# The ESCAPE Data Lake: The machinery behind testing, monitoring and supporting a unified federated storage infrastructure of the exabyte-scale

Rizart Dona on behalf of the ESCAPE project

CERN

May 20, 2021 - 25th International Conference on Computing in High-Energy and Nuclear Physics
(vCHEP 2021)

# Science Projects



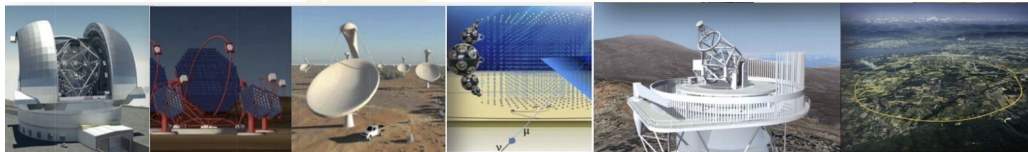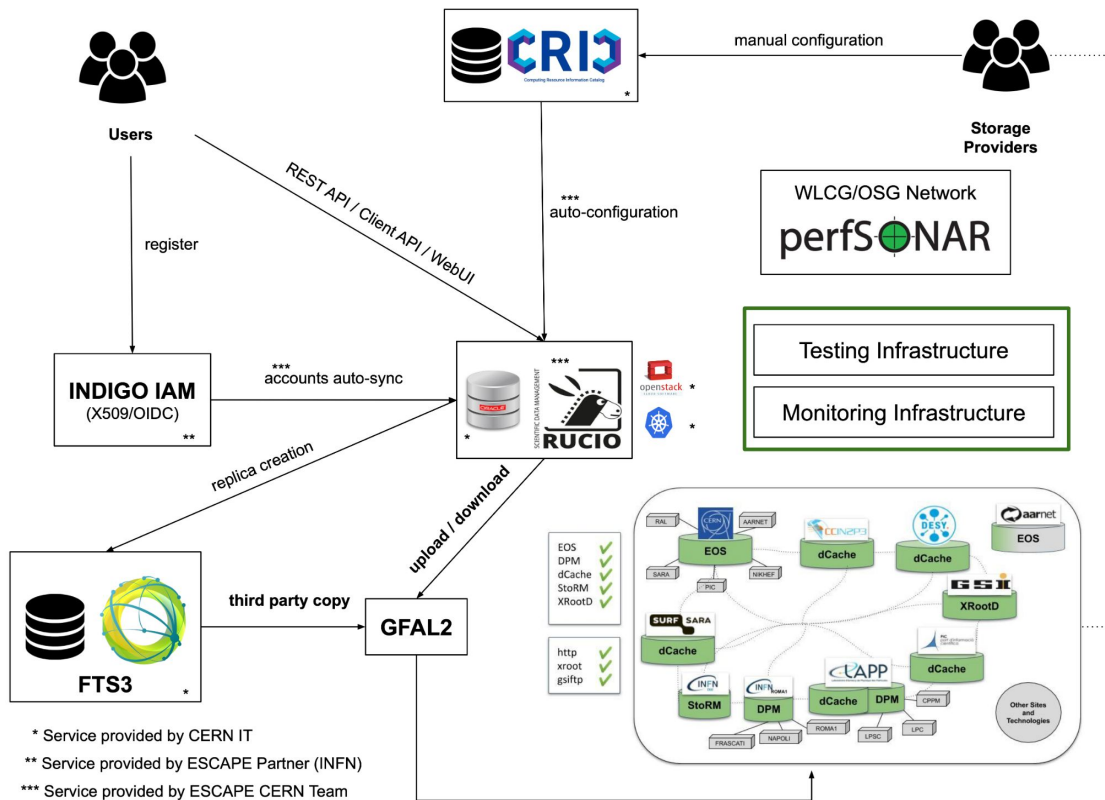# Partners



# Project Goals

- Prototype an **infrastructure** adapted to exabyte-scale needs of large science projects.
- Ensure sciences **drive** the development of EOSC.
- Address **FAIR** data management principles.

# Data Lake Architecture



- Federated data infrastructure
  - Multiple **storage/protocol** technologies
  - Combination of HEP-specific services and industry standards
  - **QoS** and file transitions, **distributed redundancy** and **data policies**
- Data Transfer Stack
  - **Rucio** → Data orchestrator, policy-driven data management
  - **FTS3** → Middleware, reliable large scale file transfer service
  - **GFAL2** → Grid file access library, multi-protocol access
- Identity and Access Management
- CRIC Information Catalogue
- perfSONAR boxes deployed in the OSG/WLCG network

