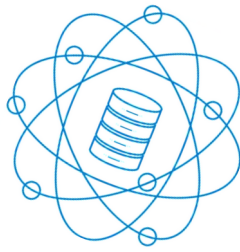


# The BDP Infrastructure for monitoring and analysing the ATLAS experiment processing activities at INFN-CNAF Tier-1

Giacomo Levrini, University of Bologna [giacomo.levrini@studio.unibo.it](mailto:giacomo.levrini@studio.unibo.it)  
Aksieniia Shtimmerman, CNAF INFN [ashtimmerman@infn.it](mailto:ashtimmerman@infn.it)



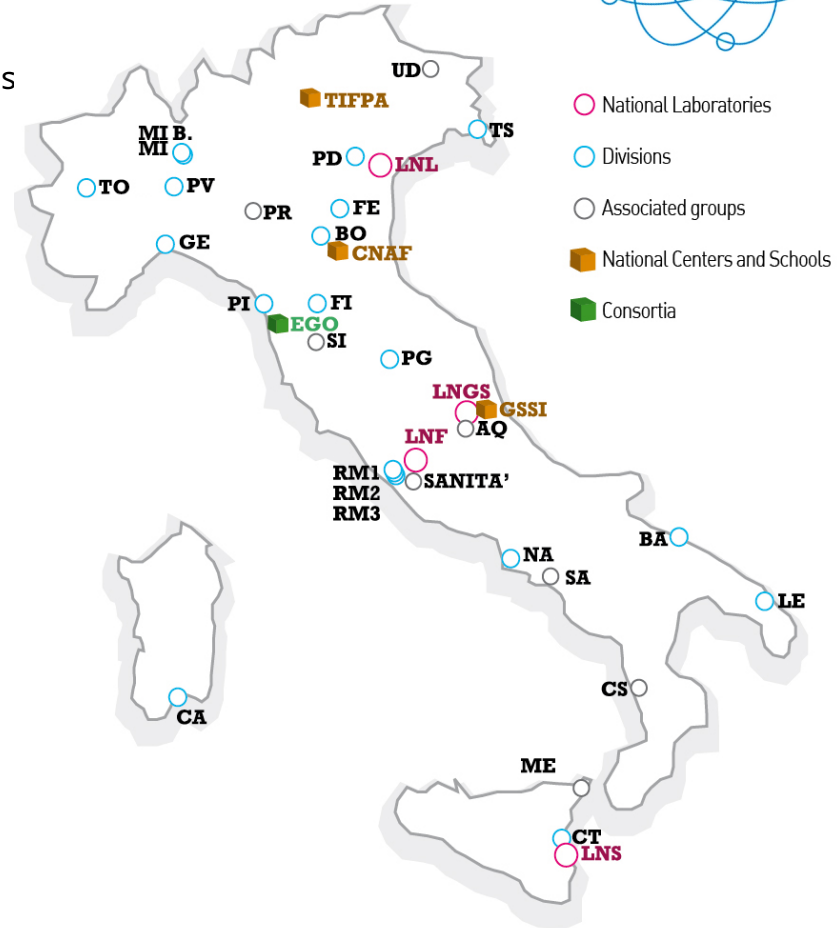


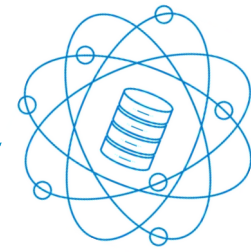
The CNAF Tier1 in the Italian supercomputing distributed infrastructure.

Italy provides Tier-1, Tier-2 and Tier-3 facilities to the ATLAS collaboration. The Tier-1, located at CNAF, Bologna, is the main center, also referred to as the regional center. The Tier-1, together with the other Italian centers, provides both resources and expertise to the ATLAS computing community and manages the so-called Italian Cloud of Computing.

