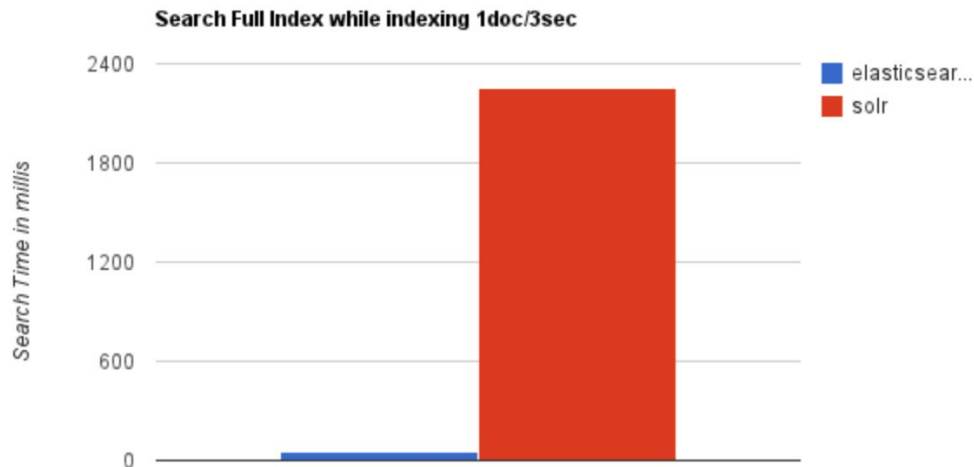
- Elasticsearch (Lucene, Java)
- Solr (Lucene, Java)
- Sphinx (C++)
- Xapian (C++)

## Почему выбрали Elasticsearch

- Поддержка шардинга и репликации на лету
- Написан на Java и базируется на Lucene
- Real-time индексация
- *Статья Ryan Sonnek*



# ES.Преимущества производительности

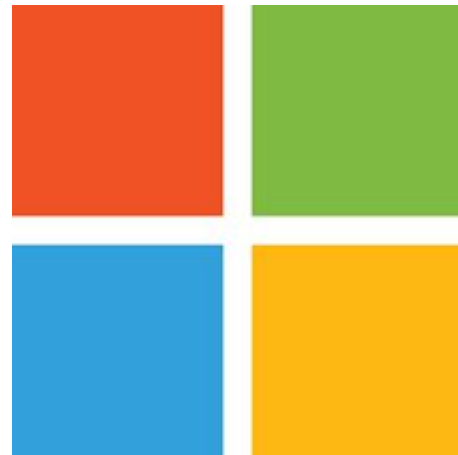


Кто использует ещё...

GitHub



UBER





## ElasticSearch. История



Shay Banon Создал в 2010.

2012 год создана компания ElasticSearch BV для платной поддержки пользователей.

В 2014 году привлечены инвестиции.

В 2015 изменение названия на Elastic

Актуальная версия 6.1.2

# ElasticSearch - основные преимущества

## Open Source

Полнотекстовый поиск

Распределенность

Отказоустойчивость

Документ-ориентируемость

Schema-free

Ответ в режиме реального времени



elastic

# ElasticSearch - основные преимущества

Open Source

**Полнотекстовый поиск**

Распределенность

Отказоустойчивость

Документ-ориентируемость

Schema-free

Ответ в режиме реального времени

1: Winter is coming.  
2: Ours is the fury.  
3: The choice is yours.



<u>term</u>	<u>freq</u>	<u>documents</u>
choice	1	3
coming	1	1
fury	1	2
is	3	1, 2, 3
ours	1	2
the	2	2, 3
winter	1	1
yours	1	3
Dictionary		Postings