---

- **Testing** infrastructure
- **Monitoring** infrastructure

\* Service provided by CERN IT
\*\* Service provided by ESCAPE Partner (INFN)
\*\*\* Service provided by ESCAPE CERN Team

rizart.dona@cern.ch

# Testing Infrastructure - Continuous Testing

- In order to make sure that all data transfer/access solutions are functioning

  properly we have **continuous testing** in place

  → Crucial step for consolidating the infrastructure

- Separate tests target each component individually and explore scenarios

  that involve both **functional testing** as well as **stress testing**

  → Breaking the testing complexity of the stack

- Configurable software has been developed and deployed to make the process **automatic**

- Data from testing is **visualized** in the equivalent Grafana dashboards

  → Part of the monitoring of the Data Lake

# Testing Infrastructure - GFAL

- All **Rucio Storage Elements (RSE)** consist of one or more endpoints that are associated with a supported protocol

- There are three types of operations that are being done concerning GFAL **functional testing**

  - **Upload** of files to all endpoints of all RSEs

  - **Download** of the files that were uploaded in the previous step

  - **Deletion** of the files that were uploaded in the first step

- This flow examines the **basic data operations** one can perform on the storage level of the Data Lake

  - Results → per storage element / per endpoint

  - Automatically pushing results in an Elasticsearch datasource

- Integrated with **CRIC**

  - Automatically fetching the configuration of the RSEs before each run

rizart.dona@cern.ch

# Testing Infrastructure - FTS

- In this case, the same endpoints as in the GFAL tests are examined
- Goal is to trigger **third party copy** transfers between all possible endpoint pairs that participate in the Data Lake

- Configurable software to trigger the transfers asynchronously and do automatic cleanup
- Extensive **error handling**
    - → testing flow will continue even if endpoints fail mid-test

Config example:



Run example:

rizart.dona@cern.ch

# Testing Infrastructure - Rucio

- An agent is used for fast ad-hoc functional testing

    ○ Uploads files to all storage elements

    ○ Performs replica creation (triggers third party copies) between all pairs with so called Rucio rules

- rucio-analysis: extensible python3 framework for yaml-based Rucio tests

    ○ Developed and deployed by the ESCAPE SKA Team

    ○ Same type of tests as the testing agent

rizart.dona@cern.ch

rizart.dona@cern.ch

# Monitoring Infrastructure - GFAL Operations

- Monitoring basic **GFAL operations** per storage element

- Users can **filter** plots by

  - Endpoints

  - Operations (upload, download, delete) - SUCCEEDED/FAILED/SKIPPED

  - Protocol (gsiftp, root, http)

- **Error messages** available for failed operations

rizart.dona@cern.ch

# Monitoring Infrastructure - FTS Transfers

- Monitoring **TPC transfers** between the endpoints

- Main highlights for users
  - Aggregated stats
    - Attempted transfers
    - Percentages
    - Job states
    - Transfer types
  - FTS transfers **efficiency matrix**
  - **Error messages** & links to log files

# Monitoring Infrastructure - Rucio Activities

- Monitoring **Rucio** specific activities

  → **replica** creation/deletion

- **Matrices** help users understand the connectivity between storage elements (similar to the FTS one)

  - Transfer successes / Throughput / Transferred volume

- Table with **error messages** and links to the actual FTS transfer logs

May 20, 2021

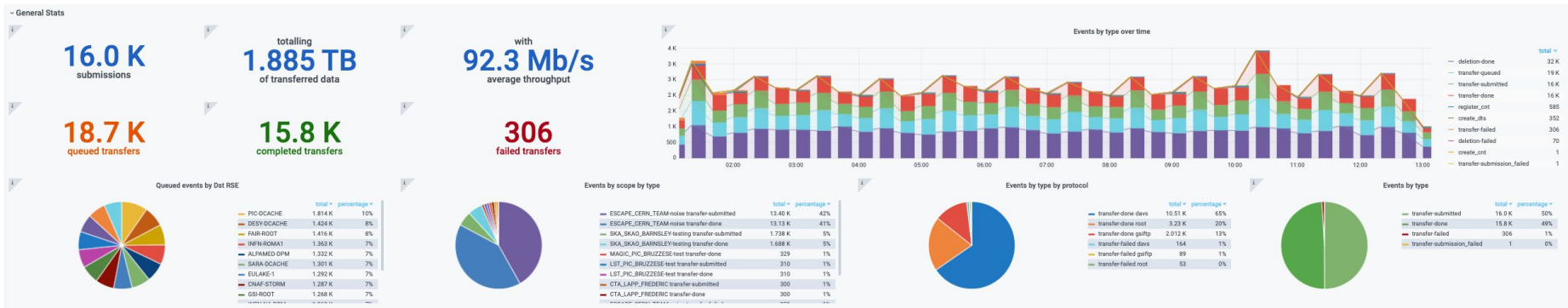rizart.dona@cern.ch

# Monitoring Infrastructure - Rucio Accounting

- Rucio **accounting data**
- Storage used / Number of files
  - Per **storage element**
  - Per **experiment**
- Captures trends over time
  - → Easy to observe changes on the storage level



