



ESCAPE DIOS Developments & Achievements

Rizart Dona and Riccardo Di Maria on behalf of the ESCAPE project
CERN

March 31st, 2022 - Webinar: ESCAPE DIOS | A federated data-lake to connect large data volumes from different communities - benefits & use cases



Science Projects

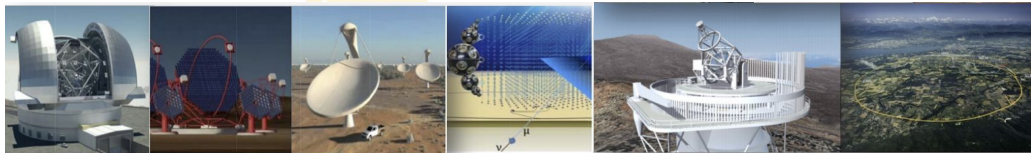


**EUROPEAN OPEN
SCIENCE CLOUD**

Horizon2020
European Union Funding
for Research & Innovation

Project Goals

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

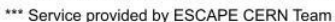


Partners



rijksuniversiteit
groningen



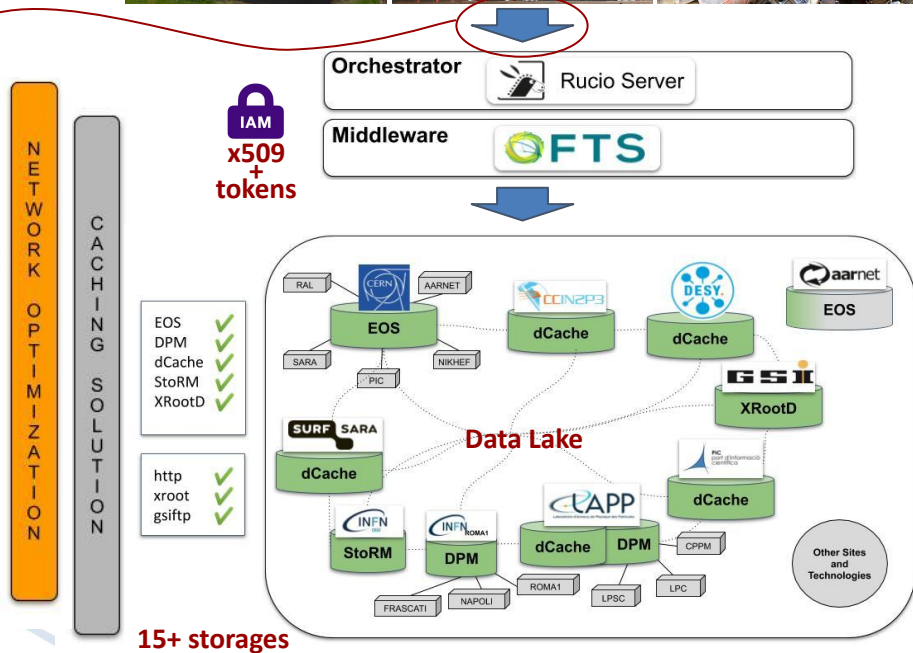


- Federated data infrastructure
 - Multiple **storage/protocol** technologies
 - Combination of HEP-specific services and industry standards
 - **QoS** and file transitions, **distributed redundancy** and **data policies**
- Data Transfer Stack
 - **Rucio** → Data orchestrator, policy-driven data management
 - **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

