



ESCAPE

European Science Cluster of Astronomy &
Particle physics ESFRI research Infrastructures

Conference

ESCAPE to the Future | 25-26 October 2022

Royal Belgian Institute of Natural Sciences | Brussels, Belgium

25 October 2022, 10:20 - 11:15

ESCAPE DIOS - Building a data lake for Open Science



ESCAPE

DIOS | Data Infrastructure
for Open Science



Xavier Espinal
CERN



Gareth Hughes
CTAO



Fabio Hernandez
LSST/Vera Rubin, CNRS / IN2P3



Yan Grande
ASTRON

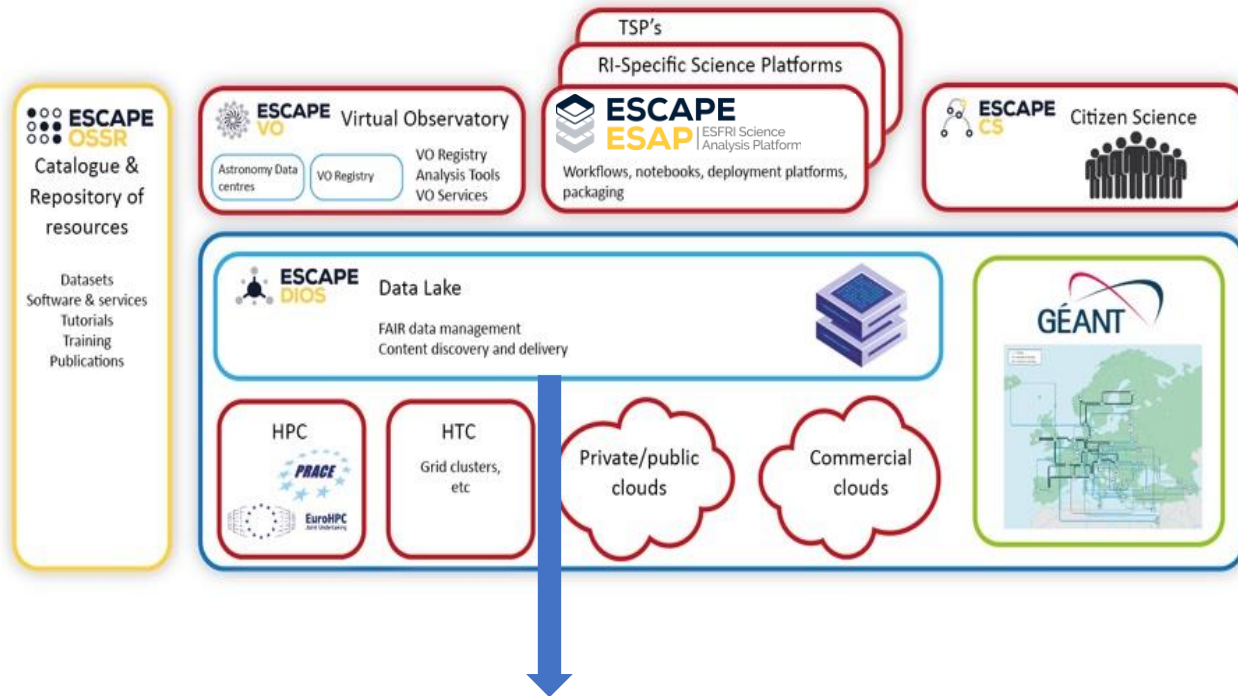


ESCAPE to the Future
25-26 October 2022
Brussels, Belgium

DIOS

Building a Data Lake for Open Science

Xavier Espinal (CERN)
on behalf of the DIOS (WP2) team



The ESCAPE Scientific Data Lake is a **reliable, policy-driven, 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

The ESCAPE Scientific Data Lake is a **reliable, policy-driven, 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

Services operated by the ESCAPE partner institutes

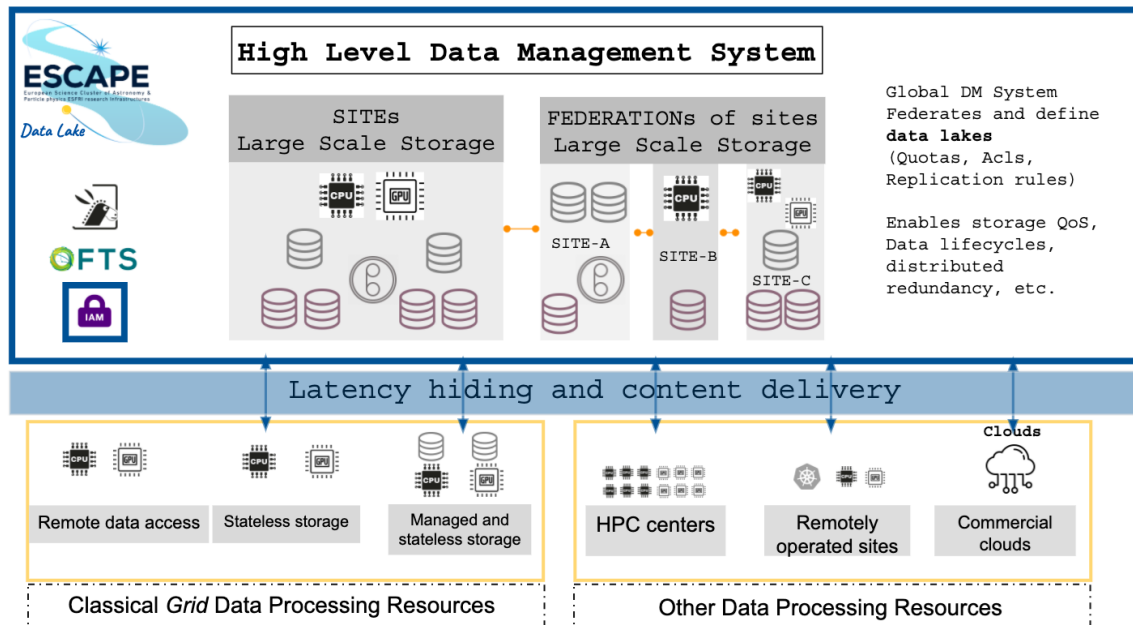
Petabyte scale storage: DESY, SURF-SARA, IN2P3-CC, CERN, IFAE-PIC, LAPP, GSI and INFN (CNAF, ROMA and Napoli)

Data management and storage orchestration (Rucio)

