# ESCAPE: A Scientific-Data Lake for Open Science

By the ESCAPE WP2/DIOS

ESCAPE EEB Extended discussions day, 28 Sep 2021

# Outline

- A Data Lake. A data-what?! (Xavier)

- Data Lake as a Service Demo (Riccardo/Muhammad)

- WP2 in the past 12 month, and what's next (Rosie)

- Academic explanation of a Data Lake does not fit with the scope and purpose of ESCAPE's **Scientific**-**Data** Lake (SDL)
  - Is Data Lake a Bad naming? *probably!* Is it a fair analogy? *absolutely!*
- According to wikipedia:

## Data lake

From Wikipedia, the free encyclopedia

A **data lake** is a system or repository of data stored in its natural/raw format,[1] usually object blobs or files. A data lake is usually a single store of data including raw copies of source system data, sensor data, social data etc.,[2] and transformed data used for tasks such as reporting, visualization, advanced analytics and machine learning. A data lake can include structured data from relational databases (rows and columns), semi-structured data (CSV, logs, XML, JSON), unstructured data (emails, documents, PDFs) and binary data (images, audio, video).[3] A data lake can be established "on premises" (within an organization's data centers) or "in the cloud" (using cloud services from vendors such as Amazon, Microsoft, or Google).

- Academic explanation of a Data Lake does not fit with the scope and purpose of ESCAPE's **Scientific**-**Data** Lake (SDL)

  - Is Data Lake a Bad naming? *probably!* Is it a fair analogy? *absolutely!*

- According to ESCAPE:

  - **The Scientific-Data Lake is a policy-driven, reliable, and distributed data infrastructure capable of managing Exabyte-scale data sets, and able to deliver data on-demand at low latency to all types of processing facilities.**

A common set of tools is being adopted by different and multi-disciplinary scientific communities to stablish a coherent orchestration layer for the Research Infrastructures (RI) and for the sites providing resources to one or many of these RIs. The orchestration provided by these common set of tools include data management, data transfer and a common Identity Management, enabling a global vision of a single data pool easing the implementation of policies, rules and data life-cycles acting on the Scientific-Data Lake infrastructure as a whole, meaning all sites are working and seen as one.

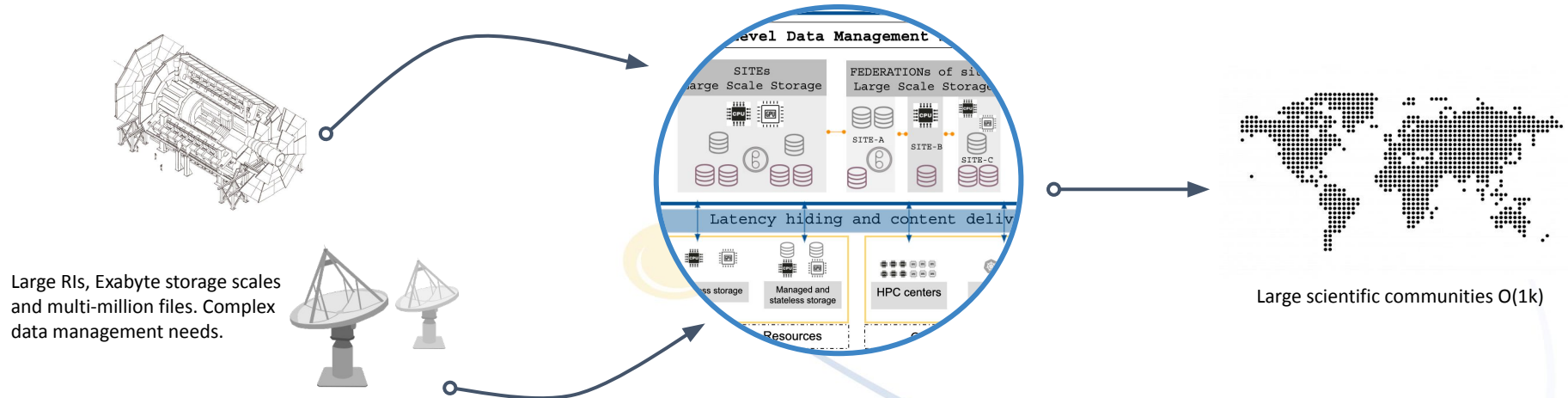# Is a Scientific-Data Lake complex to use?

**YES** and **NO**

The Data Lake provides **arbitrary** usage with an **arbitrary complexity** depending on the user requirements and needs.

A Scientific-Data Lake is **flexible** enough to address a whole range of different data management needs, let's see two different approaches.

# Is a Scientific-Data Lake complex to use?

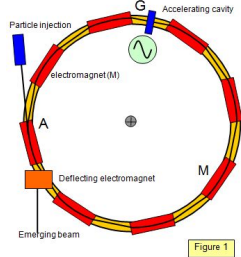Example-1: "Exabyte-scale" sciences, multi-million files, large and distributed user community

- Capability to provide fine-grained solutions on Data Management needs, harnessing different storages/sites. Capacity to implement data policies and life-cycles

- Address a data (re)processing campaigns with workloads attacking hundreds of Petabytes coherently across a large number of sites
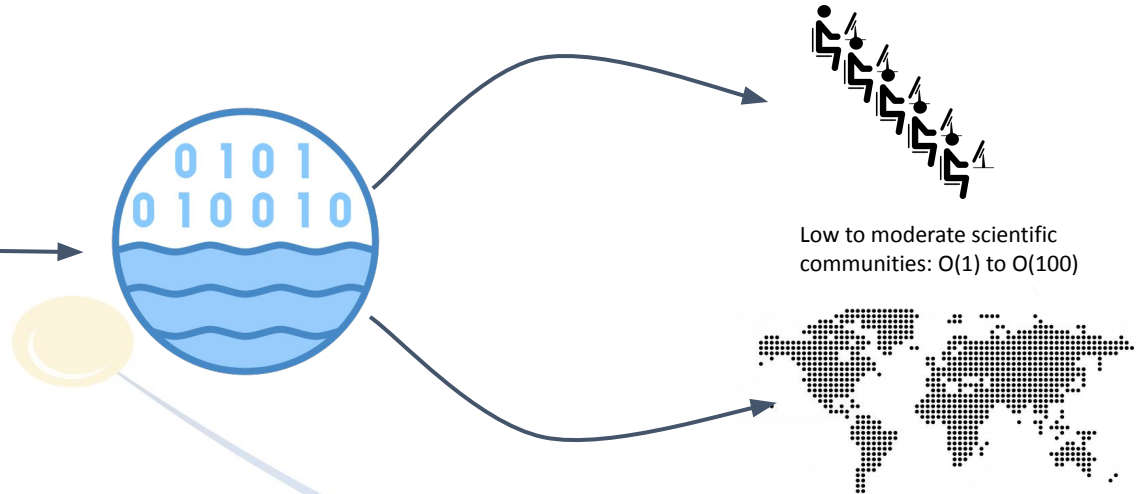


Large RIs, Exabyte storage scales and multi-million files. Complex data management needs.

