

# How to start a logging project with the Elastic Stack

Marco De Luca, Principal Solution Architect Köln@ REWE, 8. October 2019



### Please fill out the survey and WIN a Elastic Backpack!

Here is a link to the survey:

https://go.es.io/2I07X5D

Or a QR code:







### Agenda

- A very quick Elastic Stack Overview
- Step by step A suggestion on how to start
- The important foundation a common schema
- Index Templates and how to use them
- Preparing a cluster for scale
- Index Lifecycle Management
- Demo
- Q & A



## Elastic Stack Overview

•

As short as possbile 🕲







#### **Elastic Licensing Options**

|              | Free Offerings          |                              | Commercial Offerings                                  |  |
|--------------|-------------------------|------------------------------|---|--|
| SELF-MANAGED | OPEN SOURCE             | BASIC                        | Stand-alone   | ENTERPRISE                                     |
|              | Open Source<br>Features | Free proprietary<br>Features | Commercial proprietary<br>Features<br>Elastic Support | Elastic Cloud<br>Enterprise<br>Elastic Support |
|              |                         |                              |   |  |

**Commercial SaaS Offerings** 

SaaS

**Elasticsearch Service** 

Elastic App Search Service

**ELASTIC CLOUD** 

**Elastic Site Search Service** 



#### **Elastic Licensing Options**

|              | Free Offerings            |                              | Commercial Offerings                  |  |  |  |
|--------------|---------------------------|------------------------------|---------------------------------------|--|--|--|
| SELF-MANAGED | OPEN SOURCE               | BASIC                        |                                       |  |  |  |
|              | Open Source<br>Features   | Free proprietary<br>Features | Commercial pr<br>Featur<br>Elastic Su |  |  |  |
|              | Commercial SaaS Offerings |                              |                                       |  |  |  |
| SaaS         |                           |                              |                                       |  |  |  |
|              |                           |                              |                                       |  |  |  |



#### STANDALONE

#### **Elastic On-Premise**



See full list of features

**Download Basic** 

Download Open Source Software



## Step by Step

0

0

•

•

-

High Level approach

#### High Level Steps (some steps, might not be complete!)

- I. Define your ETL Pipelines
- II. Build your cluster, index templates, ILM, etc.
- **II. Build Queries and Visualizations**
- III. Configure Alarms + Machine Learning Jobs\*
- IV. Connect external Systems via REST / API



### **ETL pipeline building**

#### A bit more details

- 1. Define/Collect Data Sources (maybe the low hanging fruits?)
- 2. Define how to get the data from the source (Beats, Logstash, API, etc.)
- 3. Define your Common Schema (**ECS** is a great start)
- 4. Not a loved task, but document the above :-)
- 5. If you have data sources that cannot be easily collected you can ...
  - ... build a transformation and write a beats module for it,
- 6. Make an educated guess on the amount of data you need to collect, what is the retention of it, number of events per second/minute or day, avg. size of document



#### **Cluster Sizing and building**

- 1. Size your Elasticsearch cluster (Elastic contacts or the community can help!)
- 2. Build your templates and define your Index Lifecycle strategy
- 3. Apply your templates
- 4. Build your cluster and connect the data sources



#### Use existing Dashboards or build your own

- Beats Dashboards --> Can be installed using e.g. # filebeat setup -- dashboards
- Infra App + Logging App, Updatime --> Pre-build Applications for Infrastructure monitoring
- Self-Made Visualizations for everything not pre-build



#### **Configure Alarming and Machine Learning Jobs\***

- Define and document alarms
- Define and document ML Jobs
- Create Alarms and ML Jobs

\* Part of paid gold/platinum subscription

#### **Need to connect external Systems**

Define your interfaces

- External System → REST Interface available?
- Other interfaces possbile?
- Use ES and alarms or Logstash?



#### From Source to Analytics – A Pipeline Example





#### From Source to Analytics – A Pipeline Example





## The important foundation

A "common schema" or "The Elastic Common Schema"

#### What is the Elastic Common Schema (ECS)?

ESC is a new specification that provides a consistent and customizable way to structure your data in Elasticsearch, facilitating the analysis of data from diverse sources.