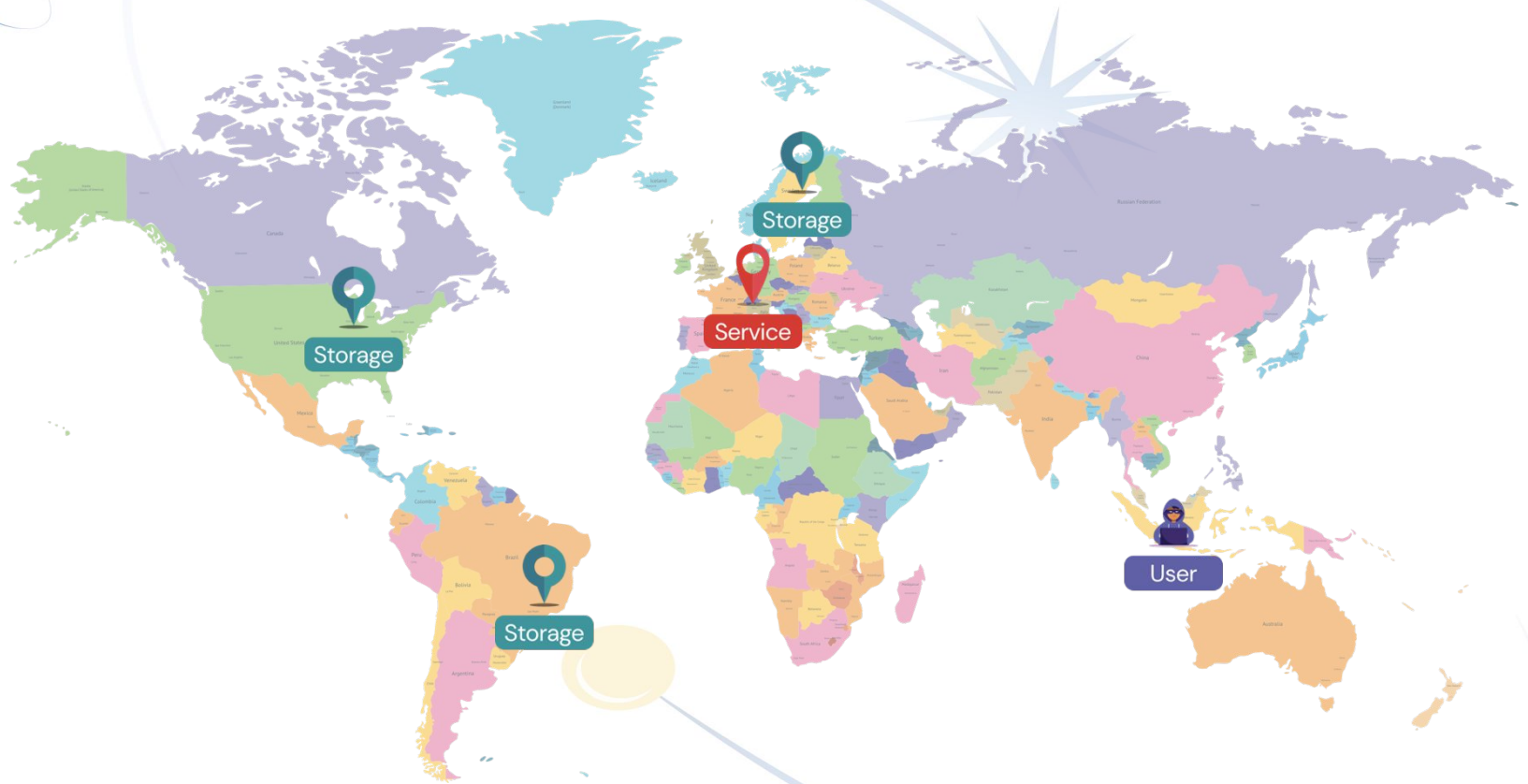
DLaaS for Open Science

- Goal: make **end-user comfortable** in embarking on a Data Lake experience
 - abstract the complexities of the Data Lake from the scientists
→ focus on doing science instead of data procurement
- An ever-increasing number of experiments are looking at Rucio Data Management system
 - **DLaaS** potentially interesting for both **aficionados** and **newcomers**