Large scientific communities O(1k)

# Is a Scientific-Data Lake complex to use?

Example 2: Low to moderate data management needs, *Cloud*-like, simple push/pull model

- User/experiment see the SDL as an *http-endpoint*. Data is *FAIR-ed* behind the scenes and readily available for the user/collaborators on an URL-like style. The **complexity is hidden**.

- Analysis model based on notebooks or few tasks running in computing centres.



Instruments, experiments or measurements with low to moderate data volumes and number of files; store, share, analyze and archive. No specific data management needs.

Low to moderate scientific communities: O(1) to O(100)

# Is a Scientific-Data Lake complex to use?

In both of these previous examples the Data Lake ensures data is *FAIR-ed*: findability, accessibility, interoperability, and reusability, easing the process for the data to be open and scientifically reproducible (analysis preservation)

In summary, the Data Lake provides **arbitrary complexity** depending on the user's need, but providing the same capabilities and reliability in all scenarios

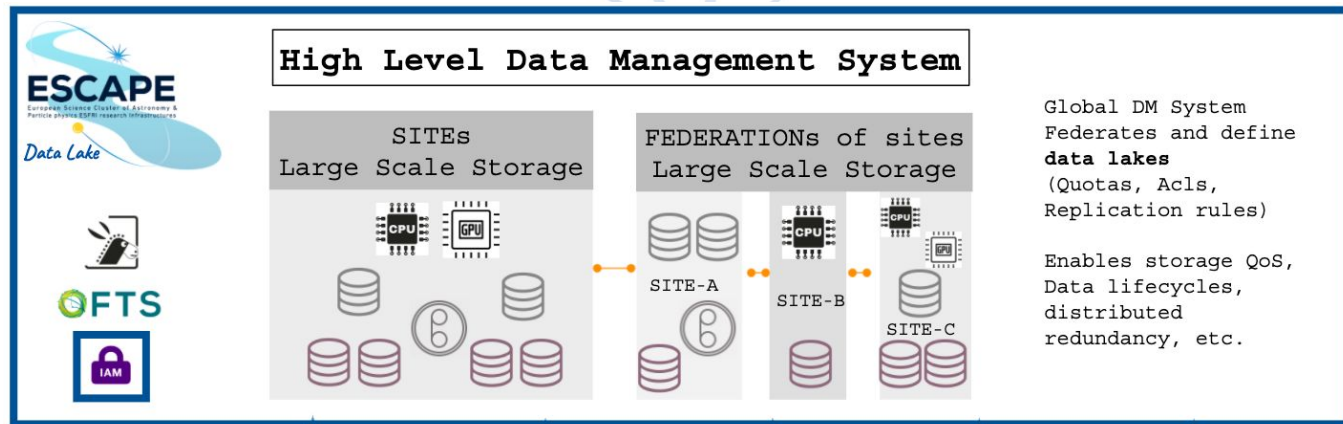# Right. But what is the reasoning behind a Data Lake?

Several upcoming challenges to address

- ## From the RIs/Experiments
  - **Exa-Scale or multi-PB is the norm** for current and future scientific machines, data volume, number of files and costs not anymore fitting on a bunch of diskservers
  - Need to **federate resources**: sites spread across a continent or several continents
  - Need to make data available to a **large user communities**
  - Need **data orchestration** tools: data "catalogue", a system to reliably and efficiently move files, information system, monitoring, data-lifecycles, set data policies, etc.
  - Integrate the growing offer of **heterogeneous computing resources**: private/public/hybrid clouds, HPCs, ephemeral resources, etc.

- ## From the sites and Infrastructure Providers
  - **Rationalise an increasing complexity**: sciences/experiments they provide resources are growing
  - Need for **homogeneous tools and services**: technologies, deployment and operations homogeneous as possible for the sciences they provide resources, leading to maintenance and operations **cost reduction**
  - **Resource optimization**. Is it worth a large scale disk storage with little I/O? or better focus on small disk stateless buffers with high I/O and invest more on C(G)PUs?

# Right. But how is the Data Lake addressing these challenges?

- The Data Lake **harnesses a set of sites**, providing each of them resources to experiments and user communities, with the goal to carry out independently well defined tasks

- The fundamental three main services providing the orchestration are:

    - A form of a **common trust relationship** (AAI), our ESCAPE IAM

    - A way to **locate data** is in place as a high level data management system (RUCIO), able to globally implement fine grained data management concepts: policies, replication rules, data life-cycles, on-demand redundancy, etc.

    - A **file transfer and data movement** mechanism (FTS) guarantee a high-level transport layer with the required protocols and interfaces with the storage systems at the sites

- The Data Lake model opens the door to storage consolidation by simplifying the scope of storage used for data processing oriented facilities
    - **Unmanaged/stateless storage**: staging areas, streaming-cache layer or buffers (data transport layers from new on) are used to transition files between the Scientific-Data Lake
    - **Ease integration** of heterogeneous computing resources: grid-like sites, HPCs, commercial, private and hybrid clouds, sporadic resources, etc.

