# Advancing Large Scale Scientific Collaborations with Rucio

Hugo González Labrador on behalf of the Rucio project

*CERN | IT Department*

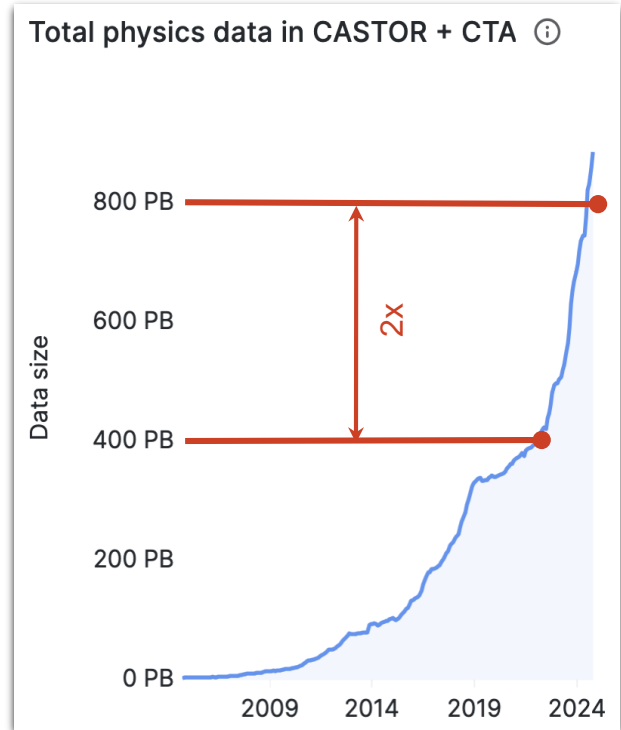CHEP 2024 | 21st October | Kraków

# Taming the data deluge

Big complex machines produce vast amounts (Exabytes) of data

> CERN experiments have generated more data in the last ~2 years than ever before

New instruments will be ready by the end of the decade with rising storage needs (peaking at EBs/year):

> *HL-LHC, SKA, KM3NET, DUNE, …*

A formal solution is required to tame these data volumes



Total physics data in CASTOR + CTA ⓘ

# Rucio in a nutshell

Rucio provides a mature and modular scientific **data management federation**

    **Seamless integration** of **scientific and commercial** storage and their network systems

        Data is stored in **global single namespace** and can contain **any potential payload**

        Facilities can be **distributed at multiple locations** belonging to **different administrative domains**

        Designed with **more than a decade of operational experience** in very large-scale data management
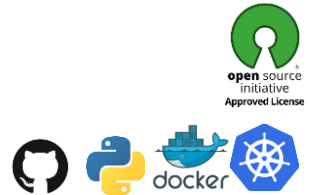
Rucio is location-aware and manages data in a heterogeneous distributed environment

    Creation, location, transfer, deletion, annotation, and access

    **Orchestration of dataflows** with both low-level and high-level policies

Rucio is free and open-source software licenced under *Apache v2.0*

Open community-driven development process

# What I like about Rucio

Reduce human operations to manage the data by expressing what you want, not how you want it
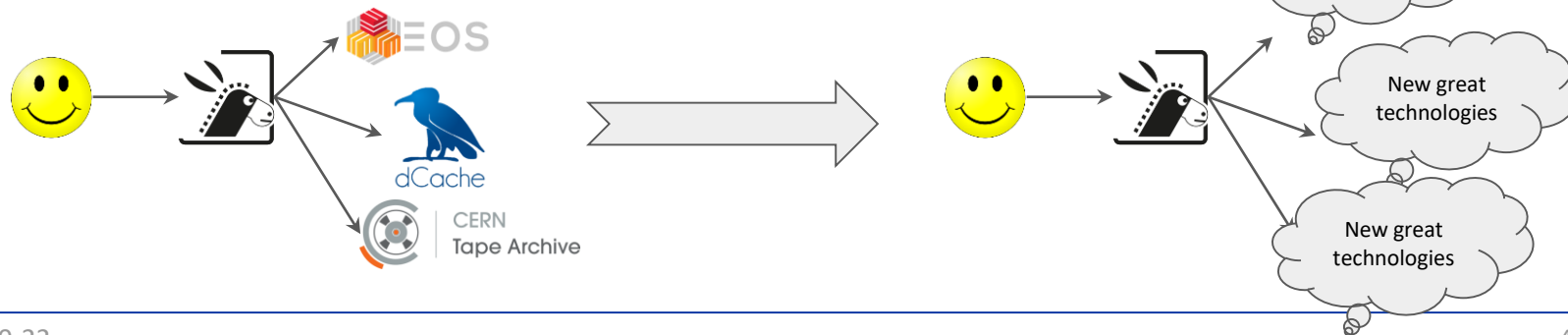
e.g., *"Three copies of this dataset, distributed across MULTIPLE CONTINENTS, with at least one copy on TAPE"*

e.g., *"One copy of this file ANYWHERE, as long as it is a very fast DISK"*

Provide another level of indirection through a scalable catalog of data that abstract the underlying:

Protocols: S3, WebDAV, XROOTD, …

Technologies: dCache, EOS, CTA, Storm ,…

# Where is Rucio being deployed?



Rucio is **used in 20+ scientific collaborations**

De-facto scientific data management tool for HEP and other sciences

Used by CERN (participating) experiments: ATLAS, CMS, Neutrino Platform (DUNE, ICARUS), …

Others: Belle II, LIGO/VIRGO, CTAO, Vera Rubin Observatory, …

NEW: SKA, KM3NET, IHEP, XENONnT…

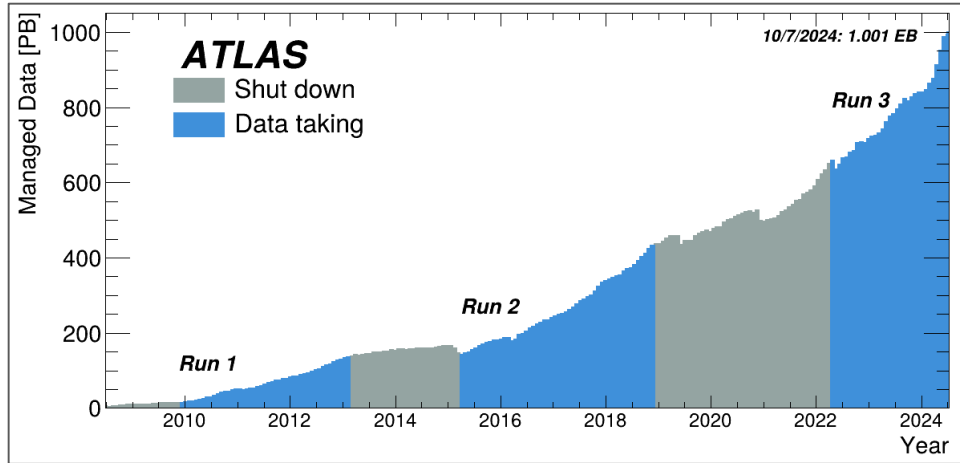Rucio manages over **2 Exabytes of data** across all its community deployments

Rucio **proven to work** to very challenging data environments

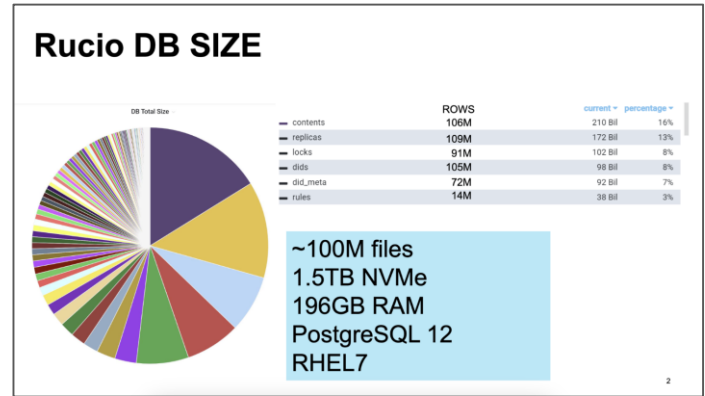# Proven to work in very challenging environments

Rucio@ATLAS: 1+ Exabyte, 120 data centres

Rucio@Belle2: ~100M files on PostgreSQL



Rucio for ATLAS Computing, Mario Lasnigg, 7th Rucio Workshop



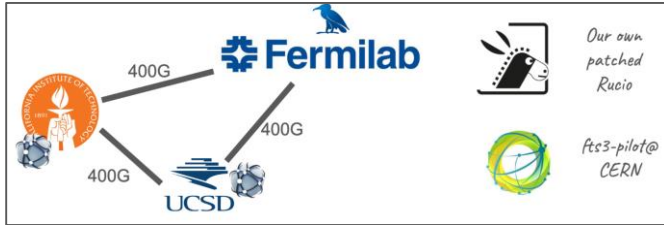PostgreSQL operation for RUCIO service at BN, Hinori Ito, 7th Rucio Workshop