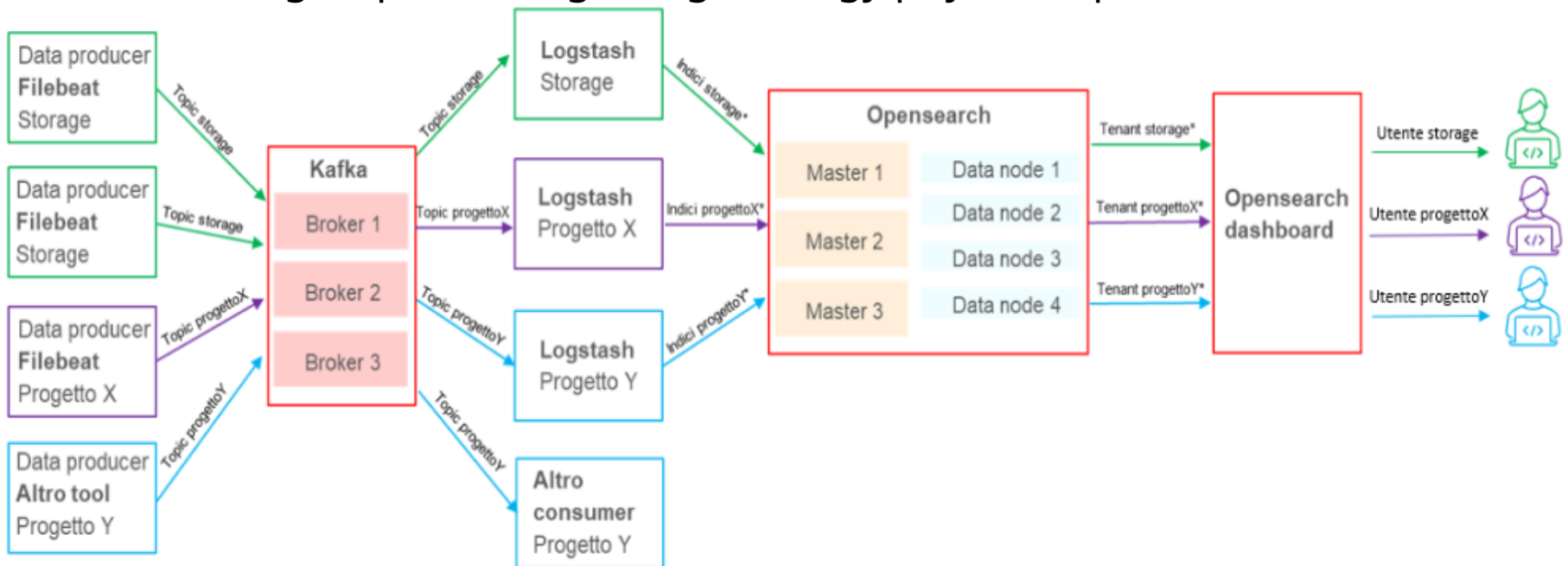
### 1 Tier1:

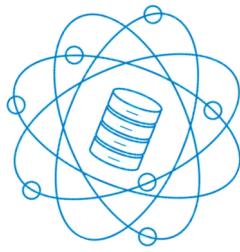
- 1.2MHSpec CPU power equivalent to 79'000 cores
- 80 PB disk space
- 65 PB tape library