# The ESCAPE Scientific-Data Lake



Users perceive one data infrastructure, federating resources and providing a common global namespace under a single entity allowing to define:

- Quotas
- Access rights
- Data Lifecycles
- Replication rules

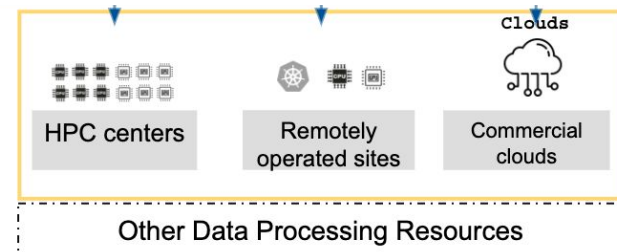This is achieved by a common **storage orchestration** and **data management system**
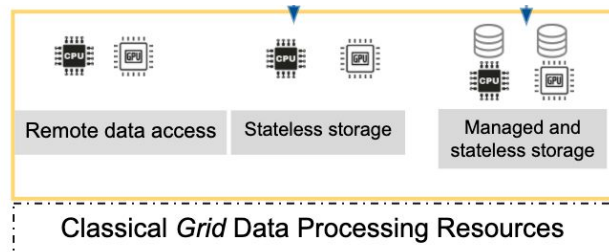
# The ESCAPE Scientific-Data Lake

Computing centers getting heterogeneous, standard *one-size-fits-all-grid-site* model is **vanishing**:

- Managed storage is not always needed nor effective
- CPU/GPU focus
- Exploit commercial/hybrid clouds resources
- Exploit HPCs
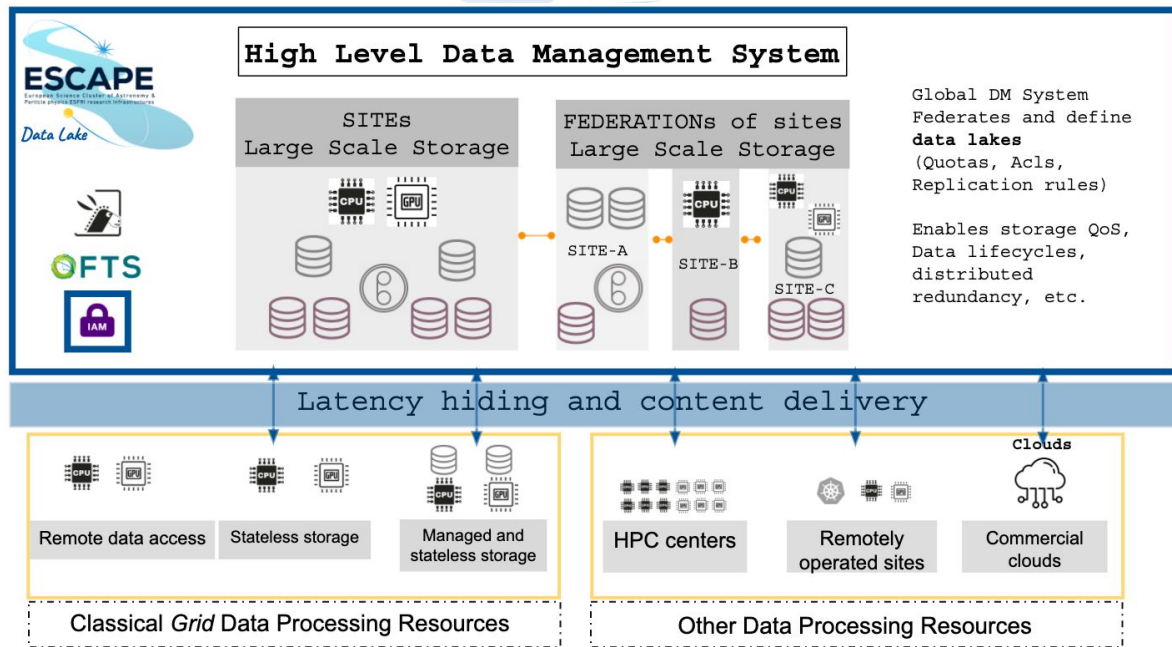- Opportunistic resources, e.g. time-limited contributions, vouchers

Compute is **stateless** but not Data; need to embrace the coming new political reality

| | | |
|---|---|---|
| Remote data access | Stateless storage | Managed and stateless storage |

Classical *Grid* Data Processing Resources

| | | |
|---|---|---|
| HPC centers | Remotely operated sites | Commercial clouds |

Other Data Processing Resources

# The ESCAPE Scientific-Data Lake

- Data management and orchestration: **Rucio**

- File transfer and data movement: **FTS**

- Content delivery and latency hiding: **XCache**

- Data Lake Information System: **CRIC**

- AAI: Indigo **IAM** (OAuth2/OIDC and legacy x509 support)

- The Data Lake **harnesses heterogeneous facilities** different types of storage systems: EOS, dCache, DPM, STORM, xrootd, http-enabled storage

- Allowing serving the data to heterogeneous facilities, from conventional Grid sites to HPC centres and Cloud providers

# A picture is worth a thousand words but a video is worth a million

Data Lake as a Service Demo (Muhammad/Riccardo)

# WP2 in the past 12 months and the future: Delivering results (Rosie)

- The past year's work
- Current focus and future plans

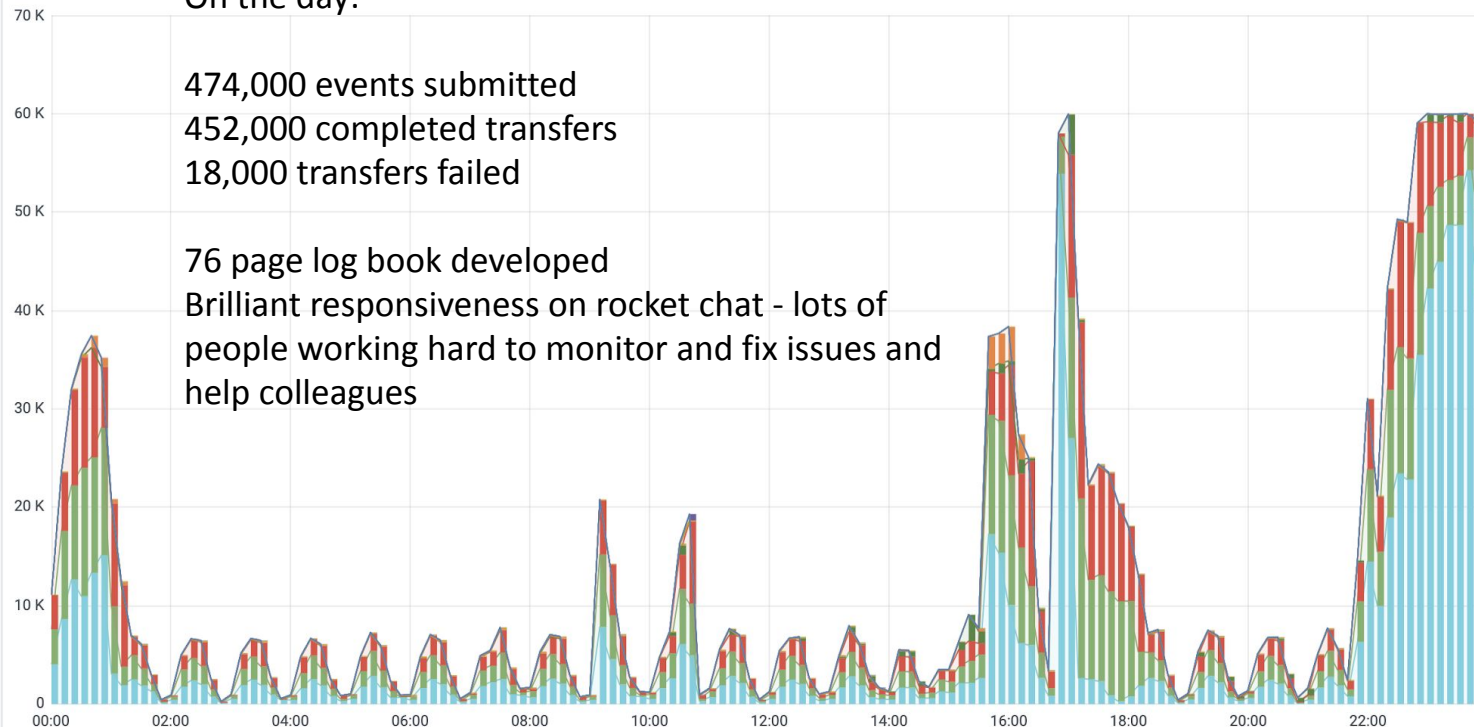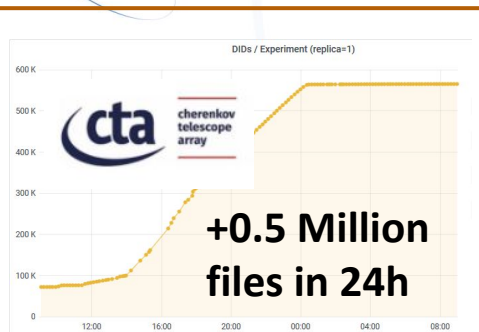# Data Lake **24-hour Dress Rehearsal** 17 Nov 2020