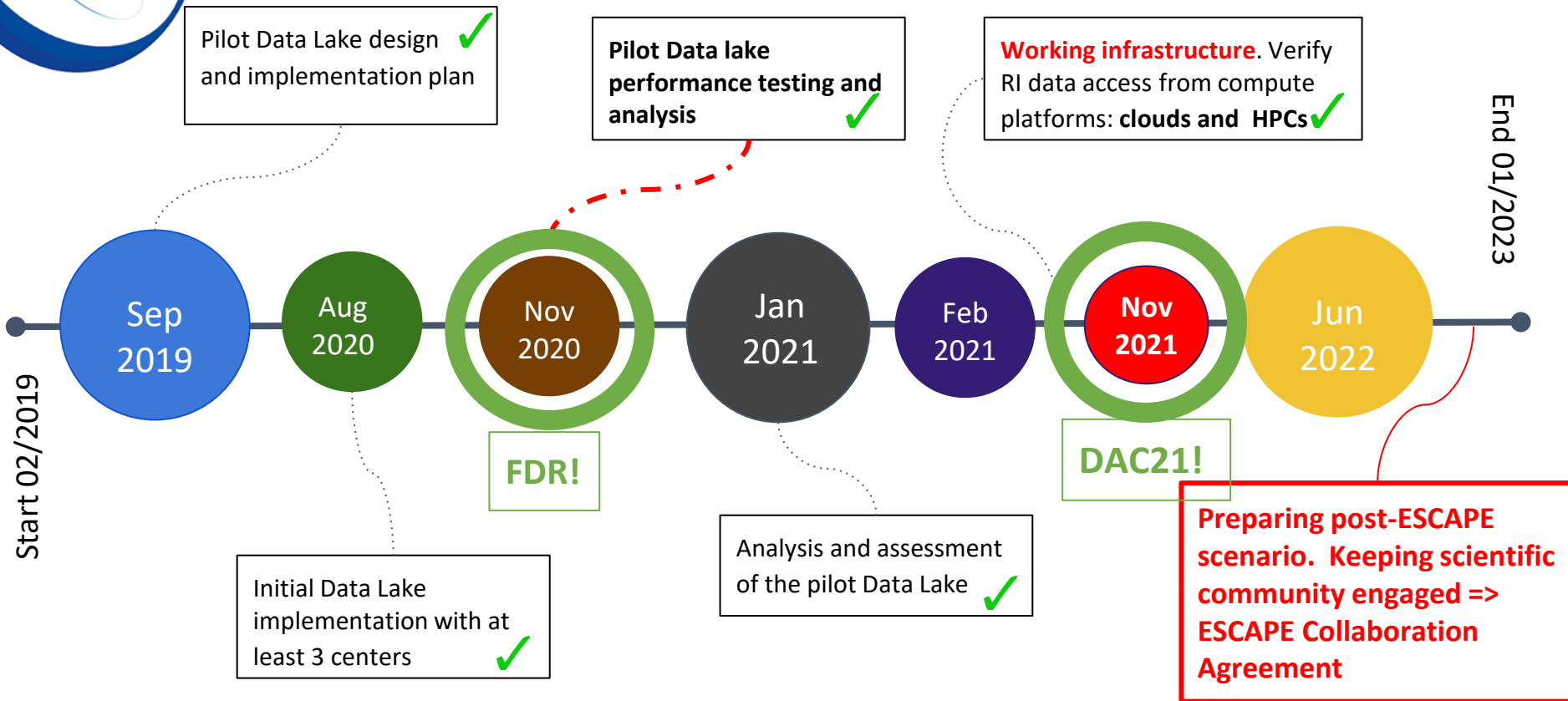
File transfer and data movement services (FTS)

Global Data Lake Information System (CRIC)

ESCAPE IAM: common Auth/Authz/IM (AAI)



ESCAPE DIOS Roadmap



From a Data Lake Pilot to a full Prototype (1/2)

WP2 work plan focused on a continuous assessment and evolution of the pilot Data Lake, with the target to meet ESFRI/RI requirements and resulting in a fully working system

- **Token-based authentication** boosted its integration in the several layers of the Data Lake infrastructure: Rucio, FTS, storages (wip) and integration with other AAI *providers*. Easing user experience with a single and global authentication point
- **Data life-cycle accommodation** ESFRI/RIs users are able to define data replication rules, lifetimes, access policies, data location and storage *quality of service* (adjusting storage cost with data value)
- **Webdav/HTTP** promoted to be the de-facto standard in the Data Lake. The widespread knowledge of HTTP protocols provide a flexible way to interact and integrate with other storage resources, also eases data access from heterogeneous compute platforms and end-user devices
- **Data Management (Rucio) Evolution and Consolidation** channeling feedback from the new scientific communities using Rucio. Discussions on extending metadata capabilities together with ESO. Two extra ESFRI/RI private Rucio instances in operation for SKA and CTA, harmonically using the same global Data Lake storage infrastructure

From a Data Lake Pilot to a full Prototype (2/2)

WP2 work plan focused on a continuous assessment and evolution of the pilot Data Lake, with the target to meet ESFRI/RI requirements and resulting in a fully working system

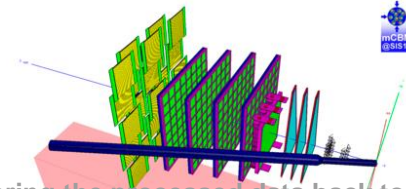
- **Enlarged Data Lake monitoring capabilities** providing real time follow up for data transfers, automated test suite results, resources usage
- **Active Deployment and Operations (DepOps) team** early in the project identified need to share expertise, organised via a well-established meeting. Crucial to consolidate the infrastructure, to foster knowledge transfer and to prepare and drive the data challenges
- **Expanded Data Lake capabilities with user environments** the *Data Lake as a Service* product provides to the end users increased data browse/download/**upload** capabilities, trigger data movement, integrate with local storage, leverage storage caches, etc. Extending functionalities of Analysis Platforms (in conjunction with WP5), and to leverage computing infrastructures (ie. local batch systems and external resource providers)
- **Integration of heterogeneous resources** has been demonstrated, Data Lake interfacing with commercial clouds, public clouds and HPCs

DIOS work plan brought together scientific communities addressing collective goals in a common data infrastructure. The various Data Challenges certified the infrastructure as a fully working system

Putting the system to work: Data and Analysis Challenges (1/3)



- Registration of RAW data acquired by the mCBM detector on FAIR-ROOT
- Ingestion and replication of simulated R3B data
- Ingestion and replication of simulated and digitised raw PANDA fallback data
- Particle-transport and digitisation of Monte-Carlo events
- Live ingestion of simulated data
- Retrieval of stored RAW data from the data-lake, processing of the data and storing the processed data back to

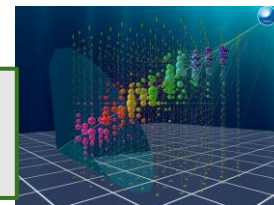


Raw data injected, stored and preserved in the DL. Data processed by users, results are stored back in the DL.