Sciences



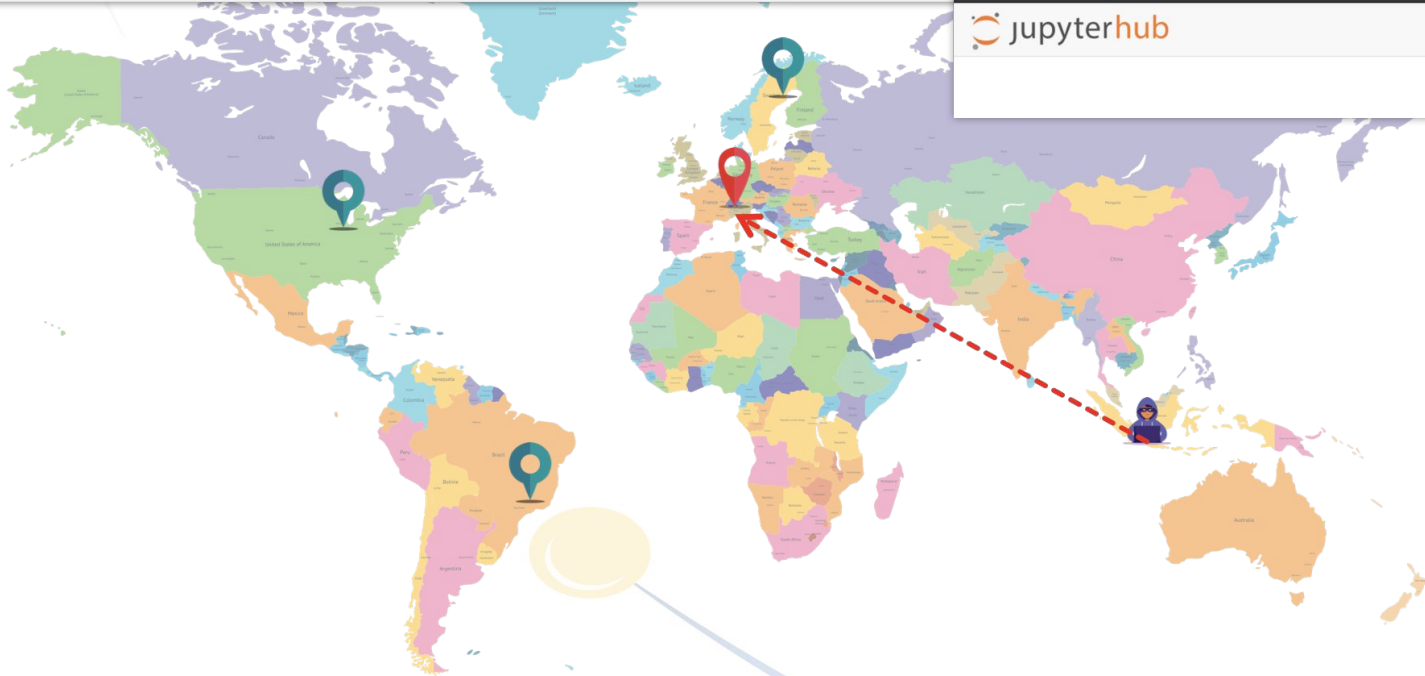
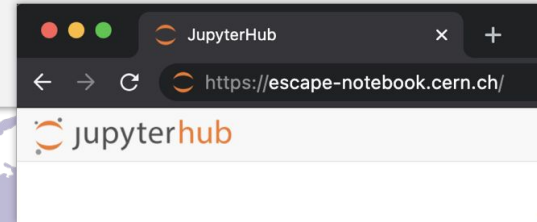
DLaaS - Overview



DLaaS - IAM (OpenID Connect) and Experiment Environments

→ Scientists Contact DataLake-as-a-Service

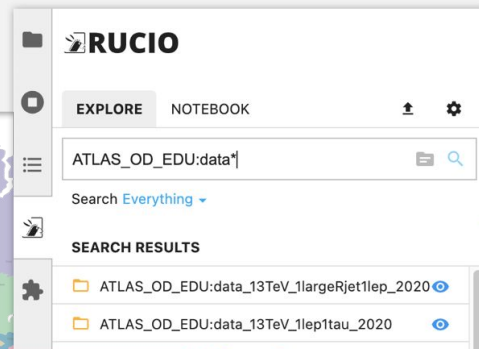
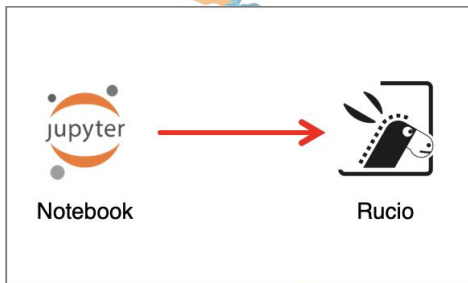
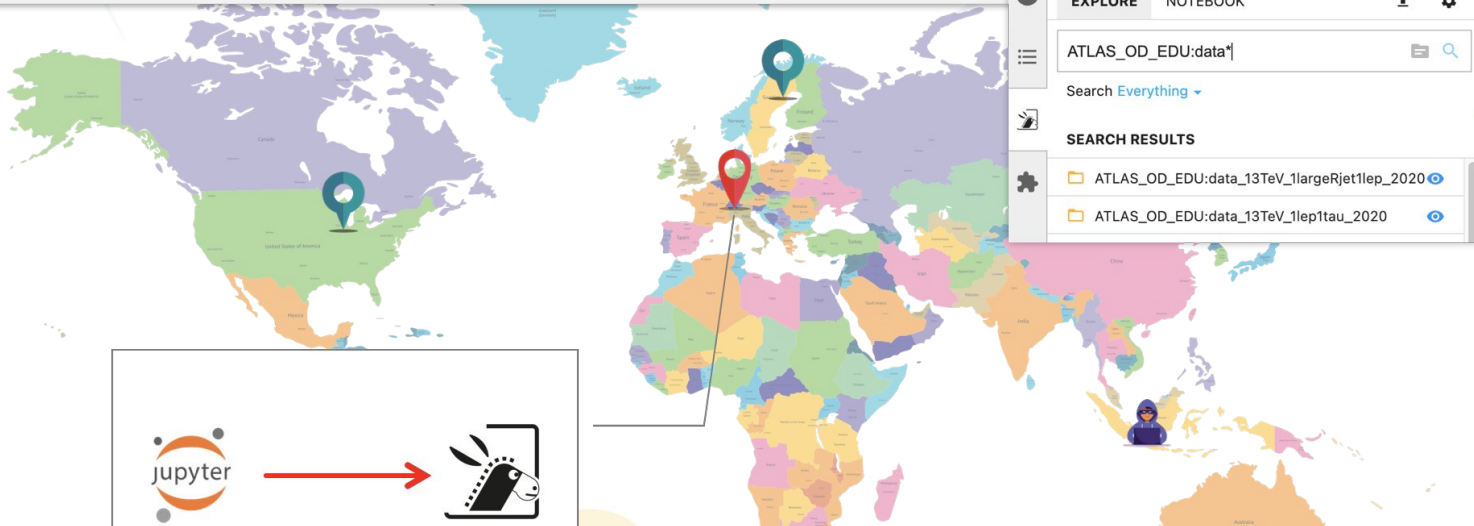
Requests are handled by Jupyter servers at CERN, Geneva