On the day:

474,000 events submitted
452,000 completed transfers
18,000 transfers failed

76 page log book developed
Brilliant responsiveness on rocket chat - lots of people working hard to monitor and fix issues and help colleagues
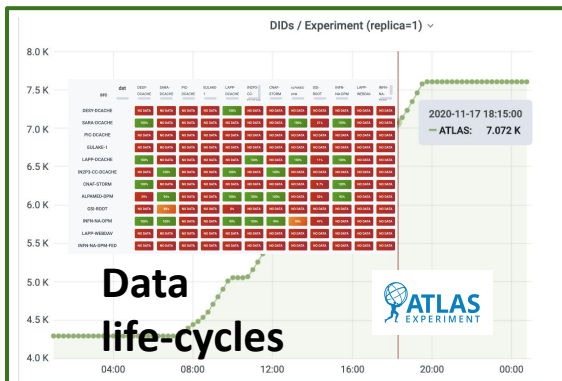
# Data Lake **24-hour Dress Rehearsal** 17 Nov 2020



**+0.5 Million files in 24h**
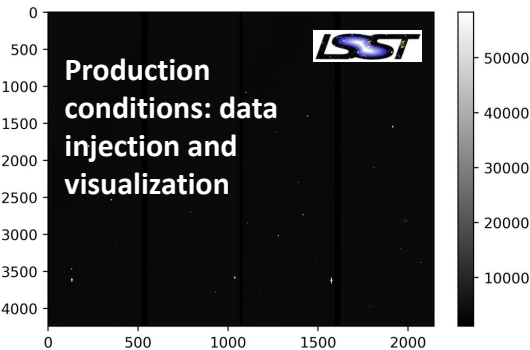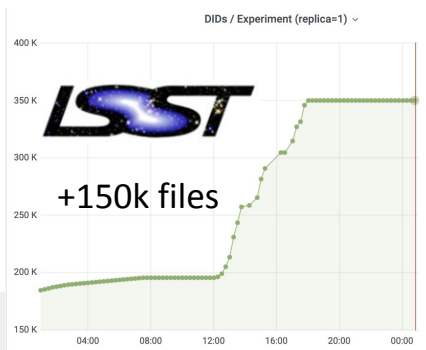
**CTA**: Simulate a night data captured from telescope in Canary Island for 6 h: ingest 500 Dataset of 10 files.



**Data life-cycles**

**ATLAS**: Storage QoS functionality tests: upload files from LAPP cluster to ALPAMED-DPM (FRANCE) and INFN-NA-DPM (ITALY), then request transfer to 1 RSE **QoS=SAFE** and 2 RSEs **QoS=CHEAP-ANALYSIS**



**From raw data to plots in the day**

**LOFAR**: astronomical radio source 3C196 made using LOFAR data. The raw visibility data was downloaded via rucio from the EULAKE-1 and processed on Open Nebula at surfsara using the container-based LOFAR software



**+150k files**

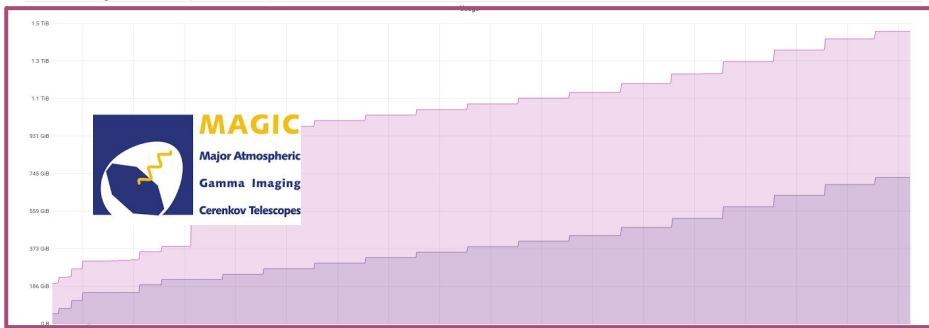**Production conditions: data injection and visualization**