- Ingestion of raw data from the storage at the KM3Net shore station to the Data Lake, and policy-based data replication across the Data Lake infrastructure

Offload data from the storage buffer in the coast, replicate across sites, run data calibration, store back. Data product ready for user consumption



- Long-haul transfer and replication, CTA-RUCIO @PIC: non-deterministic (La Palma) and deterministic (PIC) RSEs and findable in the datalake (using the CTA Rucio instance). Data is
- ⇒ **Next dedicated talk by Gareth!**

Full data re-processing workflow: data injection and distribution, data processing and results consolidation



Putting the system to work: Data and Analysis Challenges (2/3)



- Exercises (data production, replication and documentation) before and during the DAC21. Include the creation of datasets for real-kind final user analysis examples using current open-access datasets. ~200*10 = 2000 files uploaded

Large experiment demonstrating open data capabilities

(<http://opendata.atlas.cern/software/>). Testing and validating the reading access of the samples via de Jupyter rucio extension, and running multiple analysis pipelines.



⇒ Next dedicated talk by Gareth!

Data management from remote locations

- Long haul raw data ingestion and replication. Data is successfully transferred from the MAGIC telescopes to the Data Lake, file deleted on the telescope storage
- Data transfer monitored. Data can be discovered using the CTA-RUCIO instance.
- the reading access of the samples via ammapy library.



⇒ Next dedicated talk by Yan!

Full-cycle scientific data management and data processing

- Ingestion of LOFAR data from a remote site to the Data Lake. Data transfer and replication into off-site storage, after successful replication delete data at the source
- Process data in the Data Lake, as an external location, combine results with other astronomical data to produce
- Include a read-only RSE to a location outside the data lake. Get data from there into the DL.
- Extending use cases by using larger files and leveraging several QoS, running all processing in the DLaaS, requiring a the availability of specific LOFAR software in the DLaaS.

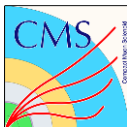


Putting the system to work: Data and Analysis Challenges (3/3)



- Data replication. Data in correct place in timely manner.
- **Long haul data replication. SKAO Rucio (Australia and South-Africa to UK RSEs), using the RUCIO SKA instance.**
- End-to-end proof of concept data lifecycle test, AUS/SA to northern hemisphere sites

Global-scale Data Management



- **Multi purpose Analysis Facility PoC with data access via DASK (workload orchestrator) leveraging computing at Marconi (HPC) and large batch clusters**
- Access control for embargo data, test in CNAF and DESY
- Content delivery and caching: XCache Protocol Translation: xroot internal vs http External for Data Lake data transfer. Performance comparisons for Analysis workflows

DL interface with local and heterogeneous resources, CDN and caching



- Simulate replication of one night's worth of raw images data between two Vera C. Rubin data facilities, perform the exercise several times. Each iteration is composed of 15TB, 800k files, ideally to be replicated

⇒ Next dedicated talk by Fabio!

Leverage telescope local storage data replication to fulfill daily data management cycles





ESCAPE to the Future
25-26 October 2022
Brussels, Belgium



Cherenkov Telescope Array Observatory (CTAO) & MAGIC Collaboration

Gareth Hughes (CTAO)
Matthias Füßling (CTAO)



Agustin Bruzzese
Jordi Delgado
Matthias Füßling
Frederic Gillardo
Gareth Hughes
Gonzalo Merino
Nadine Neyroud





- Two 17 m diameter Imaging Atmospheric Cherenkov Telescopes
- Located on the Canary island of La Palma at 2200 m a.s.l.
 - **MAGIC-I: since 2004**
 - **MAGIC-II: since 2009**
- Energy range 50 GeV – 50 TeV
- Major upgrades of both cameras electronics in 2013
- Multiple hardware upgrades

- Significant contributions to:
 - **Galactic and extra galactic science**
 - **Fundamental physics**
 - **Transients**
 - **esp. ~TeV Gamma-ray Bursts**

Cherenkov Telescope Array Observatory

- The first ground-based gamma-ray **observatory**
 - 5-10x increase in sensitivity
 - Broad energy range
20 GeV to 300 TeV
 - Will serve large **user community** data & science tools using FAIR principles
 - **Proposal-driven** observatory
- 30-year lifetime
 - 6 PB of data per year
 - Significant effort in maintenance and operational cost optimization
- Two arrays, One Observatory (whole sky)
 - Inter-site coordination
 - Uniform approach to scientific operations
 - **+ 4x data centres**
 - PIC (Spain), CSCS (Switzerland), Frascati (Italy) & DESY (Germany)



CTAO

Setup for ESCAPE tests

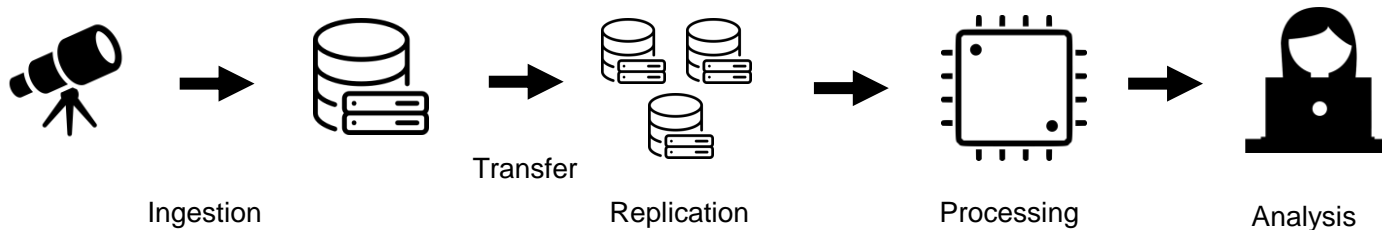
- Ability to produce & process **large volumes of data**
- Observatories are by their nature often in **remote** locations
- Data **transfer, storage and processing** are key



Koji Noda



Use Cases: Telescope to Scientist



Long-haul Transfer :

- automatic data detection
 - data transferred and replicated
 - transfer monitored
 - automatic deletion at source
- discoverable
- differing QoS

(Re)Processing:

- Rucio and DIRAC integration
- the output data is findable in the data lake

Analysis:

- find data and workflows
- using both the ESAP and DLaaS



La Palma to PIC

Raw data from telescopes to data centres

- PIC & LAPP deployed ESCAPE-Rucio RSEs
- PIC deployed CTA-Rucio instance for DAC21 Challenges
 - Deployed using Kubernetes
 - Different setups and technologies tested
 - Two different protocols tested xrootd & gridFTP
 - Deterministic and non-det RSEs
 - Monitored using elastic search and Grafana
- Tests took place between La Palma and PIC (Spain)
 - Realistic test conditions, 1 Gbps connection (now upgraded to 10 Gbps)
 - Rucio Storage Elements placed in La Palma
- Rucio RSEs tens of TBs transferred in a <24 hrs window
 - Observations take place at night
- Important requirements checked
 - Data replicated to tape storage
 - Data detected and deleted correctly