Used Storage / Experiment (replica=1)

| | current |
|---|---|
| LOFAR | 32.0 TB |
| SKA | 9.22 TB |
| CTA | 1.78 TB |
| LSST | 824 GB |
| CMS | 391 GB |
| MAGIC | 331 GB |
| ATLAS | 163 GB |
| FAIR | 2.44 GB |
| VIRGO | 62.9 MB |

| Storage used for RSEs (all replicas) | Storage used for scopes (replica=1) | Storage used for experiments (replica=1) |
|---|---|---|
| **215.63 TB** | **156.62 TB** | **44.75 TB** |

Used Storage per RSE

| | |
|---|---|
| ALPAMED-DPM 2021-05-19T08:10:45.000Z | 15.79 TB |
| AWS_WEBDAV 2021-05-19T08:10:45.000Z | 0 B |
| CNAF-STORM 2021-05-19T08:10:45.000Z | 8.91 TB |
| CNAF_CMS_TEMP 2021-05-19T08:10:45.000Z | 1.53 GB |
| DESY-DCACHE 2021-05-19T08:10:45.000Z | 23.26 TB |
| EULAKE-1 2021-05-19T08:10:45.000Z | 25.96 TB |
| FAIR-ROOT 2021-05-19T08:10:45.000Z | 15.24 TB |
| GSI-ROOT 2021-05-19T08:10:45.000Z | 1.51 TB |
| IN2P3-CC-DCACHE 2021-05-19T08:10:45.000Z | 16.42 TB |
| INFN-NA-DPM 2021-05-19T08:10:45.000Z | 15.03 TB |
| INFN-NA-DPM-FED 2021-05-19T08:10:45.000Z | 13.57 TB |
| INFN-ROMA1 2021-05-19T08:10:45.000Z | 3.84 TB |
| LAPP-DCACHE 2021-05-19T08:10:45.000Z | 13.09 TB |
| LAPP-WEBDAV 2021-05-19T08:10:45.000Z | 144.87 GB |
| ORM-INJECT 2021-05-19T08:10:45.000Z | 126.15 GB |
| PIC-DCACHE 2021-05-19T08:10:45.000Z | 59.21 GB |
| PIC-INJECT 2021-05-19T08:10:45.000Z | 235.32 GB |
| SARA-DCACHE 2021-05-19T08:10:45.000Z | 62.44 TB |

rizart.dona@cern.ch

# Monitoring Infrastructure - The FDR experience

- Full Dress Rehearsal → 24h **full scale commissioning exercise**

  ○ Experiments were called to run their workflows on the Data Lake

Please find more details in Riccardo's [presentation](#)
(May 19th, Storage session)

- The monitoring infrastructure proved to be a **critical component** during this exercise

  ○ Instrumental for **identifying issues** and speeding up **debugging**

  ○ Users could follow up their data injection activity and **analyse** on the fly how this reflects on the Data Lake

rizart.dona@cern.ch

# Conclusions & Future work

- Successfully reached the **Pilot Phase**

  - Robust architecture that **serves** the needs of the sciences

  - **Testing** infrastructure and **monitoring** capabilities

    - Essential for validating the Data Lake status

    - Allow users and **experiments** to efficiently track their activities

- Moving towards the **Prototype Phase**

  - Consolidation of topics

    - **Token** based authentication/authorization

    - **Quality of Service** mechanisms

    - **Network optimized** transfers

- The ESCAPE Data Lake is constantly **evolving**

    → In sync with the **WLCG** activities and the **DOMA** project data challenges

rizart.dona@cern.ch

# References

- ESCAPE Project, https://projectescape.eu/

- ESCAPE Data Lake Wiki, https://wiki.escape2020.de/index.php/WP2_-_DIOS

- FTS3, https://fts.web.cern.ch/fts/

- GFAL2, https://dmc-docs.web.cern.ch/dmc-docs/gfal2/gfal2.html

- Rucio, https://rucio.cern.ch/

- GFAL testing software, https://github.com/ESCAPE-WP2/Utilities-and-Operations-Scripts/tree/master/gfal-sam-testing

- FTS testing software, https://github.com/ESCAPE-WP2/fts-analysis-datalake

- Rucio testing software, https://github.com/ESCAPE-WP2/rucio-analysis

- Elasticsearch, https://www.elastic.co/elasticsearch

- Grafana, https://grafana.com

- Apache ActiveMQ, http://activemq.apache.org

- ESCAPE Grafana Org, https://monit-grafana.cern.ch/d/cHBQ2NjWz/escape-home?orgId=51

Thank you!

Questions/Comments?

rizart.dona@cern.ch