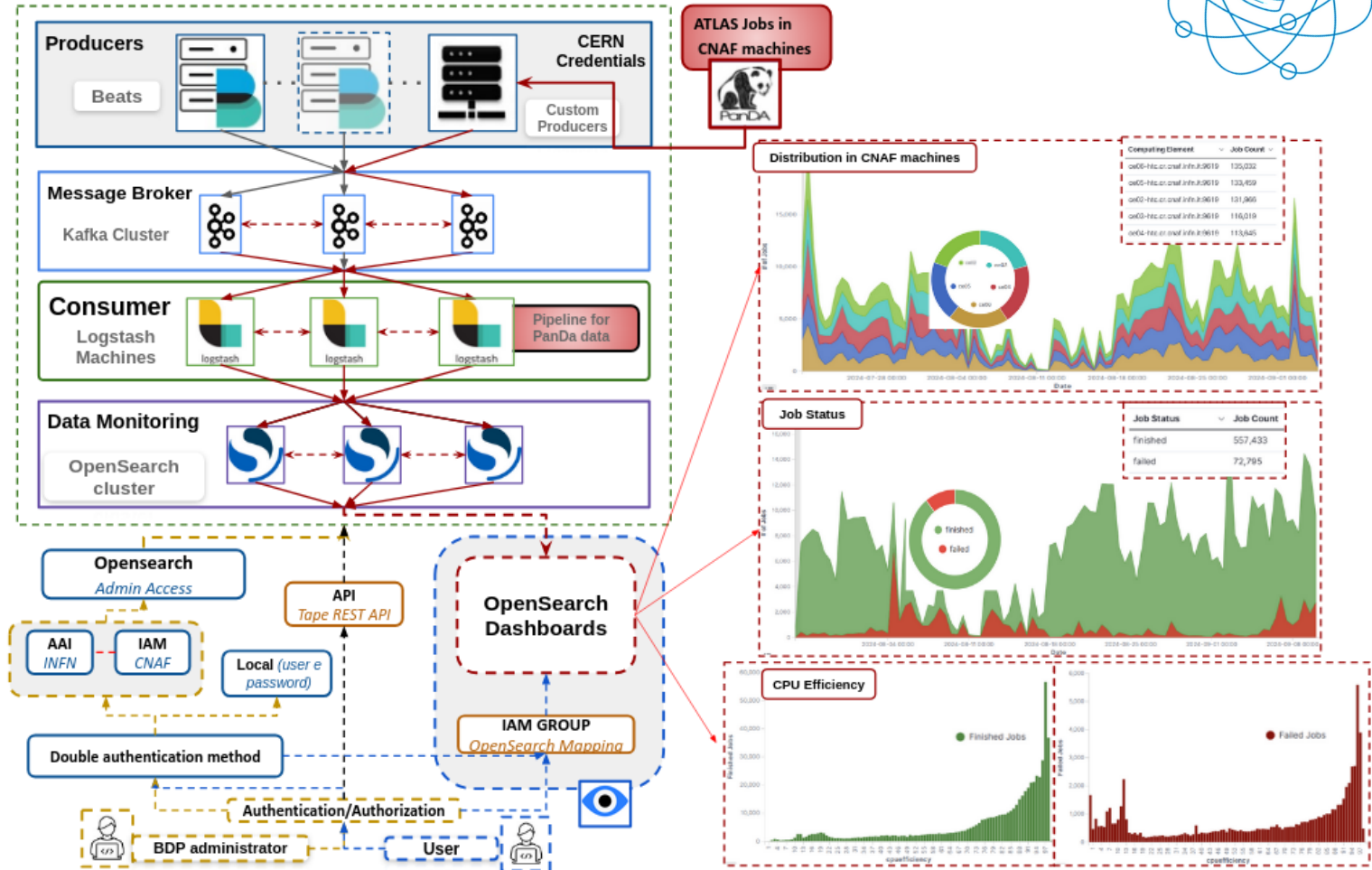


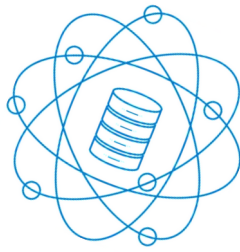
The CNAF group at INFN has implemented a Big Data Platform (BDP) infrastructure, designed for the collection and the indexing of log reports form CNAF facilities. The infrastructure is an ongoing project at CNAF and it is available for the italian groups working in high energy physics experiments.



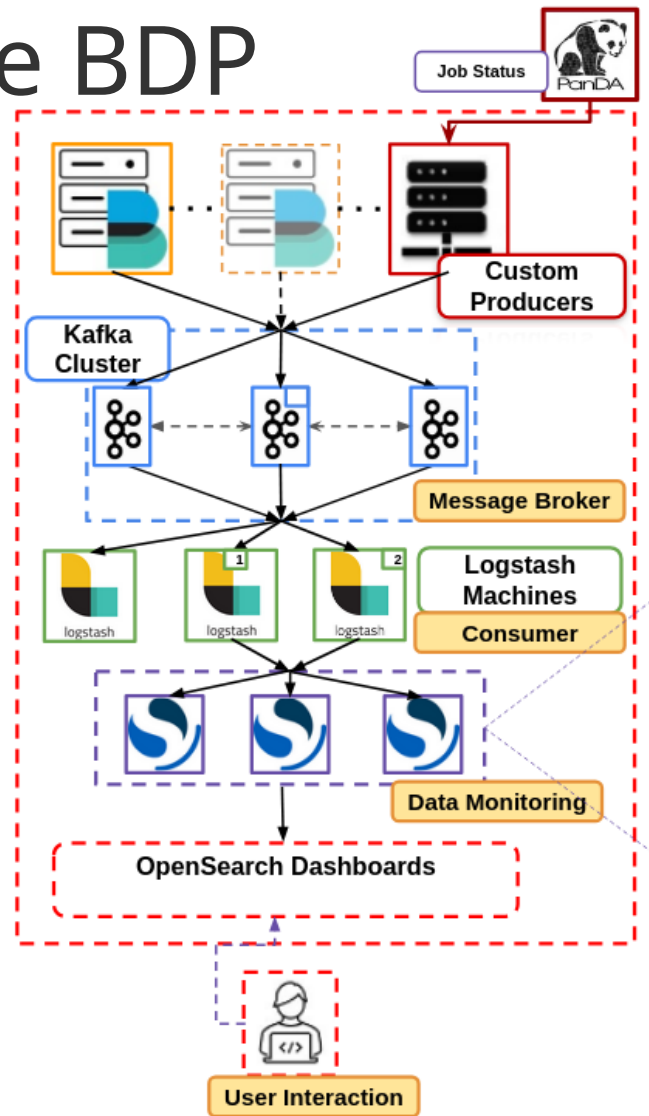


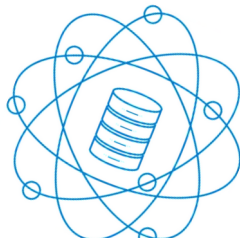
Within this framework, the first data pipeline was established for the ATLAS experiment at CERN, using input from the ATLAS Distributed Computing System PanDa.





ATLAS Job addressed to Tier-1 machines in Bologna are selected with a PanDA API, a summarizing log report is ingested in the data indexing pipeline in a compact JSON format with a custom producer.





This pipeline focuses on the ATLAS computational job data processed by the Italian INFN Tier-1 computing farm.

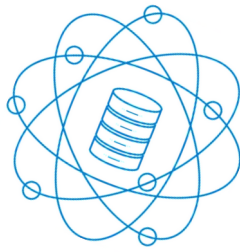
```
logstash$ cat /etc/logstash/config/pipelines.yml
```

```
pipeline.id: main
path.config: "/etc/logstash/config/*.conf"
pipeline.workers: 2
queue.drain: false
```