Long-haul transfer test were successful

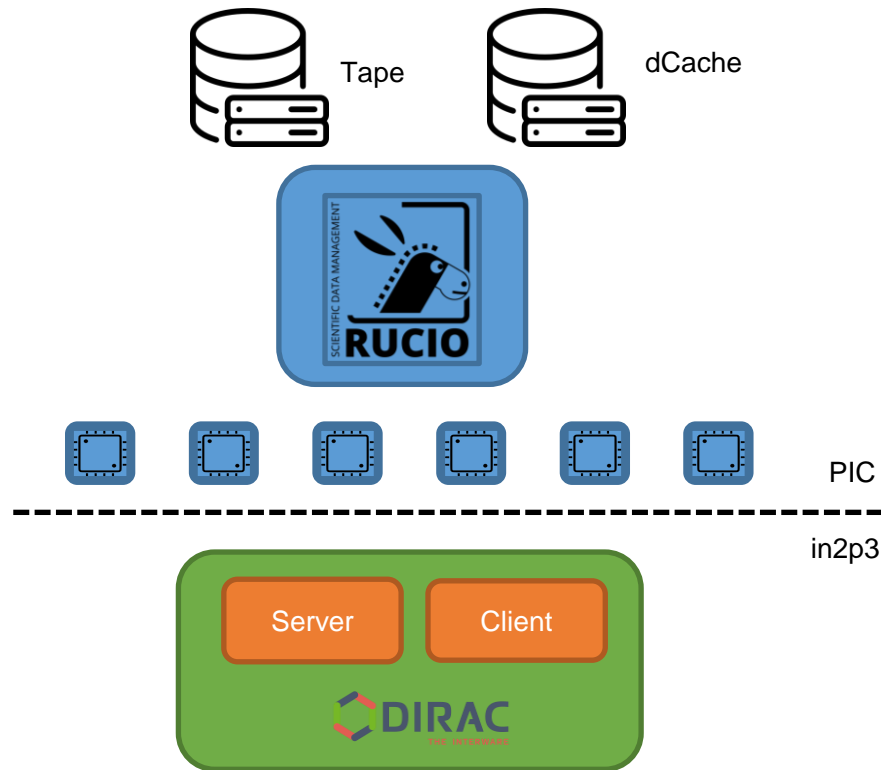


(Re)Processing DIRAC

Raw data processed to be made science-ready

- Using the CTA-DIRAC Workload Management System
 - Starting with the Belle-II Rucio-DIRAC plugin
- **Successes**, able to:
 - Ingest data using DIRAC
 - Launch test jobs on worker nodes
 - Launch CTA production scripts
 - Read data from tape
- **Futures work**:
 - Ingest using Rucio directly
 - Identify and implement the functionality required to integrate (CTA)DIRAC with Rucio

Processing Use Cases learnt a lot: Future collaboration with other ESFRIs and software developers

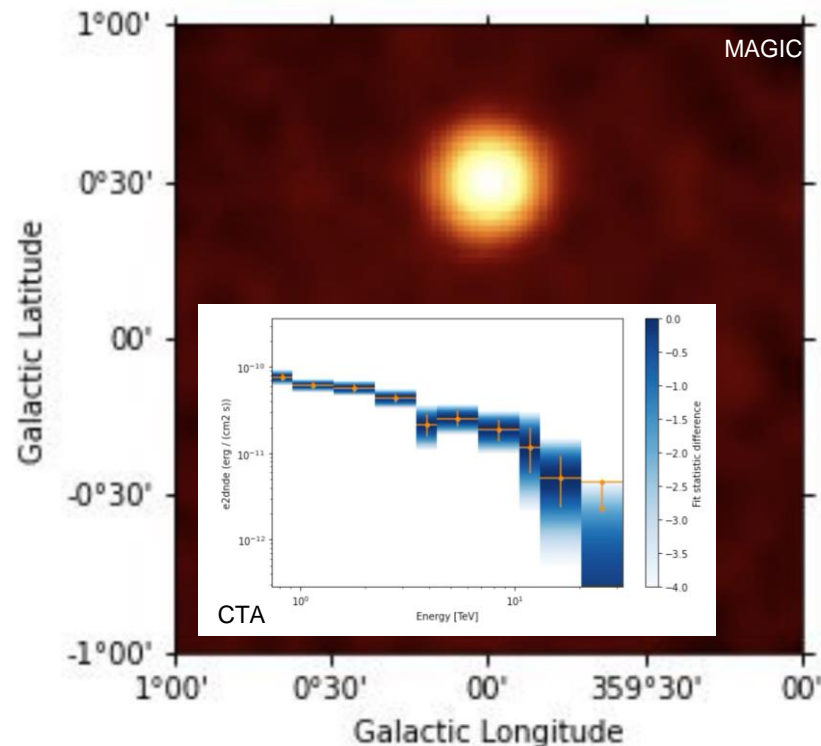




Science-Ready Data used by Scientists

- Data were analyzed interactively using both the **ESRFI Science Analysis Platform (ESAP)** and the **Data Lake as a Service (DLaaS)**
- ESAP allows workflows and facilities to be identified and selected
- DLaaS makes it easy to find and replicate data

Analysis performed on two different platforms by both **MAGIC** & **CTAO**



MAGIC upgraded the data transfer system from La Palma to PIC using Rucio client and new scripts.
The system is connected to the on-site and data transfer databases, and it's currently in production mode.



CTAO will have 4 data centers for distributed storage, processing data and to host all software services.
CTAO will evaluate Rucio as the baseline technology for its bulk data archive.

- Immediate: Tests are ongoing:

- Refine file deletion at source & second replication site
- Longer range transfer tests (Japan)
- Range of file sizes
- Priority data products
- Test new faster connection to La Palma

- Near future: common topics of interest

- Workload Management integration
 - DIRAC-Rucio collaboration
 - CTAO, Belle-II, KM3NeT, DUNE, SKAO, ...
- Interoperability
 - metadata
- Embargoed data
 - A&A and tokens

- We look forward to further work together on these topics in future collaborations

- **ESFRIs speaking with one voice benefits all**

Thank You

● CTAO

- Luisa Arrabito
- Karl Kosack
- Nektarios Benekos
- Federico Ferrini

● DIOS work package for their active support and engagement

- Xavier Espinal
- Rosie Bolton
- Riccardo di Maria
- Rizart Dona
- Alba Vendrell

● PIC Operations Team

- Vanessa Acín - Network
- Elena Planas - dCache
- Esther Acción - Enstore
- Ricard Cruz - Kubernetes
- Carles Acosta - xrootd