# ElasticSearch - основные преимущества

Open Source

Полнотекстовый поиск

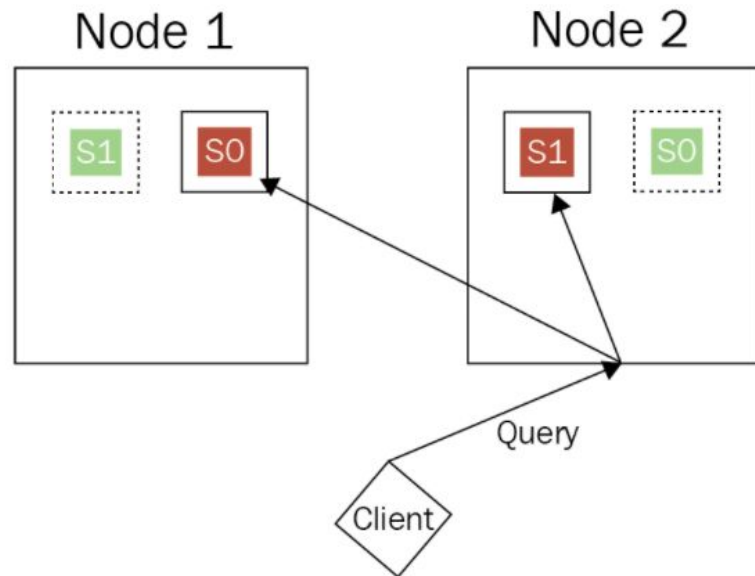
**Распределенность**

Отказоустойчивость

Документ-ориентируемость

Schema-free

Ответ в режиме реального времени





# ElasticSearch - основные преимущества

Open Source

Полнотекстовый поиск

Распределенность

**Отказоустойчивость**

Документ-ориентируемость

Schema-free

Ответ в режиме реального времени



# ElasticSearch - основные преимущества

Open Source

Полнотекстовый поиск

Распределенность

Отказоустойчивость

Документ-ориентируемость

Schema-free

Ответ в режиме реального времени

```
- _source: {  
  + OKTMOPPO: {...},  
  + ogrn: "1045300290947",  
  + loadId: 523,  
  + inn: "5321100197",  
  + accounts: {...},  
  + timeZoneOlson: "Europe/Moscow",  
  + organizationType: {...},  
  + timeZoneUtcOffset: "UTC+03:00",  
  + factualAddress: {...},  
  + fileVersion: "20171029000002_065",  
  + contractsCount: 230,  
  + contractsSum: 170859935.7299998,  
  + regionCode: "53",  
  + postalAddress:  
    "Российская Федерация, 173007, Новгородская обл, Великий Новгород г, ул  
  + headAgency: {...},  
  + IKUInfo: {...},  
  + registrationDate: "2004-11-24",  
  + okogu: {...},  
  + fax: "721-8162-662811",  
  + shortName: "ГЛАВНОЕ УПРАВЛЕНИЕ МЧС РОССИИ ПО НОВГОРОДСКОЙ ОБЛАСТИ",  
  + fz: "44",  
  + contractsYearStats: {...},  
  + contactPerson: {...},  
  + phone: "721-8162-662794",  
  + okopf: {...},  
  + fullName: "ГЛАВНОЕ УПРАВЛЕНИЕ МИНИСТЕРСТВА РОССИЙСКОЙ ФЕДЕРАЦИИ ПО ДЕЛАМ  
  kpp: "532101001",  
  + orderingAgency: {...},  
  + actual: "true",  
  + organizationRole: [...],  
  + register: "true",  
  + ОКПО: "08929072",  
}
```



# ElasticSearch - основные преимущества

Open Source

Полнотекстовый поиск

Распределенность

Отказоустойчивость

Документ-ориентируемость

**Schema-free**

Ответ в режиме реального времени

# ElasticSearch - основные преимущества

Open Source

Полнотекстовый поиск

Распределенность

Отказоустойчивость

Документ-ориентируемость

Schema-free

**Ответ в режиме реального времени**







## ElasticSearch - терминология

- **Нода** - процесс JVM, запущенный на каком-то сервере
- **Индекс** - набор документов
- **Шард** - часть индекса
- **Реплика** - копия шарда



## ElasticSearch - терминология. Цвета кластера

- **Красный** - часть индекса недоступна вообще
- **Желтый** - какие-то реплики либо находятся в состоянии миграции, либо не прикреплены к нодам
- **Зеленый** - доступно требуемое количество реплик каждого шарда, каждого индекса