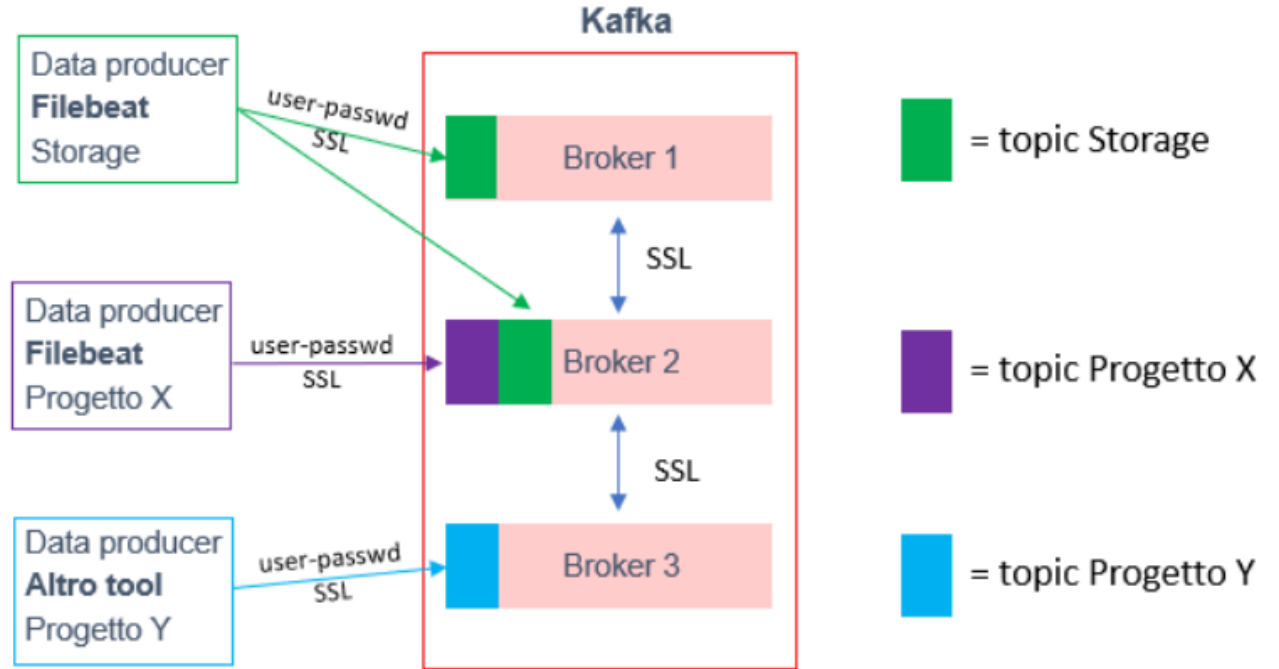
```
logstash0$ cat /etc/logstash/config/input.conf
```

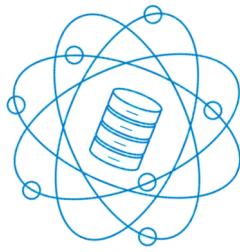
```
input {
  kafka {
    codec => json
    bootstrap_servers => "kafka01.cr.cnaf.infn.it:9192,kafka02.cr.cnaf.infn.it:9192,kafka03.cr.cnaf.infn.it:9192"
    sasl_jaas_config => "org.apache.kafka.common.security.plain.PlainLoginModule required username='*****' password='*****';"
    security_protocol => "SASL_SSL"
    sasl_mechanism => "PLAIN"
    ssl_truststore_location => "/root/truststore.jks"
    ssl_truststore_password => "*****"
    group_id => "atlas"
    auto_offset_reset => "earliest"
    topics => ["atlas"]
    id => "logstash"
  }
}
```

```
[aashtimmerman@cnlog-logstash02 ~]$ ll /etc/logstash/config
total 64
-rw-r--r-- 1 root root 2178 Jul 23 10:55 filters.conf
-rw-r--r-- 1 root root 654 Jul 15 12:03 input.conf
-rw-r--r-- 1 12011476 users 1833 Jul 19 2023 jvm.options
-rw-r--r-- 1 12011476 users 7437 Jul 19 2023 log4j2.properties
-rw-r--r-- 1 root root 342 Jul 15 11:15 logstash-sample.conf_orig
-rw-r--r-- 1 root root 31 Jul 15 11:31 logstash.yml
-rw-r--r-- 1 12011476 users 15104 Jul 11 11:11 logstash.yml_orig
-rw-r--r-- 1 root root 448 Jul 15 12:58 output.conf
-rw-r--r-- 1 root root 106 Jul 15 11:32 pipelines.yml
-rw-r--r-- 1 12011476 users 5204 Jul 11 11:10 pipelines.yml_orig
-rw-r--r-- 1 12011476 users 1696 Jul 19 2023 startup.options
```