ESCAPE to the Future
25-26 October 2022
Brussels, Belgium

**Preparing the most ambitious
astronomical survey ever
attempted with ESCAPE**

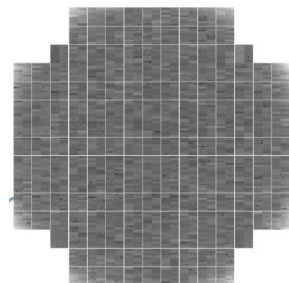
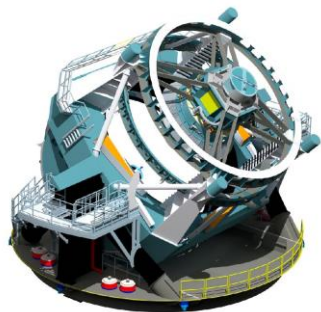
Adrien Georget, Fabio Hernandez, Lionel Schwarz

Vera C. Rubin Observatory - French Data Facility

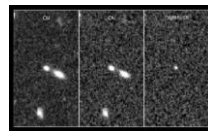
IN2P3 / CNRS computing centre, Lyon (France)



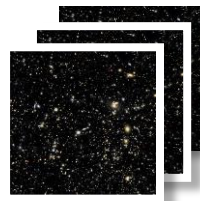
Rubin Observatory Legacy Survey of Space and Time



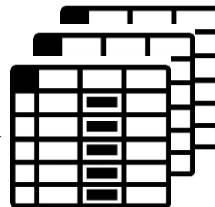
raw images



alerts



science-ready images



astronomical catalog



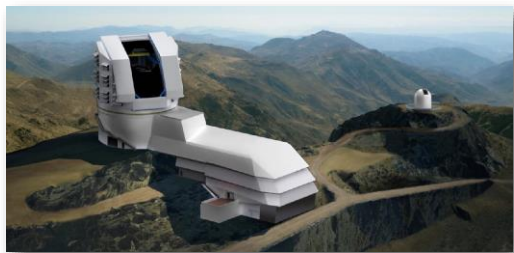
science collaborations



LSST aims to deliver a catalog of 20 billion galaxies and 17 billion stars with their associated physical properties

LSST Overview

OBSERVATORY



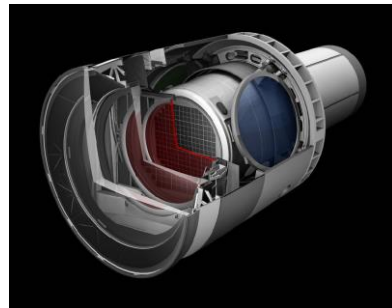
southern hemisphere |
2647m a.s.l. | stable air |
clear sky | dark nights |
good infrastructure

TELESCOPE



main mirror \varnothing 8.4 m (effective
6.4 m) | large aperture:
f/1.234 | wide field of view |
350 ton | compact | to be
repositioned about 3M times
over 10 years of operations

CAMERA



3.2 G pixels | \varnothing 1.65 m |
3.7 m long | 3 ton | 3
lenses | 3.5° field of view
| 9.6 deg^2 | 6 filters *ugrizy*
| 320-1050 nm | focal plane
and electronics in
cryostat at 173K



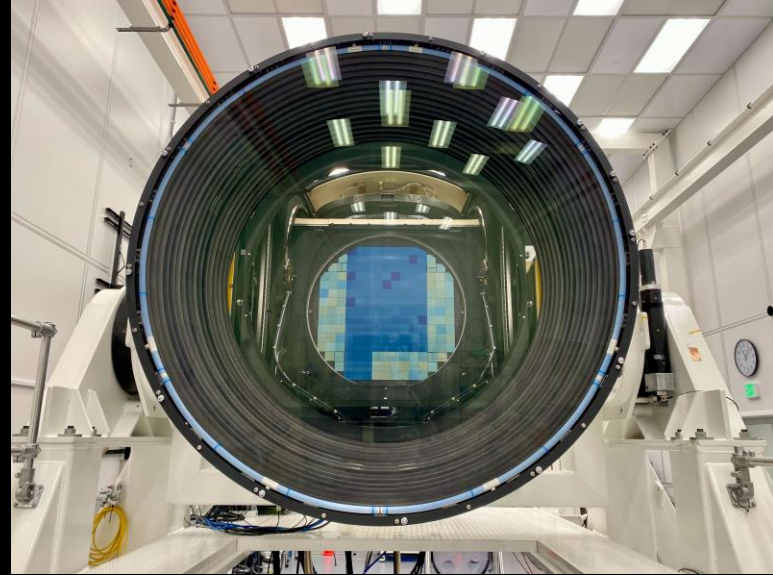
Source: Rubin Observatory

Raw data

6.4 GB per exposure (compressed)
2000 science + 500 calibration images per
night
300 nights per year
16 TB per night, ~5 PB per year

Aggregated data over 10 years of operations

image collection: ~6 million exposures
derived data set: ~0.5 EB
final astronomical catalog database: 15 PB



Source: Rubin Observatory



Cloud

EPO Data Center

Dedicated Long Haul Networks

Two redundant 100 Gb links from Santiago to Florida (existing fiber)

Additional 100 Gb link (spectrum on new fiber) from Santiago-Florida (Chile and US national links not shown)

UK Data Facility ROE, Edinburgh, UK

Data Release Production (25%)

US Data Facility SLAC, California, USA

Archive Center
Alert Production
Data Release Production (25%)
Calibration Products Production
Long-term storage
Data Access Center
Data Access and User Services

French Data Facility CC-IN2P3, Lyon, France

Data Release Production (50%)
Long-term storage

HQ Site AURA, Tucson, USA

Observatory Management
Data Production
System Performance
Education and Public Outreach

Summit and Base Sites

Observatory Operations Telescope
and Camera
Data Acquisition
Long-term storage
Chilean Data Access Center



LSST usage of the ESCAPE data lake

- Focused on **data replication**
- Performed inter-data facility replication of one night's worth of raw data, repeatedly over 5 consecutive days
 - *realistic data set: 4,000 exposures, 800k files, ~15 TB*
 - *replication time budget: 12 hours*
 - *driven by Rucio and FTS, involving storage endpoints connected to ESCAPE data lake*
 - *data flow: CERN → CC-IN2P3*
- Results
 - **reproducible replication** of the entire data set performed in less than 8 hours, **without errors**



Benefits of using ESCAPE data lake

- ESCAPE data lake infrastructure was an instrumental sandbox
 - *to lower the barrier for science projects to **evaluate sophisticated data management tools** via hands-on experience*
 - *to **accelerate adoption** of robust, proven tools and good practices by providing a ready-to-use, well-maintained, monitored, flexible infrastructure*
- ESCAPE provided a forum to share experience
 - *understanding how **other science projects manage their data** using the same tools is inspiring*
 - *getting **advice** and previews of upcoming technologies **from developers and operators** of those tools is extremely valuable*
 - *science projects **provide input to developers and operators** on their specific, sometimes atypical and very creative, use cases*