DLaaS - Data Browsing/Discovery & Access

→ Scientists Browse Data in the ESCAPE Data Lake

Requests are relayed to Rucio servers at CERN, Geneva



```
[3]: for item in hyy_20:
      print(item.pfn)
```

```
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/6f/98/data_A.GamGam.root
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/f1/3a/data_B.GamGam.root
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/45/95/data_C.GamGam.root
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/73/e3/data_D.GamGam.root
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/6d/aa/mc_341081.tth125_gamgam.GamGam.root.1
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/1b/95/mc_343981.ggh125_gamgam.GamGam.root.1
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/ff/c7/mc_345041.VBFH125_gamgam.GamGam.root.1
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/13/b8/mc_345318.WpH125J_Winc1_gamgam.GamGam.root.1
root://eoselake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/76/fd/mc_345319.ZH125J_Zinc1_gamgam.GamGam.root
```



DLaaS - Make Data Available in User Space, Data Preparation & Preservation

→ Make Available

Rucio initiates transfers from worldwide storages to CERN RSE which is serving DLaaS



RUCIO

Upload AUTHORS.md to Rucio

Please make sure that the necessary credentials are configured. You can see the upload status on the Rucio sidebar.

Destination RSE Expression:

EULAKE-1

Lifetime (in seconds):

Leave empty for indefinite

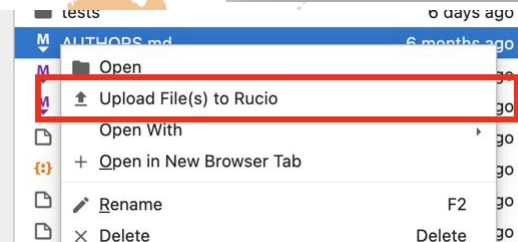
Scope:

atlas_od_

ATLAS_OD_EDU

Cancel

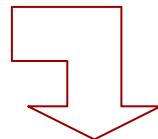
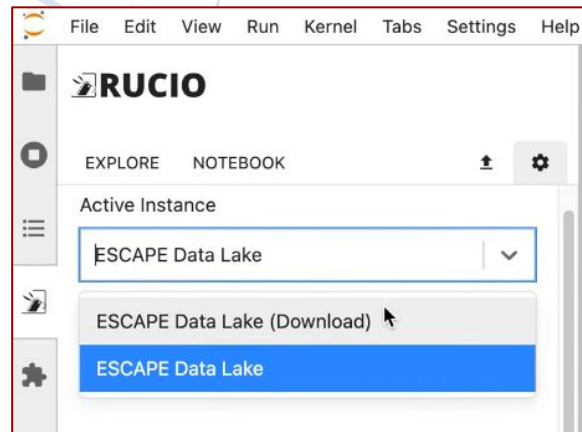
Upload



Server Options

- ☐ **Minimal environment**
Based on jupyter/scipy-notebook
- ☐ **ROOT environment**
If you need to use PyROOT
- ☒ **ROOT environment (Xcache testing)**
Run the extension in Download mode