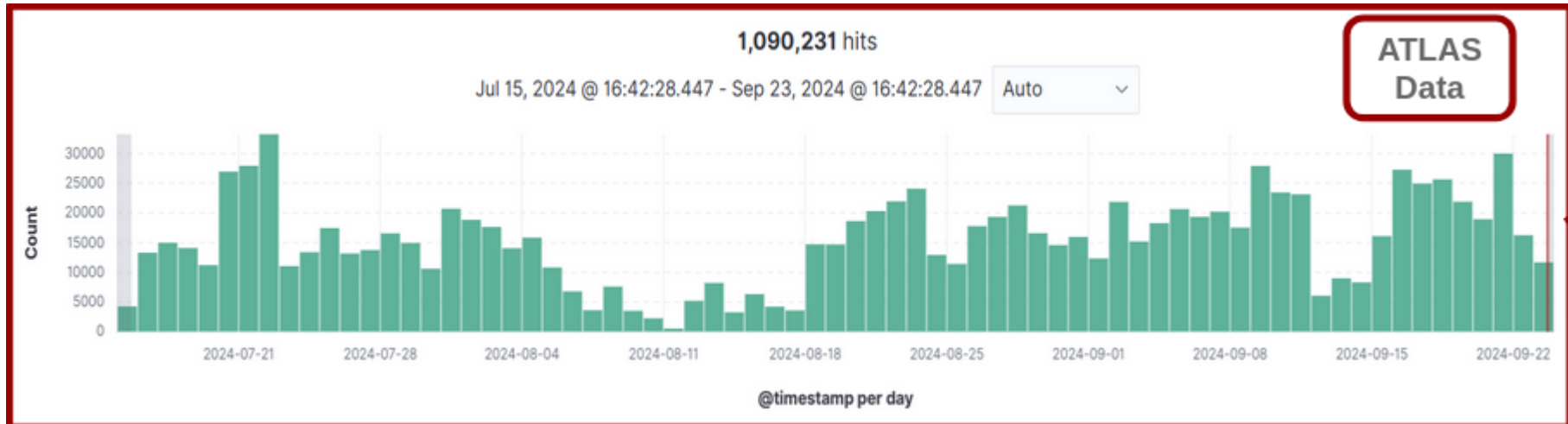


The system has been operational and effective for several years, marking our initiative as the first to integrate job information directly with the infrastructure.

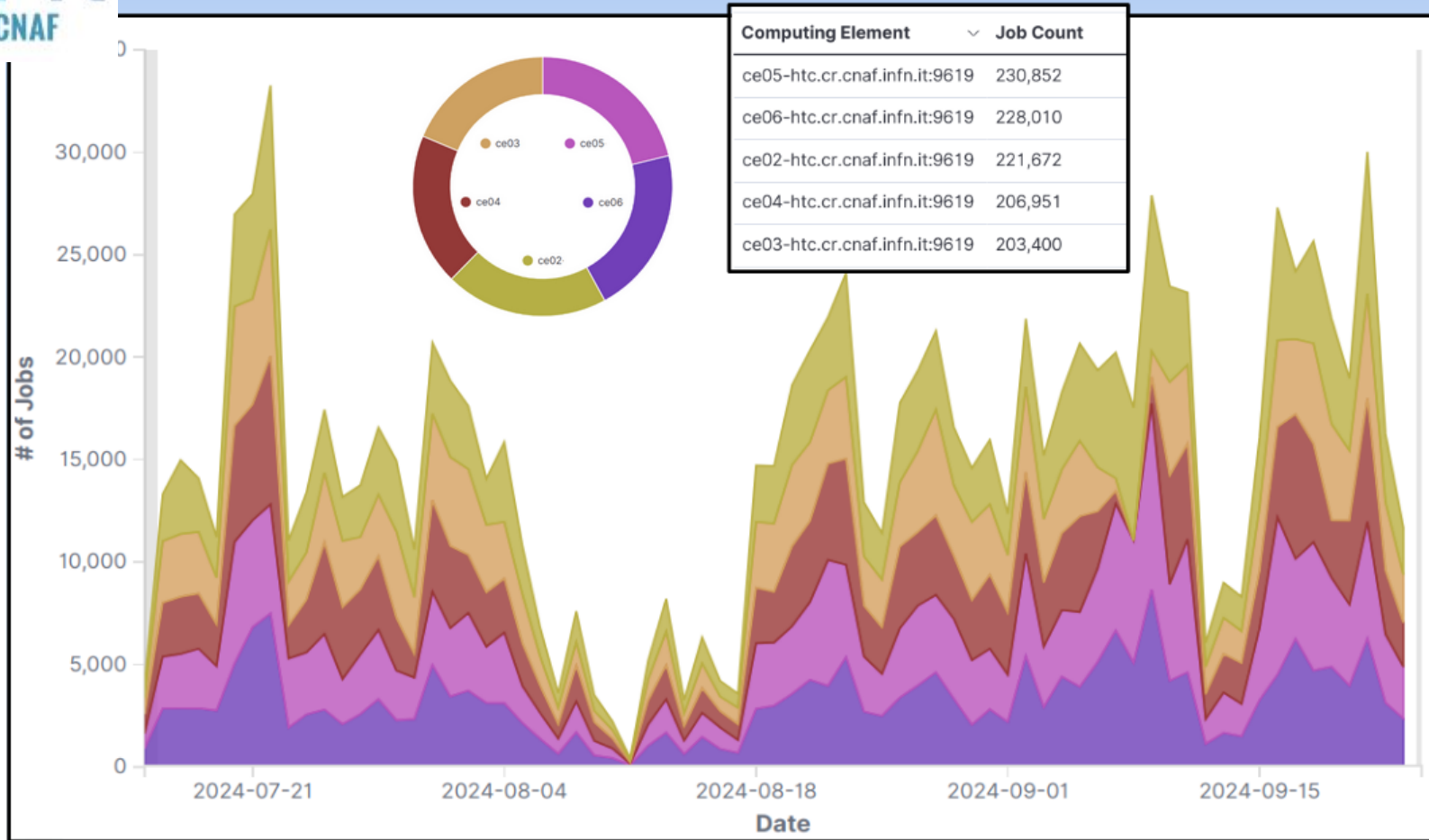
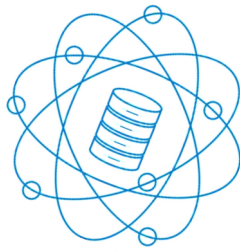