**LSST**: Simulate production conditions: ingest the HSC RC2 dataset from CC-IN2P3 local storage to the Data Lake, **at a realistic LSST data rate** (20TB/24h). Then **confirm integrity and accessibility of the data via a notebook**.
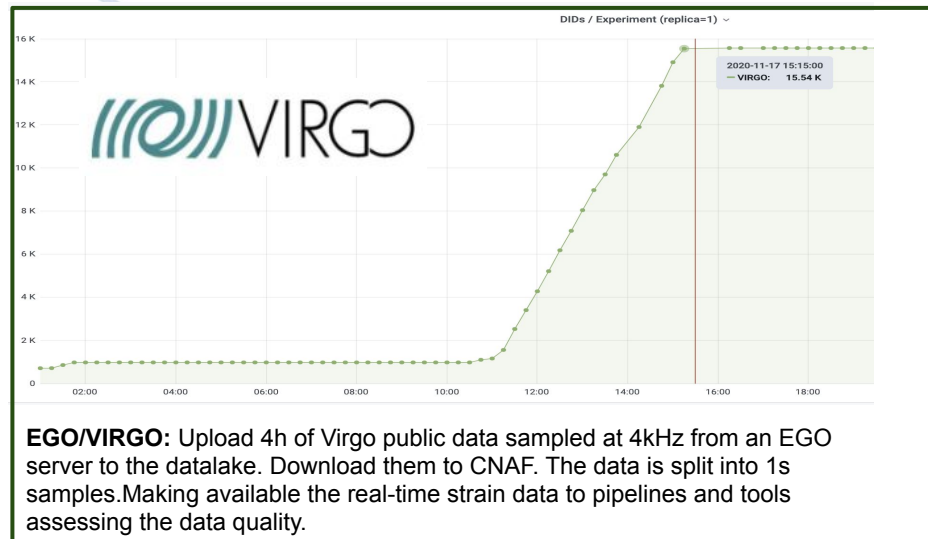
→ The image is a reconstruction drawn within a Jupyter Notebook accessing the data used in the Full Dress Rehearsal.
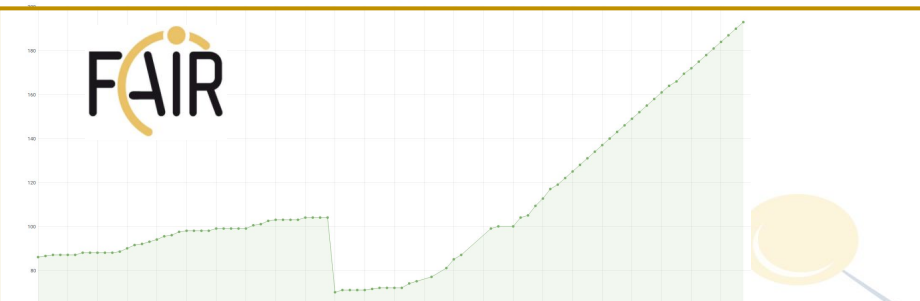
# Data Lake **24-hour Dress Rehearsal** 17 Nov 2020



**MAGIC:** Mimics a real MAGIC observation use case. Remote storage (Data Lake aware) **next to the telescope** acts as a buffer for subsequent data injection to the ESCAPE Data Lake (and local deletion after success)



**EGO/VIRGO:** Upload 4h of Virgo public data sampled at 4kHz from an EGO server to the datalake. Download them to CNAF. The data is split into 1s samples. Making available the real-time strain data to pipelines and tools assessing the data quality.



**FAIR**: Upload one 1-GB file every 10 minutes for the whole duration of the rehearsal. Request 2 replicas in QOS=SAFE and 1 replica in QOS=CHEAP-ANALYSIS. File size and QoS tagging approximate data ingestion from CBM (i.e. the FAIR experiment expected to produce the largest volume of raw data)



**SKA**: Pulsar Observations injection test. For 4 hours at any point during the 24hrs, injecting new group of files in a dataset every ten minutes. Files fall into two containers, representing different SKA Projects. 24-hr test moving data on basis of QoS class. + **2-million file dataset** test.

# Data Lake **24-hour Dress Rehearsal** 17 Nov 2020

On the day:

Mix of file sizes, dataset structures and rule types

Use of specific placement (at target RSEs) and automated placement based on QoS - which gave good resilience to downtime

Some transfers failed but rucio showed itself to be very resilient, able to recover

All FDR participants ran tests on the main ESCAPE Rucio instance

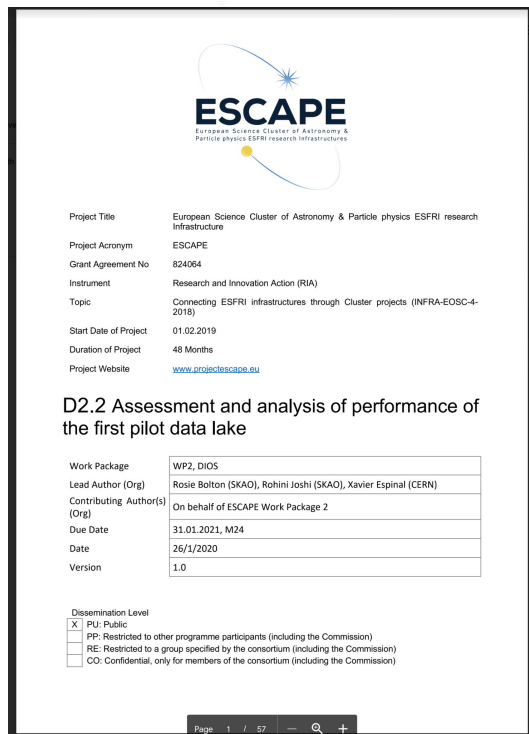Rucio Stats ☆ ⌔    2020-11-17 00:00:00 to 2020-11-17 23:59:59

asource | Monit ESCAPE (long term) ∨    rses All ∨    scopes All ∨    experiments All ∨    bin 1m ∨

### DIDs per Experiment (replica=1)

| Experiment | Number of DIDs | Number of files | Number of datasets | Number of containers | Average Filesize ↓ |
|---|---|---|---|---|---|
| LOFAR | 25.3 K | 25.2 K | 5 | 0 | 1.666 GB |
| FAIR | 194 | 192 | 2 | 0 | 1.036 GB |
| CMS | 401 | 398 | 3 | 0 | 1.026 GB |
| MAGIC | 13.5 K | 824 | 12.6 K | 18 | 573 MB |
| ATLAS | 7.604 K | 6.952 K | 652 | 0 | 235 MB |
| LSST | 350 K | 350 K | 13 | 0 | 18.5 MB |
| CTA | 564 K | 563 K | 1.458 K | 0 | 9.273 MB |
| SKA | 2.736 Mil | 2.703 Mil | 33.0 K | 25 | 3.259 MB |
| VIRGO | 15.6 K | 15.6 K | 10 | 0 | 86.4 kB |

# Nov 2020 FDR -> ESCAPE Deliverable 2.2

**FDR and pilot data lake assessment gave us a "wish list" for topics to address during the main data lake assessment phase**

