



Scientific data management with Rucio

Martin Barisits (CERN)

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

Principally developed by and for the ATLAS Experiment, now with many more communities

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

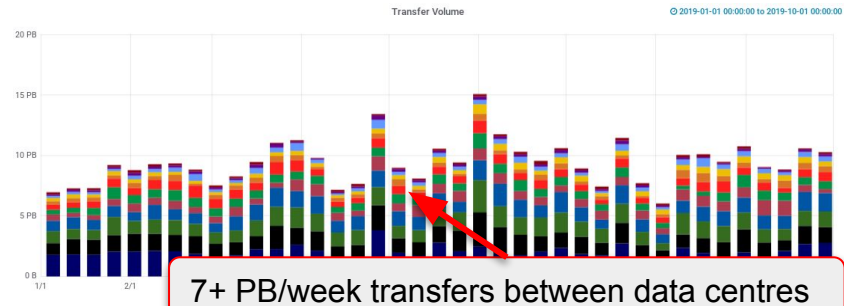
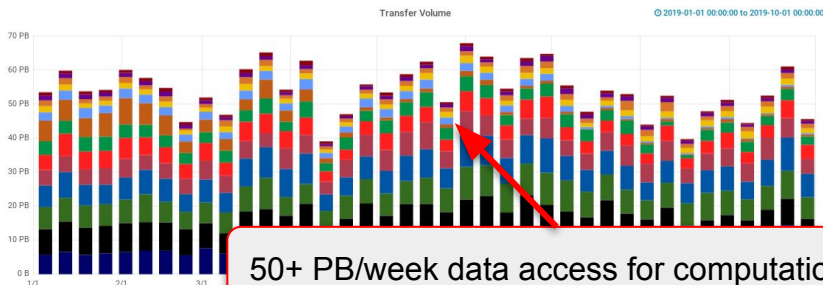
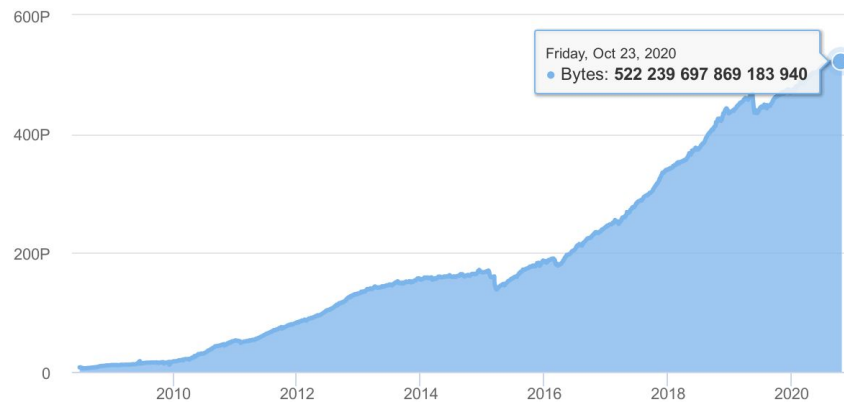
Open community-driven development process



A few numbers to set the scale

- 1B+ files, 500+ PB of data, 400+ Hz interaction
- 120 data centres, 5 HPCs, 2 clouds, 1000+ users
- 500 Petabytes/year transferred & deleted
- 2.5 Exabytes/year uploaded & downloaded

Increase 1+ order of magnitude for HL-LHC



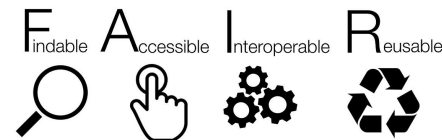
Rucio main functionalities



Provides many features that can be enabled selectively

More advanced features
↓

- **Horizontally scalable catalog** for files, collections, and metadata
- Transfers between facilities including **disk, tapes, clouds, HPCs**
- **Authentication and authorisation** for users and groups
- **Many interfaces** available, including CLI, web, FUSE, and REST API
- **Extensive monitoring** for all dataflows
- Expressive **policy engine** with rules, subscriptions, and quotas
- Automated **corruption identification and recovery**
- Transparent support for **multihop, caches, and CDN dataflows**
- **Data-analytics based flow control**



Rucio is not a distributed file system, it connects existing storage infrastructure over the network

No Rucio software needs to run at the data centres

Data centres are free to choose which storage system suits them best

Community



Advanced European Network of E-infrastructures
for Astronomy with the SKA



Science & Technology
Facilities Council



European Science Cluster of Astronomy &
Particle physics ESFRI research Infrastructures