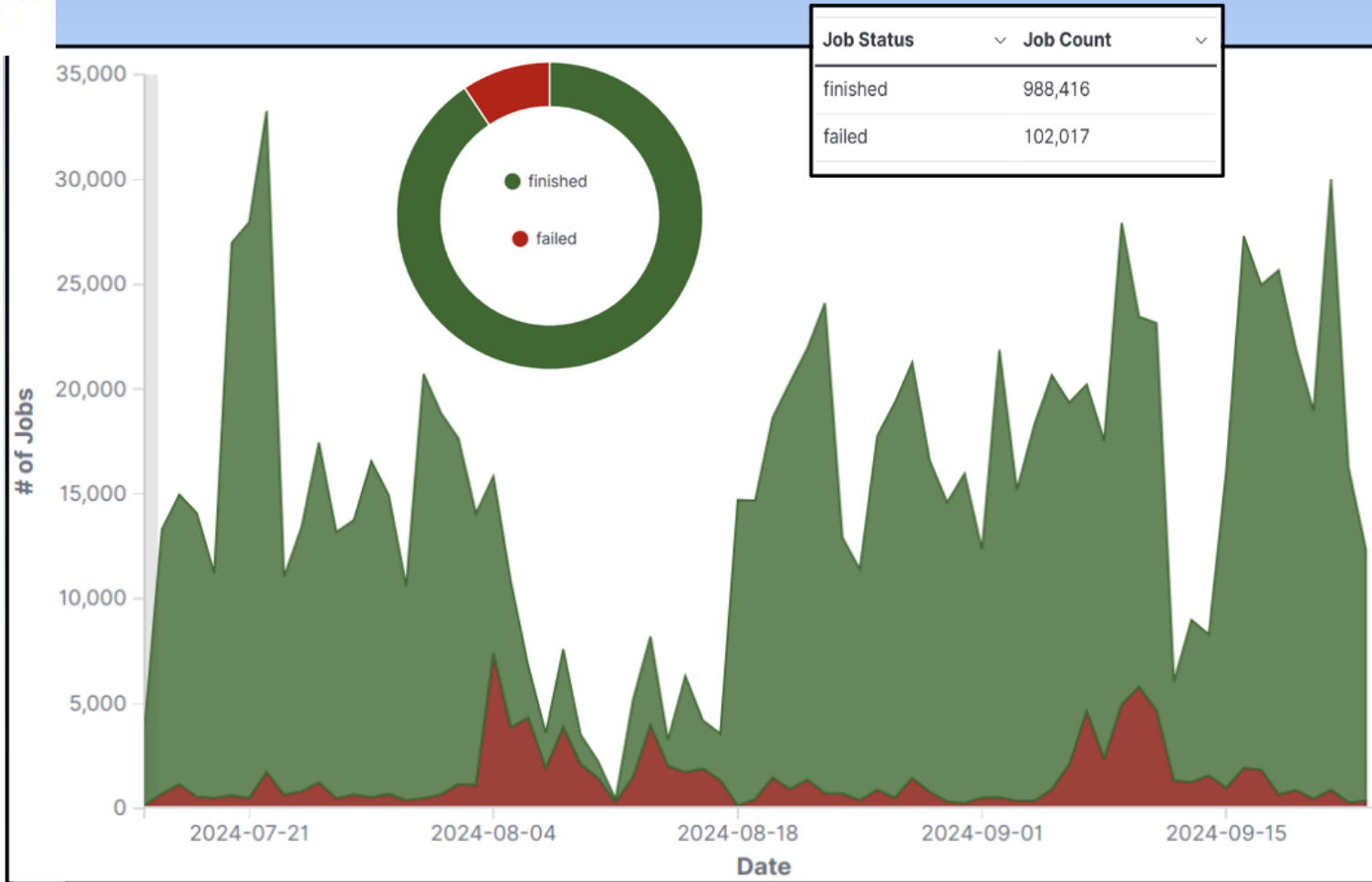
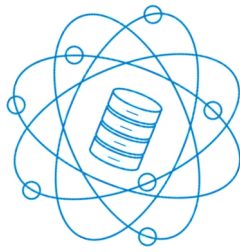


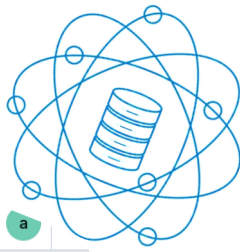
The save ATLAS computational job data and analysis of the jobs usage Big Data Platform of CNAF is the most important factor when trying to analyse the overall efficiency of the ATLAS computational job in the Italian big data center.