- Better resilience of IAM
- File deletion understanding
- Testing and monitoring improvements
- Support for embargoed data / policy package
- QoS use improvements
  - including Tape use
- Rucio subscriptions
- Token-based auth
- Rucio technical developments
- Data corruption testing
- Complete onboarding of interested experiments
- Increase Rucio experience
- Rucio collaboration workshop
- Operations improvements
- Documentation "How to"

## ESCAPE

European Science Cluster of Astronomy &
Particle physics ESFRI research infrastructures

| Project Title | European Science Cluster of Astronomy & Particle physics ESFRI research Infrastructure |
|---|---|
| Project Acronym | ESCAPE |
| Grant Agreement No | 824064 |
| Instrument | Research and Innovation Action (RIA) |
| Topic | Connecting ESFRI infrastructures through Cluster projects (INFRA-EOSC-4-2018) |
| Start Date of Project | 01.02.2019 |
| Duration of Project | 48 Months |
| Project Website | www.projectescape.eu |

### D2.2 Assessment and analysis of performance of the first pilot data lake

| Work Package | WP2, DIOS |
|---|---|
| Lead Author (Org) | Rosie Bolton (SKAO), Rohini Joshi (SKAO), Xavier Espinal (CERN) |
| Contributing Author(s) (Org) | On behalf of ESCAPE Work Package 2 |
| Due Date | 31.01.2021, M24 |
| Date | 26/1/2020 |
| Version | 1.0 |

Dissemination Level

| | |
|---|---|
| X | PU: Public |
| | PP: Restricted to other programme participants (including the Commission) |
| | RE: Restricted to a group specified by the consortium (including the Commission) |
| | CO: Confidential, only for members of the consortium (including the Commission) |

Page 1 / 57

# "FDR" to "DAC21" (no longer a "Dress Rehearsal")

Format of targeted testing period was powerful in the FDR stage; we will repeat this as a focused week for assessment of the Data Lake prototype:

DAC21 = Data Analysis Challenge 2021

November 22-26 2021

**Experiment-led tests** determined by the use cases relevant for our partners

Focus on ensuring sustainability of the expertise and testing technical advances made during the past year

Three "Activity" areas, which cut across original task definitions

A1: Data Injection, Replication and Processing

A2: Multi-Rucio, Data Lake Operations;  A3: QoS, AAI developments

# A1- Highlights

- Full involvement of ESCAPE experiments for this activity(!)
  - Definition of tests to run / things to try for DAC21 assessment
- DataLake-as-a-Service (DLaaS)
  - Data preparation and access (+caching layer!)
  - Working on coupling with ESAP-WP5 as sharing common goal
- Additional tools interacting with Rucio - e.g. DIRAC workflow manager (CTA)
- Workflows will use underlying storage with **additional QoS**:
  - **Tape** (DESY, PIC, SurfSARA, CNAF) and **Erasure Coding** from CERN and SurfSARA.
- Both Deterministic and non-deterministic RSEs will be tested in DAC21
- Long haul transfers
  - On going work with MAGIC to send data from La Palma (Canary Islands) to PIC
  - SKA transfers between Australia, South Africa and Europe

# A2- highlight areas

Increasing the expertise

Improved monitoring

Technical developments

Preparing Rucio instances and underlying Data Lake infrastructure for DAC21

# Increasing expertise: Data Lake swimming club

Focus on knowledge building in ESFRI partners

As the main ESCAPE Rucio instance matures technically, we also support our partners in establishing their own Rucio instances

Offer deployment, testing and monitoring expertise into those new instances

*"Just how far can we take this lake analogy?"

Also continue to support all users on the main ESCAPE instance, testing advanced DL features.

# Data Lake swimming club

## MAGIC + CTA

Long-running MAGIC instance collaborating with CTA colleagues as use case very similar

Transfers from La Palma to Barcelona, mimic real-life data management pattern, take over as production system if successful

## SKAO

SKA Rucio instance set up as long-term test bed. Include RSEs in SKAO host countries.

Build deployment expertise (i.e. break things!)

## ASTRON / LOFAR

Planned rucio instance for data management involving MeerKAT data (SA) and ASTRON users, and data between Groningen and Dwingeloo for one of our science projects.
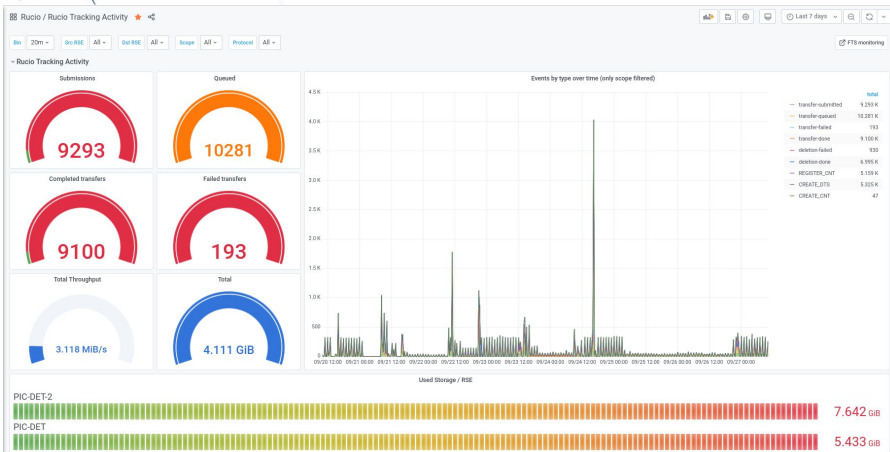
## LSST (IN2P3)

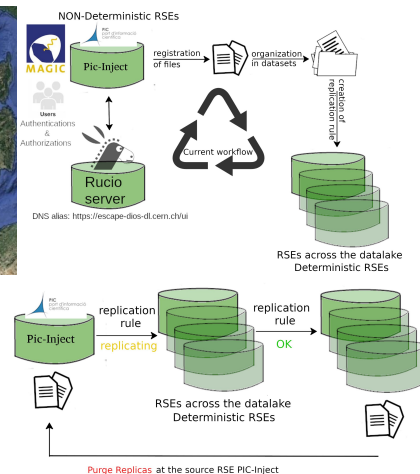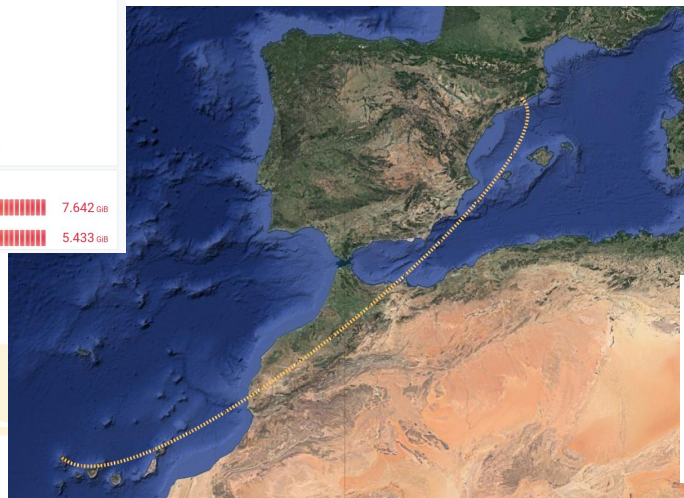LSST Rubin science platform SW assumes specific hierarchical filesystem.