#### Ingest almost anything, from almost anywhere





#### What to make of data from all over?

| Domain      | Data Sources                     | Timing                    | <b>Collection Methods</b>                         |
|-------------|----------------------------------|---------------------------|---|
| Network     | NetFlow,<br>PCAP, Zeek           | Real-time, Packet-based   | Filebeat, Packetbeat,<br>Logstash NetFlow module  |
| Application | Log                              | Real-time, Event-based    | Filebeat<br>Logstash                              |
| Cloud       | API, Log                         | Real-time, Event-based    | Beats<br>Logstash                                 |
| Host        | Signature Alert,<br>System State | Real-time, Asynchronous   | Auditbeat, Winlogbeat,<br>Filebeat Osquery module |
| Active      | Scanning                         | User-driven, Asynchronous | Vulnerability scanners                            |



#### Normalize data

10.42.42.42 - - [07/Dec/2018:11:05:07
+0100] "GET /blog HTTP/1.1" 200 2571 "-"
"Mozilla/5.0 (Macintosh; Intel Mac OS X
10\_14\_0) AppleWebKit/537.36 (KHTML, like
Gecko) Chrome/70.0.3538.102 Safari/537.36"



#### Normalize data with Elastic Common Schema (ECS)

| Searching <i>without</i> ECS   | Searching with ECS    |
|--|-----------------------|
| <pre>src:10.42.42.42 OR client_ip:10.42.42.42 OR apache2.access.remote_ip:     10.42.42.42</pre> | source.ip:10.42.42.42 |
| OR context.user.ip:10.42.42.42<br>OR src_ip:10.42.42.42  |                       |



#### Uniform analysis, no matter the data source



#### How ECS is structured

| Field Sets           |  | odit |
|----------------------|--|------|
| Field Set            | Description  |      |
| Base                 | All fields defined directly at the top level                                 |      |
| Agent                | Fields about the monitoring agent.   |      |
| Autonomous<br>System | Fields describing an Autonomous System (Internet routing prefix).            |      |
| Client               | Fields about the client side of a network connection, used with server.      |      |
| Cloud                | Fields about the cloud resource.   |      |
| Container            | Fields describing the container that generated this event.                   |      |
| Destination          | Fields about the destination side of a network connection, used with source. |      |
| DNS                  | Fields describing DNS queries and answers.                                   |      |
| ECS                  | Meta-information specific to ECS.  |      |

#### **Container Fields**

Field



type: keyword



#### How to use ECS

and define additional new fields

- Current 1.2 version of ECS has almost 400 fields
- You do **NOT** use all
- Beats use ECE per default!  $\rightarrow$  NO work todo!
- Fields not available  $\rightarrow$  Define your own and do the mapping
- Example ECS with User Defined Fields



## Index Templates

... and how to use them

•

### ECS brings its index template for fields definition, but

- need to add settings
- Important settings:
  - number\_of\_shards → Primary shards
  - number\_of\_replicas → Replicas
  - Dynamic → Set to False
- Shard number and size depends on workload, e.g.
  - − Logging  $\rightarrow$  30-60GB per shard
  - Search → 5-20GB per shard
- API to shrink and merge shards
  - used to consolidate shards



#### No ECS – Adjust the Dynamic setting

```
curl -XPUT localhost:9200/_template/template_1
{
    "index_patterns" : ["bar*, te*"],
    "order" : 0,
    "settings" : { ...},
    "mappings" : { ...},
    "mappings" : { ...}
    "docs" : {
        "dynamic": false,
        "properties": {...}
      }
    }
}
```



### **Dynamic template settings**

Elasticsearch Template

- "dynamic": true
  - Newly detected fields are added to the mapping. (default)
- "dynamic": false
  - Newly detected fields are ignored.
- "dynamic": strict
  - If new fields are detected, an exception is thrown and the document is rejected. New fields must be explicitly added to the mapping.



## Preparing a cluster for scale

Some important considerations

### Scaling the cluster using shard settings



```
Write operation: 1x
Read operation: 2x
```

Max Index Size (Logging Workload): approx. 60GB



### Scaling the cluster using shard settings

```
Node 2
                                          Node 1
                                                                                   Node 3
PUT _template/template_1
  "index patterns": ["te*", "bar*"],
  "settings": {
    "number of shards": 3,
    "number of replicas": 1
                                         P0
                                              R2
                                                              R0
                                                                   P1
                                                                                  P2
                                                                                        R1
  },
  "mappings": { ... }
```

Write operation: 3x Read operation: 6x

34

Max Index Size (Logging Workload): approx. 3 x 60GB



#### How to scale further?

Many way, it all denpends on the workload

- add more RAM (if you are not using the maximum per node)
- add more nodes
- scale your indexes by adding shards
- use the roll over API
- use ILM to free up resources (use HOT, WARM, COLD)
- When loading large amount of data, delete replicas first
- and **DO NOT** forget the hardware, especially disk I/O
- De-normalize your data