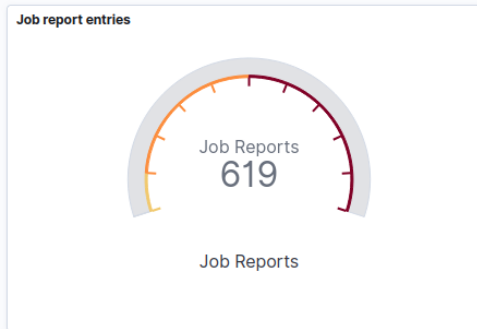






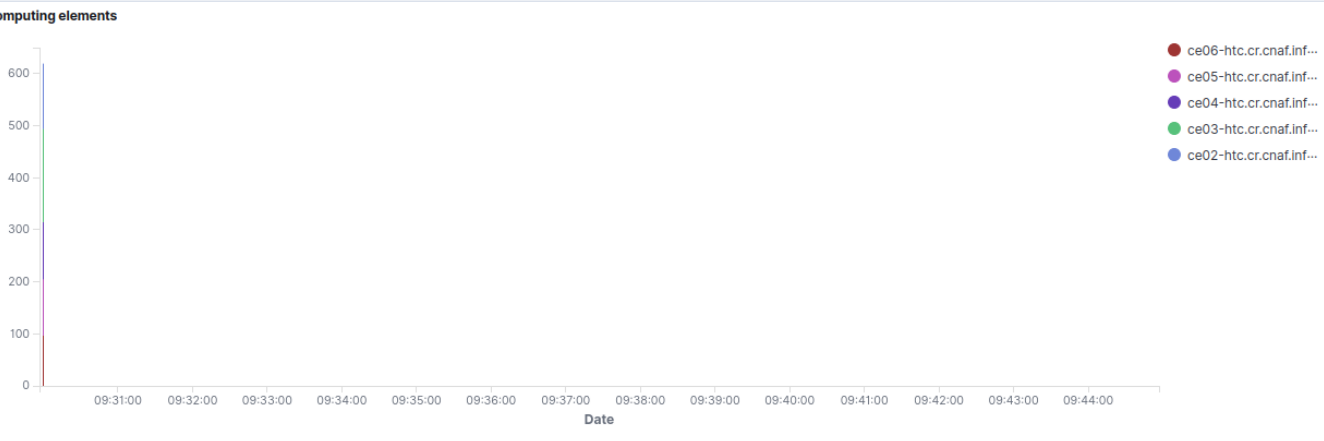
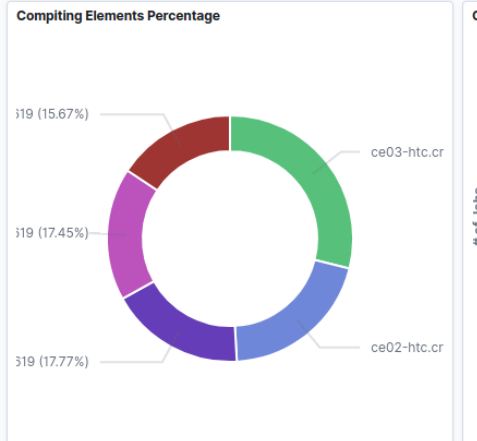


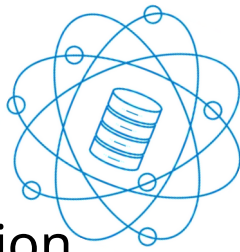
+ Add filter



### Computing Machines

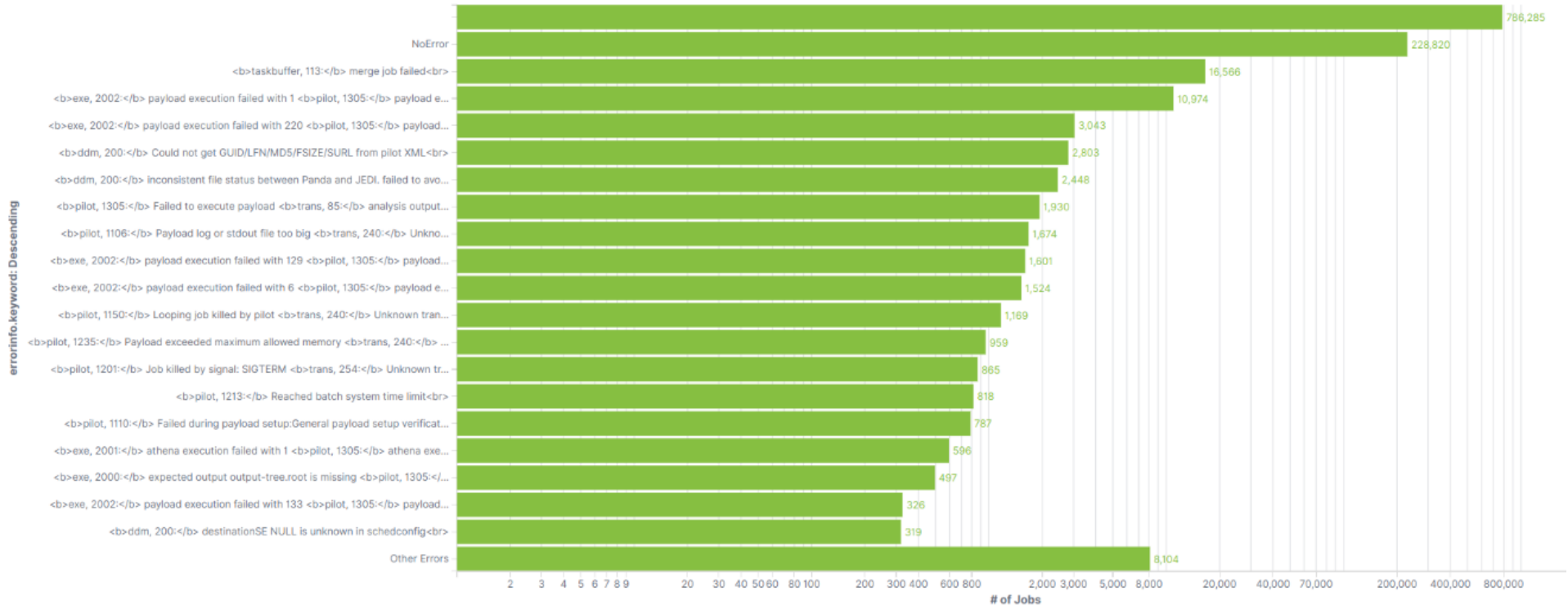
Computing Eleme...	Job Count
ce03-htc.cr.cnaf.infn.it	179
ce02-htc.cr.cnaf.infn.it	125
ce04-htc.cr.cnaf.infn.it	110
ce05-htc.cr.cnaf.infn.it	108
ce06-htc.cr.cnaf.infn.it	97