Test instance set up to enable exploration of method to control namespace mapping in Rucio.
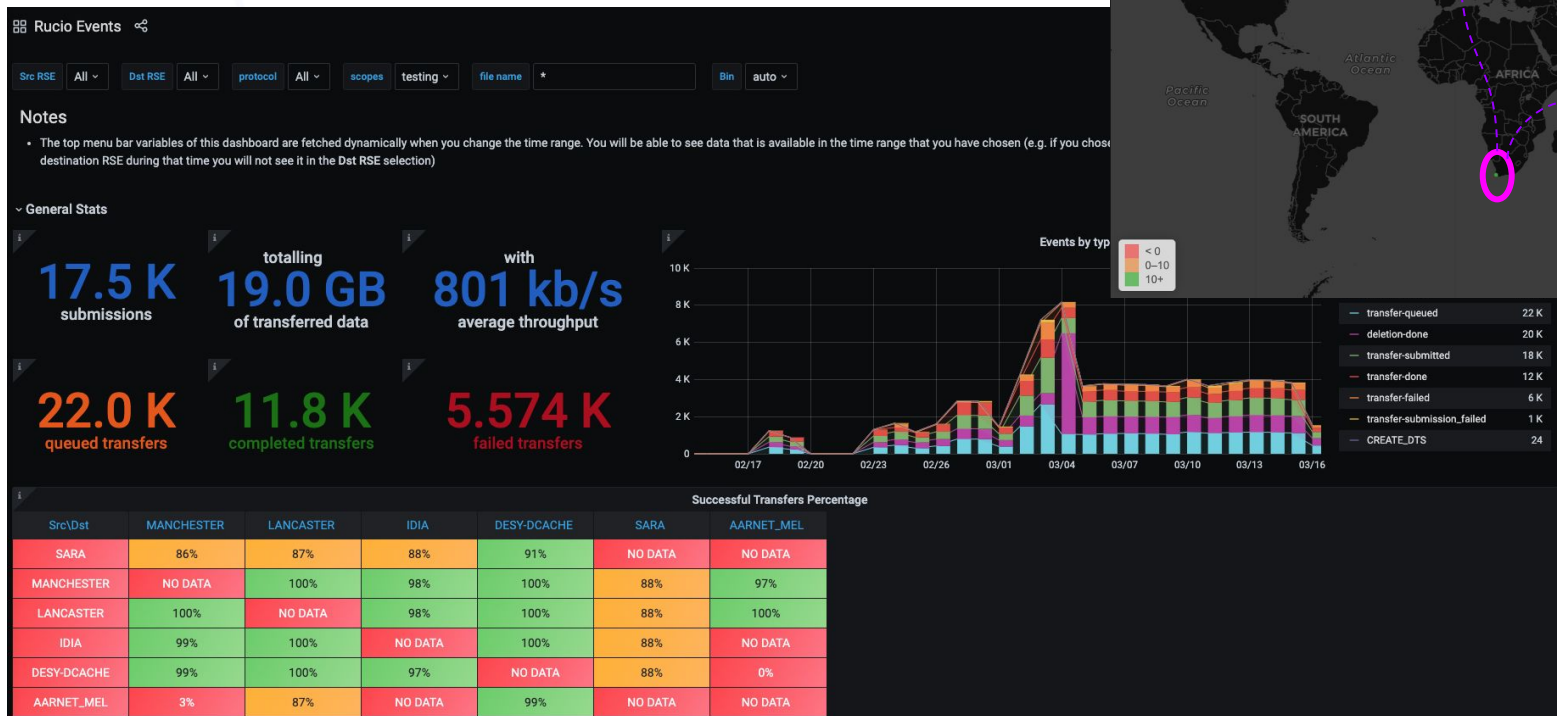
# PIC-CTA Rucio:
## Orchestration of rule creation and automatic removal



- K8s managed Rucio instance deployed at PIC
- RSE set up an RSE on la Palma
- Scripted discovery of pending files to be transferred from a query in the onsite database
- Rucio registration, transfer to a PIC endpoint, and deletion of successfully transferred files from injection site

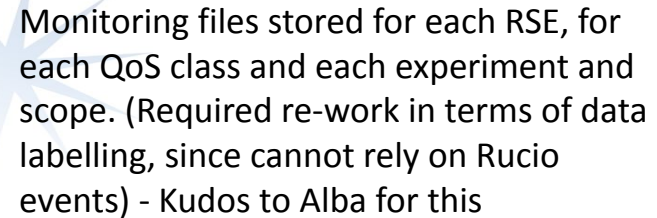- Next steps, update onsite database with transfer times
- Extended real-life tests
- If OK, migrate real systems onto this instance

# SKAO Rucio



*Manual data transfers with Rucio began Feb 18th 2021 and automated tests running hourly across all sites since Feb 23rd*

# Improved monitoring: ESCAPE Rucio



Monitoring files stored for each RSE, for each QoS class and each experiment and scope. (Required re-work in terms of data labelling, since cannot rely on Rucio events) - Kudos to Alba for this