Post-ESCAPE activities

- LSST deployed early production instances of Rucio and FTS at the US data facility
 - *currently performing transatlantic data replication exercises among the 3 data facilities contributing to Rubin*
 - *ongoing work for deploying a platform for logging and monitoring of data replication activities*
- Developing mechanisms for integrating Rucio to the LSST-specific data management components

ESCAPE to the Future

25-26 October 2022
Brussels, Belgium

DIOS: The LOFAR view

Yan Grange, Vishambhar Nath Pandey



The team

- Yan Grange
- Vishambhar Nath Pandey

“small team, big ambitions”

Many thanks to the other ASTRON ESCAPE participants (Klaas Kliffen, John Swinbank) and Hanno Holties for reflection, discussion and input.

LOFAR in one slide



LOFAR data management

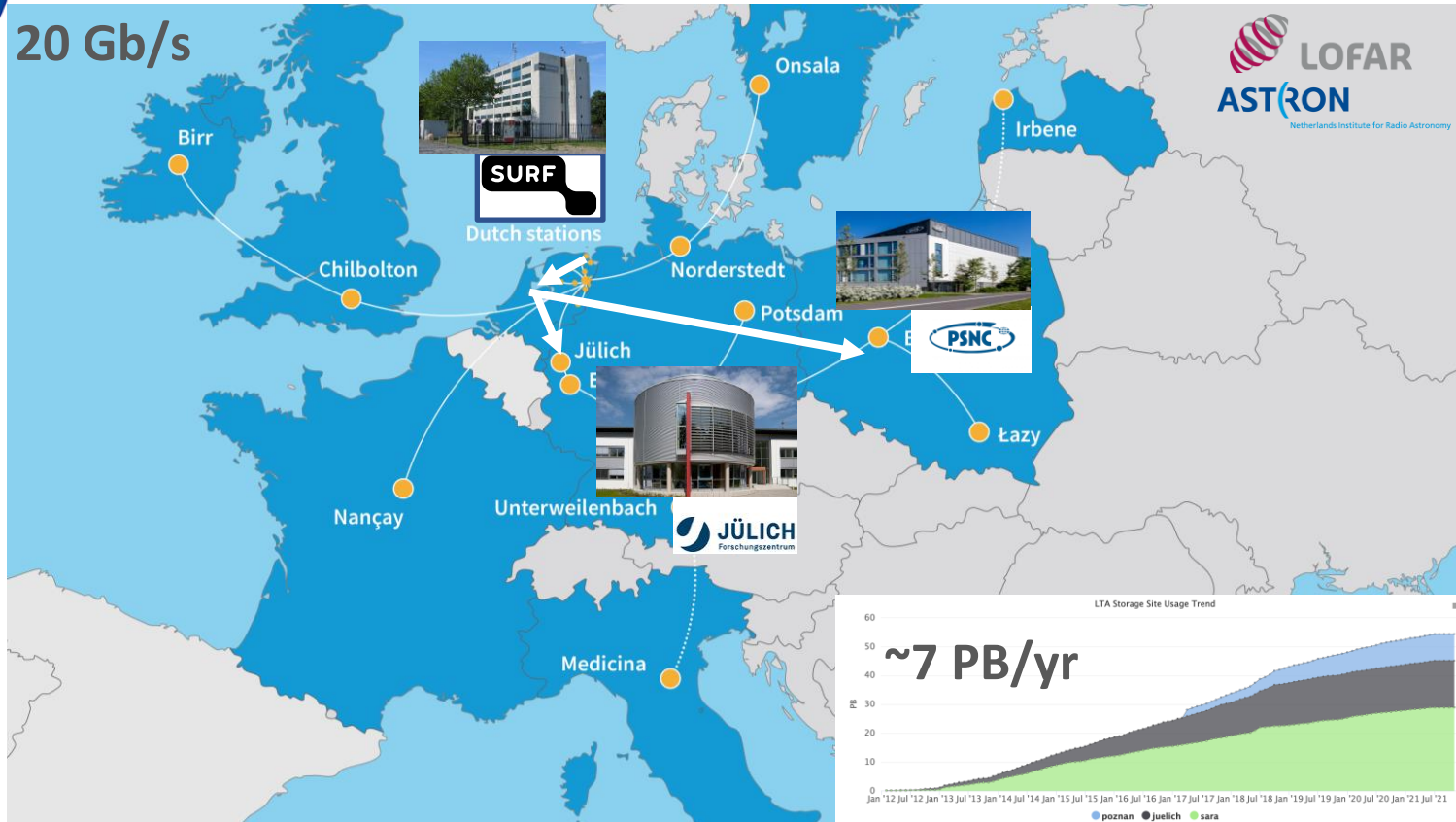
228

Gb/s



LOFAR data management

20 Gb/s



LOFAR in DIOS

DAC21



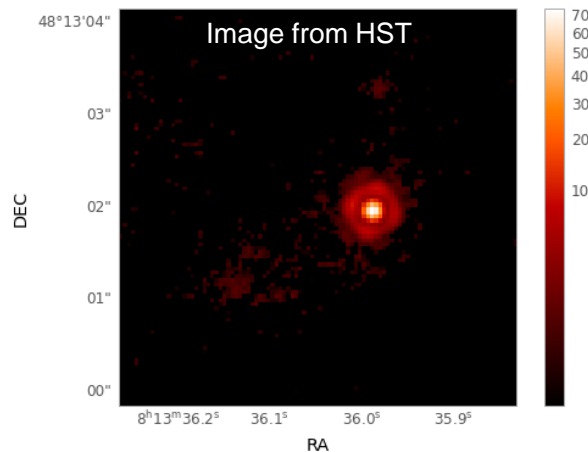
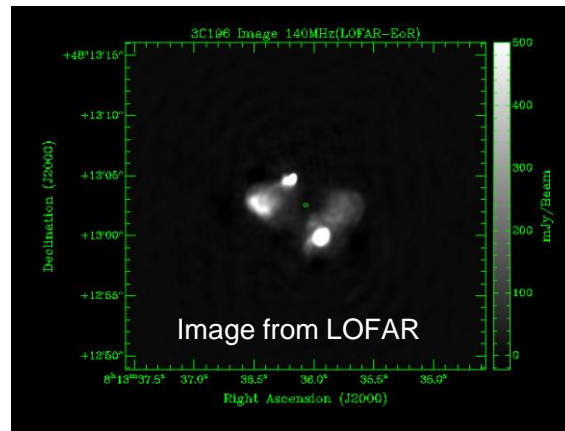
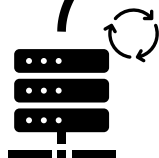
1 day of ingest