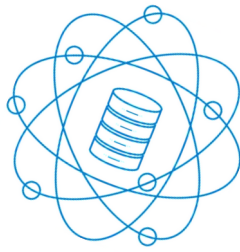




Different type of errors have been written inside the logs. Data collection will proceed along to have more entries with different errors.

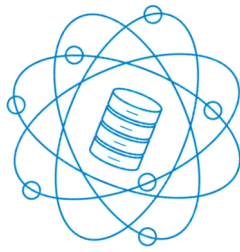
● # of Jobs





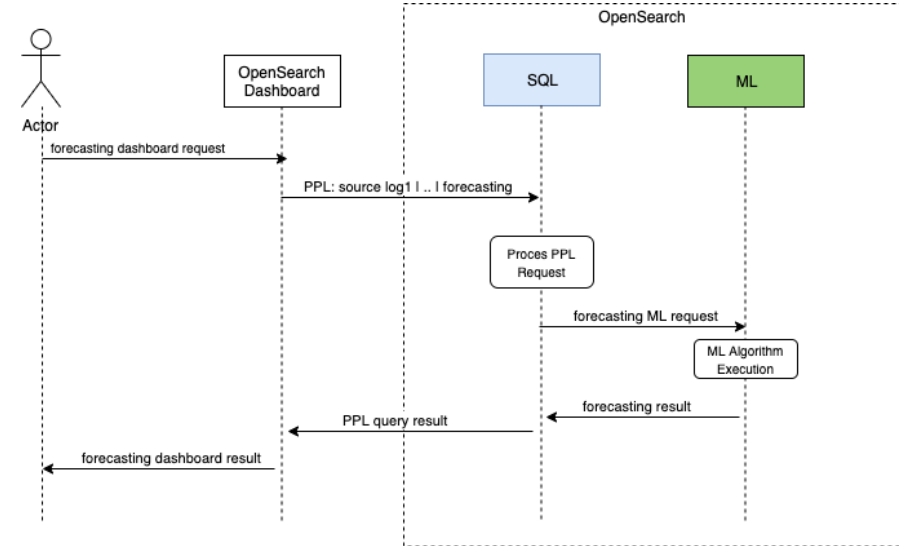
## Users and Admins access to Opensearch is granted via many authorization options:

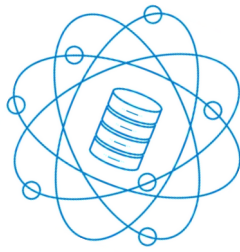
- Opensearch Dashboard or API;
- Double authentication method :
  - user and password;
  - local IAM CNAF and AAI INFN ([iam.cnaf.infn.it](http://iam.cnaf.infn.it));
- Authorization with IAM group (associated with permissions):
  - User IAM groups are automatically mapped to OpenSearch backend role;
  - Every User assigned to a given project (topic) may read or write inside the branch of the corresponding assigned topic, granting access to indices inside OpenSearch.



## Machine Learning

The main goal for the future evolution of the data ingestion inside the BDP is to start to use a ML plugins in OpenSearch, whose scope will be to investigate failed jobs, categorizing the reason of the failures and eventually seeking for anomalies inside the machines. We would like to implement and train ML algorithm, whose final scope will be to read and categorize the whole data throughput coming from PanDa.





## Contacts:

**Giacomo Levrini**, [glevrini@bo.infn.it](mailto:glevrini@bo.infn.it)

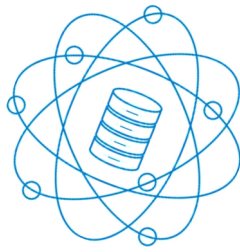
**Aksieniia Shtimmerman**, [ashtimmerman@cnafe.infn.it](mailto:ashtimmerman@cnafe.infn.it)

**Antonio Falabella**, [antonio.falabella@cnafe.infn.it](mailto:antonio.falabella@cnafe.infn.it)

**Enrico Fattibene**, [enrico.fattibene@cnafe.infn.it](mailto:enrico.fattibene@cnafe.infn.it)

**Diego Michelotto**, [diego.michelotto@cnafe.infn.it](mailto:diego.michelotto@cnafe.infn.it)

**Giusy Sergi**, [giusy.sergi@infn.cnafe.it](mailto:giusy.sergi@infn.cnafe.it)



# Questions?