Also now (since FDR) have moved onto longer-lived DB infrastructure

Automated tests continue to run to confirm the health of the ESCAPE Rucio instance and all the components it uses
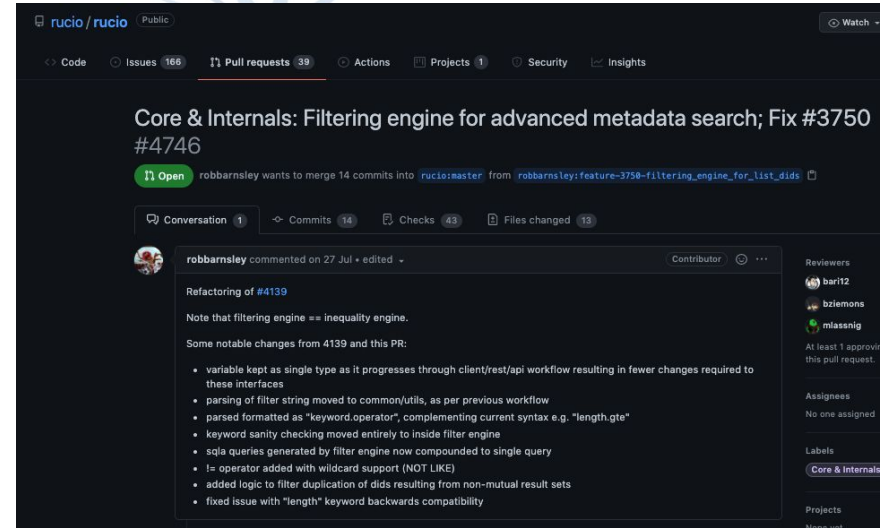
# **Technical development example**

Rucio dev work on Metadata functionality

    Makes metadata more searchable, able to used ranges ( <, >) to select data products

    Moving closer to data discovery functionality

(SKA) Tests written that use subscriptions and metadata to place new data according to predefined rules

# Preparing for DAC21

Renewed energy in DepOps meetings and in ensuring reliability of storage and systems

Once more, getting ready to "look" like a production instance for a short period of time (1 week)

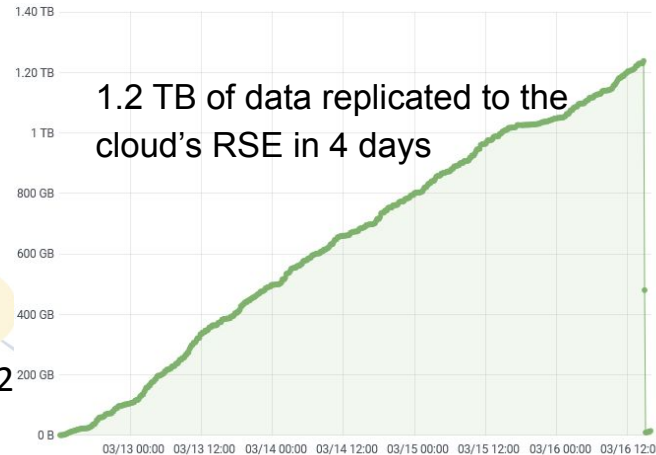We know the drill - FDR has prepared us for this!

# A3-Highlights

- Demonstrate high-availability (HA) IAM deployment.
  - Prove HA deployment by simulating single-server failures.
- Demonstrate token-based access to the Data Lake:
  - Show that the Data Lake works with industry-standard OpenID-Connect.
  - ESFRI communities will be able to upload data and read it without using X.509 certificates.
  - Allow the integration of token-based services; e.g., Datalake-as-a-Service.
- Implement continuous monitoring of Data Lake AuthN/Z configuration:
  - Check that AuthN and AuthZ work consistently across the infrastructure
- Demonstrate fine-grained authorization for selected use cases
  - Controlled Access to embargoed data sets for CMS and CTA.
  - AuthZ policies are centrally managed, but enforced on all storage systems.
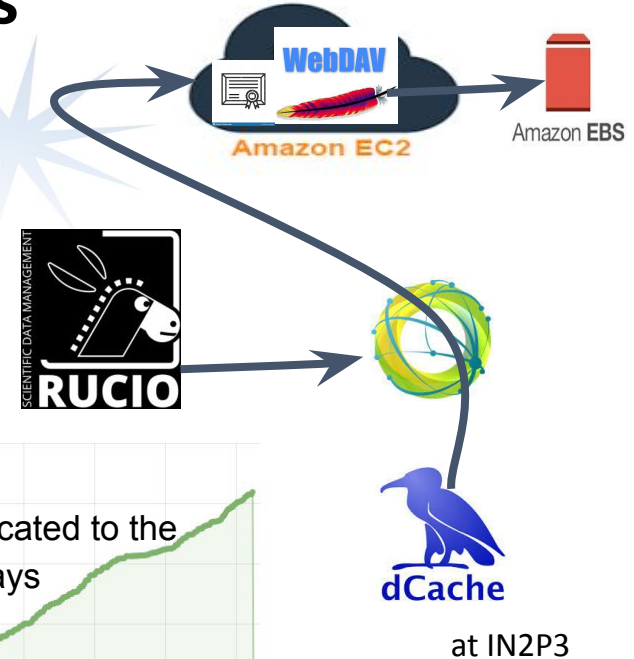
# Data Lake and commercial Clouds



## Set up

1. Set up Domain Name
2. Associate a domain name to the EC2 instance
3. Generate a grid certificate for the given domain name
4. Install the software acting as a RSE* (apache with WebDav plugin) on the EC2 instance
5. Create and map a storage capacity using the Elastic Block Store (EBS) service.
6. Declare the rse in CRIC
7. **Execute a replication rule from dCache to AWS**

*Rucio Storage Element

1.2 TB of data replicated to the cloud's RSE in 4 days

at IN2P3

### Result: IT WORKS!

Ingress cost is very low and negligible compare to EC2

EBS cost is expensive, 80 USD/TB-Month

# Data Lake and HPC (FENIX, PRACE/GEANT/SKA/CERN)

- FENIX infrastructure: https://fenix-ri.eu/
  - Bringing together data repositories and scalable supercomputing systems: JULICH, BSC, CINECA, CSCS, CSC, CEA and EBRAINS
- Collaboration ongoing to enable ESCAPE Data Lake from/to FENIX infrastructure. Three main activities identified:
  - **Identity Management**: ESCAPE IAM and FENIX AAI system as trusted IdPs
  - **Data pre-placement**: the prerequisite to access HPC storage is to use swift protocol. Juelich and CERN developed and implemented swift protocol in FTS.
  - **Software and data repository access via containers**, repositories, etc. need to be understood and addressed
- A data transfer demonstrator with CSCS HPC ongoing. Accounts and storage space created. Token delegation/orchestration ongoing.

# DAC21 and Deliverable 2.3

- DAC21 is our focussed assessment time for the pilot data lake

- Opportunity to explore the use cases and perspectives of different ESFRI partners

- Will once again work as a team to capture the output of the assessments for each ESFRI

**Thanks to all committed people in DIOS: for keeping the focus, promoting teamwork spirit and a daring involvement**

# Come on in, the water's lovely!