Start



```

jovyan@jupyter-muhilmy:/scratch/muhilmy$ rucio list-file-replicas ATLAS_LAPP_JEZEQUEL:data.root --protocol root
+-----+-----+-----+-----+-----+
| SCOPE          | NAME          | FILESIZE | ADLER32 | RSE: REPLICAS |
+-----+-----+-----+-----+-----+
| ATLAS_LAPP_JEZEQUEL | data.root    | 4.660 kB | ef084c63 | EULAKE-1: root://xcache-redirector.cern.ch//root://eoseulake.cern.ch:1094//eos/eulake/tests/rucio_test/eulake_1/ATLAS_LAPP_JEZEQUEL/bd/8f/data.root |
| ATLAS_LAPP_JEZEQUEL | data.root    | 4.660 kB | ef084c63 | ALPAMED-DPM: root://xcache-redirector.cern.ch//root://lapp-testse01.in2p3.fr:1094//dpm/in2p3.fr/home/escape/rucio/lapp_dpm/ATLAS_LAPP_JEZEQUEL/bd/8f/data.root |
+-----+-----+-----+-----+-----+

```



Achievements

- ESCAPE managed to **pilot** and **prototype** a **Data Lake** infrastructure fulfilling functional data management needs of flagship ESFRIs from several scientific disciplines
 - successful assessments of the Data Lake in [2020](#) and [2021](#)
 - sensible technologies choice from WLCG environment and LHC experiments
 - community engagement & contribution
 - Astro-particle Physics, Electromagnetic and Gravitational-Wave Astronomy, Particle Physics, and Nuclear Physics pursuing together **FAIR** and **open-access** data principles
 - collaboration with other communities, e.g. PaNOSC, ExPaNDS, CS3MESH4EOSC, EOSC-Future
 - **DataLake-as-a-Service** hides the complexities of the Data Lake from the end users
- ESCAPE end in 2022 → addressing long term sustainability



Thank you!

Questions/Comments?

Backup Slides →



References

- ESCAPE Project, <https://projectescape.eu/>
- ESCAPE Data Lake Wiki, https://wiki.escape2020.de/index.php/WP2_-_DIOS
- [ESCAPE Data Lake: Next-generation management of cross-discipline Exabyte-scale scientific data](#), CHEP 2021
- [The ESCAPE Data Lake: The machinery behind testing, monitoring and supporting a unified federated storage infrastructure of the exabyte-scale](#), CHEP 2021
- [Data Lake as a Service for Open Science](#), ISGC 2022



DataLake-as-a-Service for Open Science

The screenshot shows the Rucio UI web interface. The browser address bar displays the URL: `escape-rucio-webui.com.cn/rucio/ui/Rule/7f6870c9f1da4d829d49cac3301ca644`. The page title is "Rule metadata". The interface shows a table of rule metadata for a specific rule ID.

Field	Value
account	ruhling
activity	User Subscriptions
copies	1
created_at	Wed, 15 Sep 2021 07:24:54 UTC
id	7f6870c9f1da4d829d49cac3301ca644
id_type	FILE
expires_at	Wed, 15 Sep 2021 08:24:54 UTC
grouping	DATASET
ignore_account_limit	False
ignore_availability	False
locked	Active
locks_ok_cnt	1
locks_replicating_cnt	0
locks_stuck_cnt	0
name	ggfH125_gamgert.root
notification	NO



DLaaS - As It All Started

- An idea presented at CS3 2020 by the Rucio team [[1](#)]
- Development of a “Rucio JupyterLab Extension” as part of GSoC 2020 [[2](#),[3](#)]
- A long time has passed, many things have happened...
 - CERN OpenLab Summer Student to concretise the effort in 2021
 - [deployment](#), [docker-images](#), [documentation](#)
 - DataLake-as-a-Service (DLaaS) in production-like phase
 - extensively exploited during ESCAPE “Data and Analysis Challenge” in November 2021 by SKA, MAGIC, CTA, ATLAS, KM3NET, LOFAR, FAIR
→ [3rd ESCAPE DIOS Workshop](#)
 - EU projects e.g. EOSC-Future and CS3Mesh4EOSC/ScienceMesh
 - other communities e.g. EGI



- Deployed in Kubernetes, using Zero-to-JupyterHub Helm chart → <https://escape-notebook.cern.ch>
- OAuth authentication using ESCAPE IAM (X509 still supported)
- Rucio JupyterLab Extension in Replica mode (i.e. TPC to local storage) used
 - download mode still possible (if configured)
 - connected to ESCAPE Data Lake
 - automatically pre-configured to use OpenID Connect
 - 2 FUSE mounts to EULAKE-1 (EOS)
 - ESCAPE RSE in r-only
 - additional RSE in r/w: SCRATCH
 - making files available
aka creating a replication rule to move files to EULAKE-1