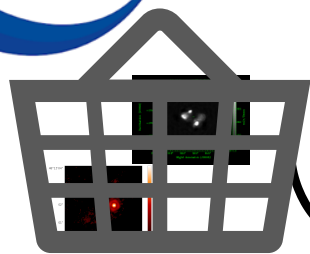
ESCAPE
DIOS | Data Infrastructure
for Open Science



ESCAPE
VO | Virtual
Observatory



LOFAR in DIOS (2)

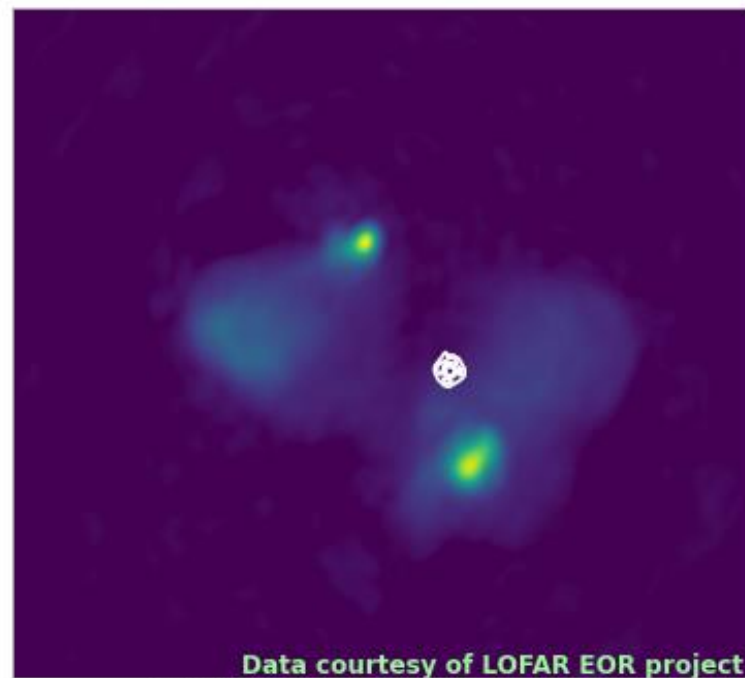


48°13'10"

DEC

05"

00"



8^h13^m36.5^s

36.0^s

35.5^s

RA

Data courtesy of LOFAR EOR project

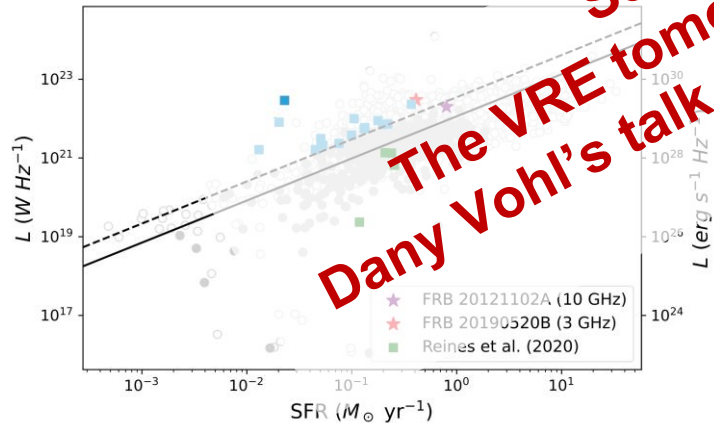
ESCAPE input to LOFAR

- LOFAR operational use case has similarities to the Data Lake but limited set of storage locations compared to e.g. particle physics.
- Rule-based data management with life cycle could be of use of a future operational model of our archives.
- Content delivery and caching for access to remote data may be useful addition to processing at larger compute systems that are further away from the data archive.
 - User processing that may happen on different types of hardware environments

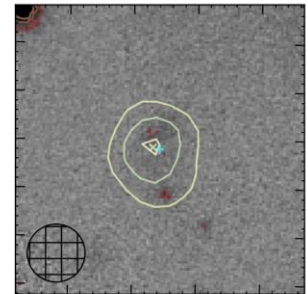
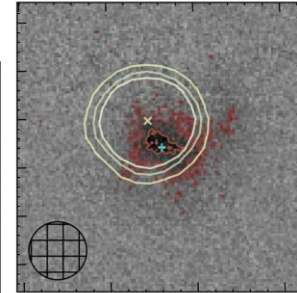
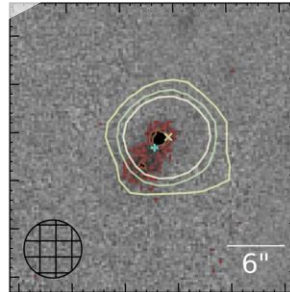
Targeted search for **compact radio sources** coincident with **dwarf galaxies** above expected **SFR**

LoTSS 2nd data release (Shimwell *et al.* 2022; 144 MHz)

- **> 4 million radio sources** over $\sim 5500 \text{ deg}^2$ covered
- **6 arcsec resolution** for RMS sensitivity of $20 \mu\text{Jy/beam}$
- **0".2 astrometric accuracy** (comparable to optical surveys)
- Point source completeness to 90% at 0.8 mJy/beam



See two talks tomorrow:
The VRE tomorrow 9:30 by Enrique Garcia,
Dany Vohl's talk in the Extreme Universe session at 10:55!



Vohl et al. In prep.

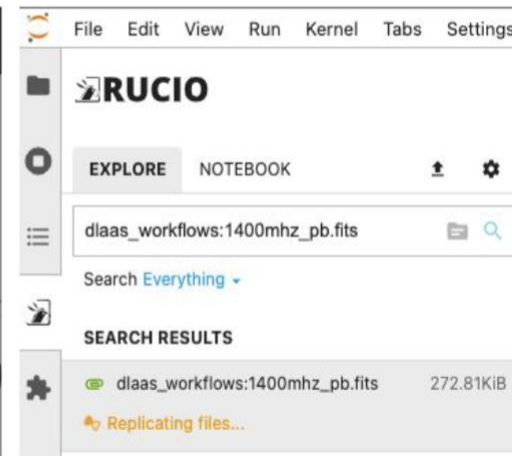
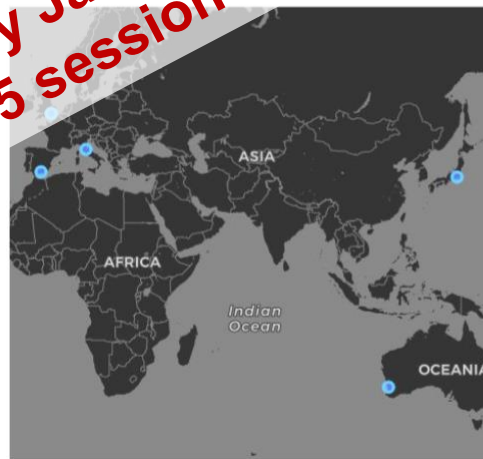
Tables from the VO (CEVO)
 Large table handling, (ESAP?)
 Images on the Data Lake (DIOS)