#### What to avoid?

- TOO many shards and indexes
  - Often comes from small indexes rolled over every day
  - If daily ingest is in MB or small GB numbers, roll over every week or month with that index.
  - Different index, different treat!
- Slow I/O
- Trying to mimic a SQL database schema --> De-normalization is key



#### Hardware Infrastructure

Even cloud is based on hardware! ;-)

- Considerations that satisfy the majority of use cases
  - Disk/Storage
    - SSD internal / SAN Allflash volumes → HOT Data
    - SAS / SATA Drives / SAN Volumes (SAS/SATA) → WARM Data
    - SATA Drives / Volumes → COLD Data
  - Memory
    - Scale over nodes to improve queries by caching more data (Search)
    - RAM to Disk Rations: 1:30 for Hot, 1:200 for Warm and 1:500+ for Cold Data
  - CPU
    - Compute intensive queries, Ingest Pipelines, Machine Learning Jobs
    - Rule of thump, 4-8 cores per 64GB RAM
    - ECE → 16 cores per 128GB RAM



## Index Lifecycle Management

Good by curator – Welcome ILM policy

### **Rollover API**

Top things off without spilling over

```
# Add > 1000 documents to logs-000001
POST /logs_write/_rollover/new-index-name
```

```
"conditions": {
    "max_age": "7d",
    "max_docs": 1000,
    "max_size": "5gb"
}
```





#### **Even more goodness – Frozen Indices**

In cases Warm has cooled down enough  $\odot$ 

#### • Frozen Indices

- Data you only search once in a while
- Great for e.g. archived or forensic data
- Index is still be searchable, unlike Snapshots
- Index is closed and will be opened for searches
- Not supposed to be for high query load
- Data needs to be mapped to memory at query time
- Benefits
  - Much higher disk to JVM heap ratios possible (1:500+)
  - Does not require Java Heap for its transient shared memory and others are moved to persistent storage
  - Will "survive" upgrades, unlike snapshots!



#### **Even more goodness - ILM**

#### Goodby curator, welcome ILM

- Index Lifecycle Management
  - Policy based rollover, delete with phases for \_ hot, warm, cold

Aanaate index

Name

logstanh-0 Inextant-1

umi-gcp-pub

apm-6.5.4-onbox

apro-6.5.4 span-1

um pro-pubsub

Delete index

Add lifecycle policy

my index

index\_text1

Rows per page: 10 14

logstanh-7

bile

- Created in the Management Section of \_ Kibana
- or out of Index Managemen

#### Name Policy name my-ilm-policy A policy name cannot start with an underscore and cannot contain a question mark or a space. Hot phase Anw This phase is required. You are actively querying and writing to your index. X Enable rollover For faster updates, you can roll over the index when it gets too big or too old. The new index created by rollover is added to the index alias and designated as the write index. Learn about rollover Index priority Index priority (optional) Set the priority for recovering your indices after a node restart. Indices with 100 higher priorities are recovered before indices with lower priorities. Learn more Warm phase You are still querying your index, but it is read-only. You can allocate shards to less performant hardware. For faster searches, you can reduce the number of shards and force merge segments. X Activate warm phase Index management old phase Update your Elasticsearch indices individually or in bulk. ou are querying your index less frequently, so you can allocate shards on INDEX OPTIONS ignificantly less performant hardware. Because your queries are slower, you an reduce the number of replicas. Show index settings Status Activate cold phase Show index mapping Show index stats 0.041 open belete phase Edit index settings. open bu no longer need your index. You can define when it is safe to delete it. Cince index open. Force merge index X Activate delete phase Refrech index



Show (SON

#### Create an index lifecycle policy

Use an index policy to automate the four phases of the index lifecycle, from actively writing to the index to deleting it. Learn about the index lifecycle.

Clear index cache Save as new polic Cancel 0.041 Flush index open Freeze index

0041



•

c

.

#### What you see is what you get!

## **Questions?**

The floor is yours.

## **ECS Resources:**

Blog post on migrating Beats data to ECS elastic.co/blog/migrating-to-elastic-common-schema-in-beats-environments

Webinar on migrating data to ECS elastic.co/webinars/introducing-the-elastic-common-schema

**Technical documentation for ECS** 

elastic.co/guide/en/ecs/current/index.html

GitHub repo for ECS

github.com/elastic/ecs



### Please fill out the survey and WIN a Elastic Backpack!

Here is a link to the survey:

https://go.es.io/2I07X5D

Or a QR code:









+



## Thanks You!