# ElasticSearch - терминология

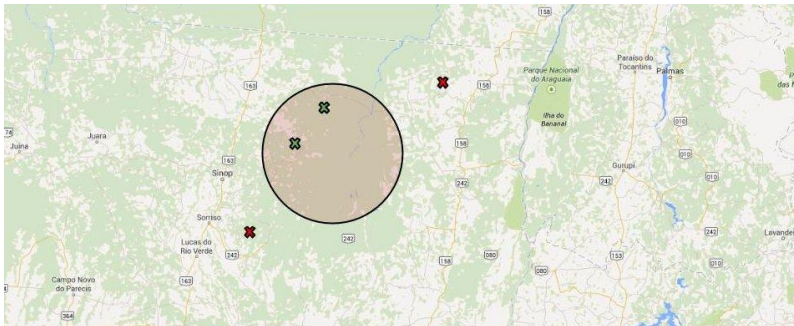
- **Стемминг** — это нахождение основы слова для заданного исходного слова. Основа необязательно совпадает с морфологическим корнем слова.
- **Лемматизация** — приведение слова к нормальной (словарной) форме. Для существительных это именительный падеж и единственное число.
- **Корпус** — в лингвистике корпусом называется совокупность текстов, собранных в соответствии с определенными принципами, размеченных по определенному стандарту и обеспеченных специализированной поисковой системой. Это может быть и разделение по стилям и жанрам, разделение по эпохе написания, по форме написания.
- **Параллельный корпус** — это один или более текстов на двух языках, сопоставленные между собой парами, когда в каждой паре оба предложения несут один и тот же смысл.
- **Стоп-слова, или шумовые слова**, — предлоги, суффиксы, междометия, цифры, частицы и подобное. Общие шумовые слова всегда исключаются из поискового запроса (кроме поиска по строгому соответствию поисковой фразы), также они игнорируются при построении инвертированного индекса.
- **N-грамма** — последовательность из n элементов. С семантической точки зрения это может быть последовательность звуков, слогов, слов или букв.



## ElasticSearch. Запросы

- Обычные запросы (Term, Match/Not Match, Prefix, Bool, Wildcard, Range)
- Географические запросы (Geo-Shape, Geo-Bounding, Geo-Distance, Geo-Polygon)
- Boosting запросы
- Span запросы
- Нечеткие запросы
- Aggregations
- Nested Queries

# ElasticSearch. Geo-Distance



```
{
  "filtered" : {
    "query" : {
      "field" : { "text" : "restaurant" }
    },
    "filter" : {
      "geo_distance" : {
        "distance" : "12km",
        "pin.location" : {
          "lat" : 40,
          "lon" : -70
        }
      }
    }
  }
}
```



## ES. Примеры запросов

GET /\_search

```
{  "query": {  
    "query_string" : {  
        "fields" : ["content", "name"],  
        "query" : "this AND that"  
    }  
}}
```

```
{
  took: 6,
  timed_out: false,
  - _shards: {
    total: 5,
    successful: 5,
    failed: 0,
  },
  - hits: {
    total: 1448930,
    max_score: 1,
    - hits: [
      - {
        _index: "suppliers",
        _type: "2017-11-17",
        _id: "AV_JHUSOYDJB_ZMeRHC7",
        _score: 1,
        - _source: {
          organizationName: "СИМИКЯ",
          - contracts223YearStats: {
            - 2016: {
              contractsCount: 1,
              contractsSum: 1236
            }
          },
          inn: "761602991051".
        }
      }
    ]
  }
}
```

## ES. Пример ответа



# ElasticSearch. Анализаторы

Изнутри каждый анализатор представляет собой своеобразный конвейер, состоящий из нескольких обработчиков:

- Символьной фильтрации
- Токенизации
- Фильтрации полученных токенов





# ElasticSearch. Анализаторы

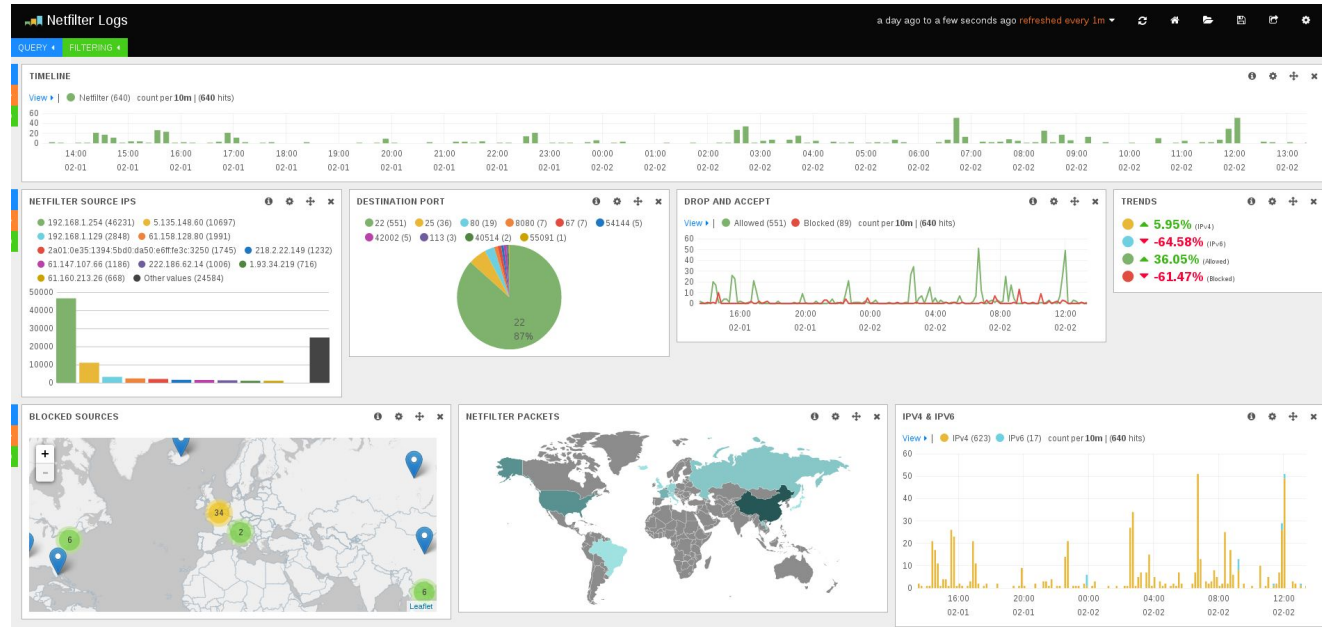
```
"settings": {  
  "analysis": {  
    "analyzer": {  
      "case_insensitive": {  
        "tokenizer": "keyword",  
        "filter": ["lowercase"]  
      }  
    }  
  }  
}
```