SKA Regional Centres

- LOFAR has formally been a member as an SKA pathfinder
- ASTRON involved in the SKA Regional Centre prototyping process.
- Several ESCAPE technologies investigated in the process (ESAP, Rucio, DLaaS, IAM)

Transfer success site matrix

Src\Dst	SPSRC_STORM	JPSRC_STORM	IMPERIAL	CNAF	AUSSRC_STORM	STFC_STORM
STFC_STORM	100%	100%	100%	100%	100%	NO DATA
SPSRC_STORM	NO DATA	100%	100%	100%	100%	100%
JPSRC_STORM	100%	NO DATA	100%	100%	100%	100%
IMPERIAL	100%	100%	NO DATA	100%	61%	100%
CNAF	100%	100%	100%	NO DATA	100%	100%
AUSSRC_STORM	100%	100%	100%	100%	NO DATA	100%



The ESCAPE community

- The collaboration within the project led to contacts with development teams of tooling (e.g. dCache, Rucio)
 - Comparing workflows with other science use cases as also been very enlightening (similarities, differences)
- The experiments with the Data Lake have been a good exercise to think about how the concept of rule-based data management in a hybrid storage environment fits within the LOFAR use case.
- Through the SKA work, the collaborations may lead to knowledge and technical developments that feed back in to LOFAR.

Thanks to all the project members that worked
with us, especially the WP2 team!

DIOS and “ESCAPE to the FUTURE”

- In a challenging future ESCAPE and DIOS acted as an **anchor** for the Scientific Community. An all-together spirit addressing **RI’s upcoming needs** in Data Management, Access and Analysis
- ESCAPE DIOS is providing a fully working framework to **address and test** novel data management and data access tools and models, giving the opportunity to **influence and steer** its development
- ESCAPE Collaboration Agreement establishes a **fundamental framework** for further joint projects and collaborations. With the right approach for success: RIs, sites and service providers working **together**

Action-1

Keep the current engagement with RIs by focusing efforts on well identified common needs ^(a), e.g.:

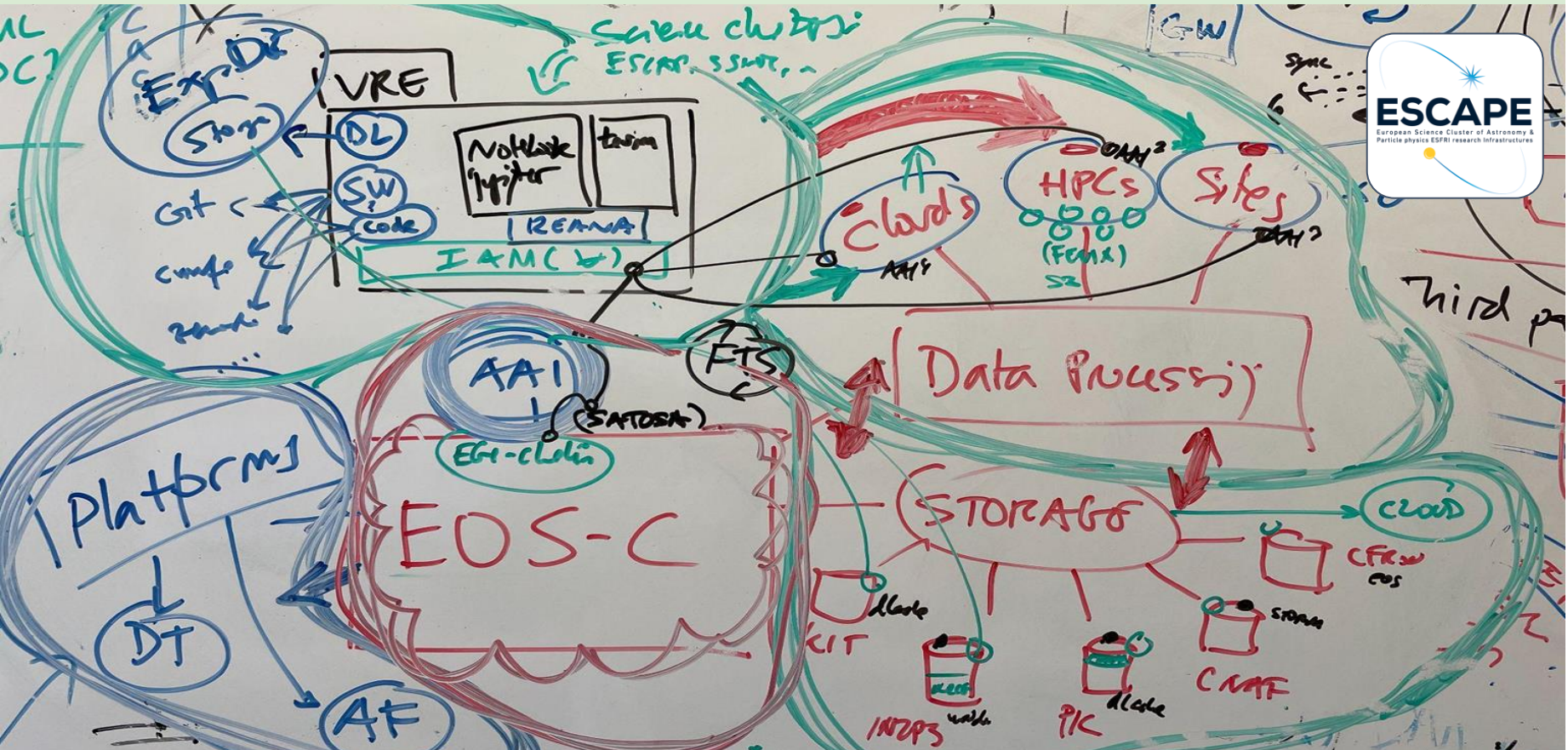
- High level Data Management and File transfer services
- Data Lake model: low latency data delivery, integration with Analysis Platforms and Analysis Frameworks
- Consolidation of a global AAI “framework”

Action-2

Expand collaborations and foster involvement with other Scientific Communities. Maintain and strengthen collaborations with related EC initiatives and projects

^(a)Topics identified together with the DIOS community during the [3rd DIOS Workshop](#) (March 2022)

Thank you! - looking forward to the exciting challenges ahead



*This is a whiteboard scribble during an informal discussion in my office last week