# EPL660: Information Retrieval and Search Engines – Lab 8

Παύλος Αντωνίου

Γραφείο: B109, ΘΕΕ01

**University of Cyprus
Department of
Computer Science**

# Hands on ElasticSearch

- Elasticsearch v7.8.0 installed on VM

- Kibana installed on VM

- Python client libraries for Elasticsearch installed
  - elasticsearch
    - more general but hides less the complexities of the API calls
  - elasticsearch-dsl
    - focused on the search capabilities and is more friendly for sending queries to ElasticSearch

- Activate Elasticsearch
  - sudo service elasticsearch start

- Activate Kibana
  - sudo service kibana start

# Hands on ElasticSearch

- Install Elasticsearch on Windows
  - Download zip via https://www.elastic.co/guide/en/elasticsearch/reference/current/zip-windows.html

- Unzip and run \bin\elasticsearch.bat to start ES



```
Command Prompt - elasticsearch.bat                                    —  □  ✕

Microsoft Windows [Version 10.0.19041.572]
(c) 2020 Microsoft Corporation. All rights reserved.

C:\Users\Pavlos>cd Downloads

C:\Users\Pavlos\Downloads>cd elasticsearch-7.9.3

C:\Users\Pavlos\Downloads\elasticsearch-7.9.3>cd bin

C:\Users\Pavlos\Downloads\elasticsearch-7.9.3\bin>elasticsearch.bat
```

- Python libraries (if anaconda is in place):
  - `conda install -c conda-forge elasticsearch`
  - `conda install -c conda-forge elasticsearch-dsl`

# Hands on ElasticSearch

- Check if Elasticsearch is working:
  - Run elasticsearch_test.py file in Spyder or Python IDLE
    - http://localhost:9200

# Hands on ElasticSearch

- Check cluster health:

  - http://127.0.0.1:9200/_cat/health?v



```
epoch        timestamp cluster       status node.total node.data shards pri relo init unassign pending_tasks max_task_wait_time active_shards_percent
1597588373 14:32:53   elasticsearch green           1         1      6   6    0    0        0             0                  -               100.0%
```

  - Elasticsearch provides a handy "traffic lights" classification of cluster health:

  🔴 RED: Some or all of (primary) shards are not ready

  🟡 YELLOW: Elasticsearch has allocated all of the primary shards, but some/all of the replicas have not been allocated. Your cluster is fully operational.

  🟢 GREEN: Elasticsearch is able to allocate all shards and replicas to machines within the cluster.
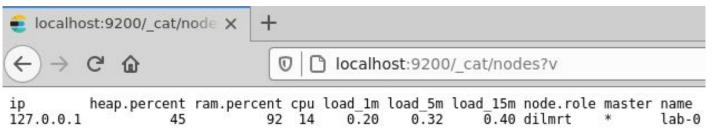
# Hands on Kibana

- Check if Kibana is working:
  - http://localhost:5601

# RESTful API Calls

- Access ElasticSearch via Restful API on browser
  - View nodes: http://127.0.0.1:9200/_cat/nodes?v



  - View all indices: http://127.0.0.1:9200/_cat/indices?v
  - View shards: http://127.0.0.1:9200/_cat/shards?v
  - View segments: http://127.0.0.1:9200/_cat/segments?v

# Today's lab

- Datasets
  - `20_newsgroups`: Text from 20 usenet groups on various topics, a classic corpus in IR evaluation, from [here](#).
  - `novels`: A number of random novels and other texts in English from the Gutenberg project, with a tendency towards late 19th and early 20th centuries.

# Today's lab

- Download lab8.zip and unzip it
- Create folder e.g. `/home/ubuntu/datasets`
- Move 20_newsgroups.tar.gz  and novels.zip into the datasets folder and unzip them
  - `tar xzvf 20_newsgroups.tar.gz`
  - `unzip novels.zip`
- In this lab:
  - You will learn how to use the ElasticSearch database
  - How to index a set of documents
  - How to ask simple queries about indexed documents
- Go through Lab8.pdf to run the examples
- Submit results to Moodle by Nov. 19 @ 15:00