Rucio@CMS, see talk: "Recent Experience with the CMS Data Management System" | 23 Oct 2024, 13:30
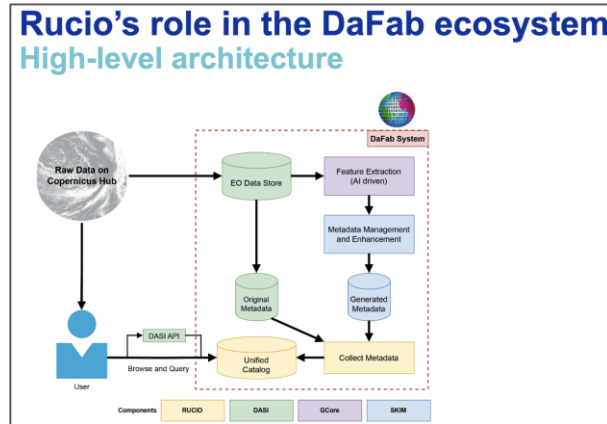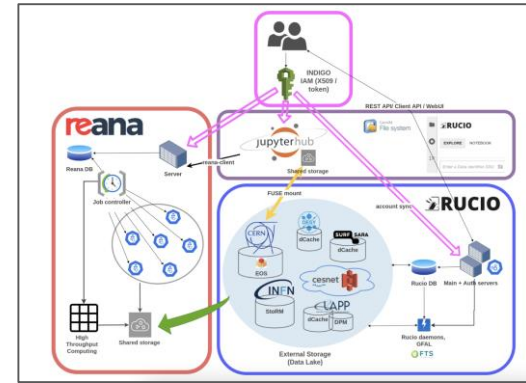
# Breeding ground for innovation

Rucio is an enabler technology in various EU and US scientific projects, for example:



**Integration between Rucio and Sense, Diego Davila, 7th Rucio Workshop**



**DaFab: Extending Rucio for Enhanced Earth Observation Data Management, Dimitris, 7th Rucio Workshop**



**Data discovery, analysis and reproducibility in Virtual Research Environments, Enrique, CHEP 2024**

# Challenges and opportunities

**Authentication**: mix of x509/Macaroons/"Tokens":

Requires smooth transition as it affects many other interlinked systems (transfer tools, RSEs,…)
DC2024 was a good first battle ground. Refer to "The WLCG Data Challenge" plenary tomorrow

**Integration with workflow management solutions:** For significant deployments, it may be beneficial to automate the transfer of the data to the closest computing facility where the user launches the jobs. Dirac interware may be the best candidate as it's already proven to work for Belle2

**Integration with external catalogues:** While Rucio provides an scalable catalog for metadata, often times the experiments have their own catalogues (pre-Rucio) and a synchronization job needs to be added
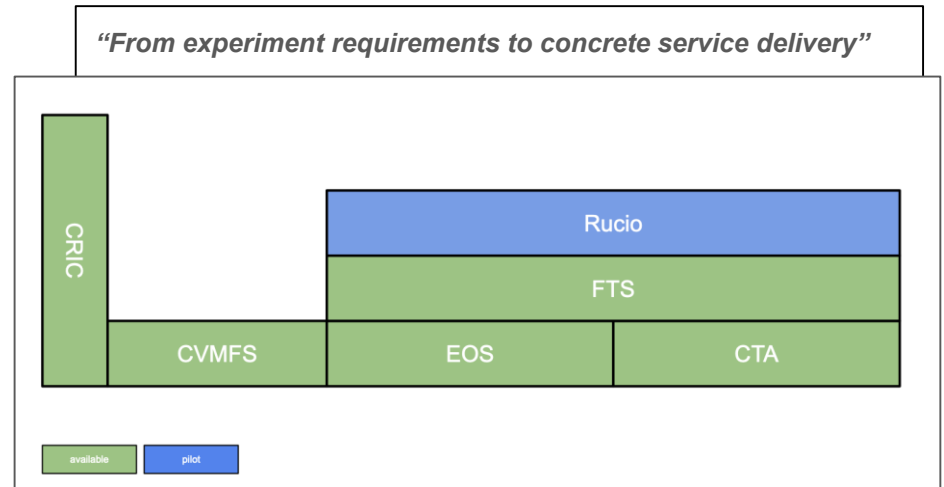
# Increasing cooperation in the Rucio community

The CERN IT department is a new member of the Rucio community since early 2024, giving a hand in 3 axis:

1. Code development

2. DevOps for ATLAS and CMS

3. Managed Rucio service for Small and Medium Experiments (pilot)



*"From experiment requirements to concrete service delivery"*

# Summary

Rucio is an open, reliable, and efficient data management system

      **Supporting the world's largest scientific experiments,** but also a good match for smaller sciences

      Extended continuously for the growing needs and requirements of the sciences

Strong cooperation between physics and multiple other fields

      Diverse communities have joined, incl. astronomy, atmospheric, environmental, …

      **Mutually supportive community**

Benefit from advances in both scientific computing and industry

      **Lower the barriers-to-entry** by keeping control of data in scientist hands

      Seamless integrations with scientific infrastructures and commercial entities

      Detailed monitoring capabilities and easy deployment have proven crucial

# Additional information

| | | |
|---|---|---|
| Website | | http://rucio.cern.ch |
| Documentation | | https://rucio.cern.ch/documentation |
| Repository | | https://github.com/rucio/ |
| Images | | https://hub.docker.com/r/rucio/ |
| Online support | | http://rucio.cern.ch/doc../join_rucio_mattermost/ |
| Developer contact | | rucio-dev@cern.ch |
| Journal article | | https://doi.org/10.1007/s41781-019-0026-3 |
| Twitter | | https://twitter.com/RucioData |

Hugo González Labrador
hugo.gonzalez.labrador@cern.ch