# EGI DataHub

## *Data as a Service – Distributed Data Management*

**Baptiste Grenier**

*EGI Foundation*

# Motivation
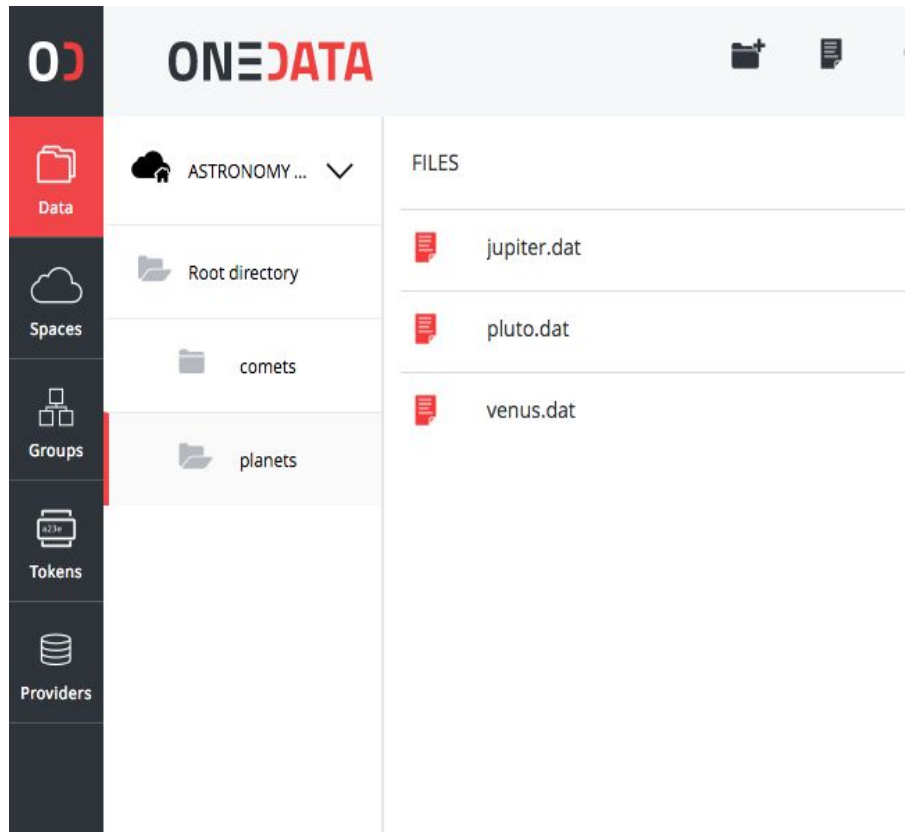
- Putting up a (scalable) distributed data infrastructure needs specific expertise, resources and knowledge

- No easy way to discover and transfer data

- No easy way of making data (publicly) accessible without transferring it to a sharing service

- No easy way of combining multiple datasets from different data providers

- Users need to access data locally and from compute resources

www.egi.eu  @EGI_eInfra

# EGI DataHub: components and concepts

- **EGI DataHub**: a **Onedata Onezone**, the **federation** and **authentication** service. SSO with all the connected storage providers (Oneprovider) through **EGI Check-in**

- **Oneprovider**: **data management** component deployed in the data centres, **provisioning data** and **managing transfers**. A default one is operated for EGI by CYFRONET.

- **Space**: a **virtual volume** where **users organize data**. A space is supported by one or multiple Oneproviders

- **Oneclient**: a client providing access to the spaces through a FUSE mount point (**local POSIX access**)

- **Web interfaces** and **APIs** are also available

# On the client side

*Web interface and Oneclient on the CLI*

# Replica management

*File distribution across providers*

# Metadata management

*Attaching metadata to files*

www.egi.eu   @EGI_eInfra

2/26/18

# File popularity and smart caching

- Transparent data access service

- Doing smart caching of remote storage

- Federating data sources/providers

- Publishing datasets

- Notebooks with DataHub

# DataHub for transparent data access



EGI DataHub for transparent data access

- Clients uses one or more providers to access data
- Data can be accessed over multiple protocols

# Smart caching of remote storage



- Site A hosts data and computing resources
- Site B hosts only data
- Site X uses data from A and B without pre-staging
- Pre-staging can also be done using APIs
- Data is accessed locally "à la" POSIX with FUSE

www.egi.eu | @EGI_eInfra

# Federation of service providers



Federation of service providers

- Heterogenous backend storage
- Common interfaces (Web, REST, POSIX, CDMI)
- Common AAI with Check-in
- Discovery of Datasets in the EGI DataHub

# Publishing and discovery of datasets



**Publication of datasets**

- PID minting
- Publishing, discovery and access to datasets

# Notebooks with DataHub



Noteboooks with DataHub, B2HANDLE and B2FIND

- Collecting and analysing dataset specificities
  - Number of files
  - Size of files

- Preparing a pilot
  - Designing and validating usage model
  - Integrating Onedata with existing resources

- Validating the pilot

- Deploying a production setup
  - Ensuring hardware requirements are sufficient
    - RAM, CPU, Disk, Network,...
    - Storage backend

- Preferred model: using docker containers
  - Using docker-compose
  - Packages for Ubuntu 16.04 and CentOS 7 also available

**The work of the EGI Foundation**
*is partly funded by the European Commission
under H2020 Framework Programme*

www.egi.eu | @EGI_eInfra

2/26/18 | 15

- **Powerful-enough Oneprovider**
  - RAM: 32GB
  - CPU: 8 vCPU
  - Disk: 50GB SSD
  - To be adjusted for the dataset and usage scenario

- **For high IOPS**
  - High-performance backend storage (CEPH)
  - Low latency network

- **POSIX mounting**
  - Oneprovider close to the Oneclient

# Links

- EGI DataHub
  - https://datahub.egi.eu/
  - https://community.egi.eu/c/egi-services/datahub
  - https://egi-datahub.readthedocs.io/
  - https://wiki.egi.eu/wiki/EGI_Federated_Data

- System requirements
  - https://onedata.org/docs/doc/system_requirements.html

- Official Onedata documentation
  - https://onedata.org
  - https://onedata.org/#/home/documentation
  - Getting started
    - https://github.com/onedata/getting-started
  - Source code: https://github.com/onedata

**Thank you
for your attention.**

*Questions?*

**www.egi.eu**
@EGI_eInfra