## Стек ELK

- ElasticSearch
- Logstash (Inputs(rabbitmq, syslog, file, redis), Filters, Outputs(file, elasticsearch..))
- Kibana
- ...



# Kibana

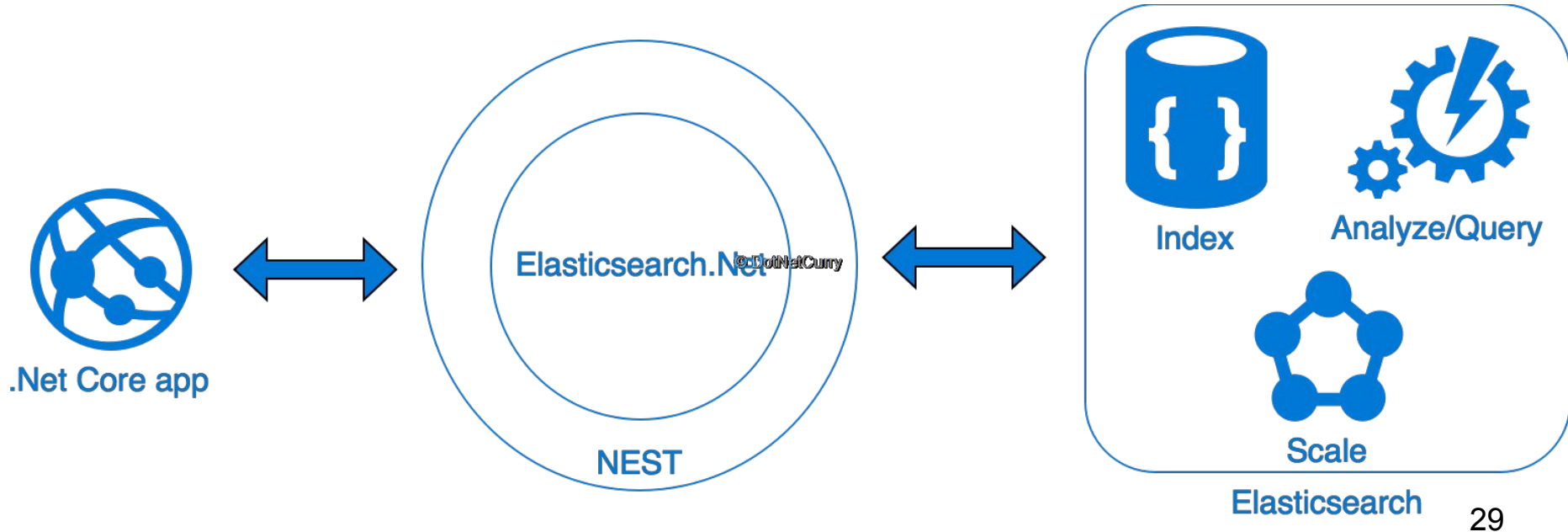




## NEST - ...

“NEST is a high level client that has the advantage of having mapped all the request and response objects, comes with a strongly typed query DSL that maps 1 to 1 with the Elasticsearch query DSL, and takes advantage of specific .NET features such as covariant results. NEST internally uses, and still exposes, the low level Elasticsearch.Net client.” [www.elastic.co/](http://www.elastic.co/)

# NEST. Cxema





**Демо**



## Эпилог.

