

EGI technical platforms for advanced computing

Tiziana Ferrari
Technical Director, EGI Foundation



www.egi.eu

EGI-Engage is co-funded by the Horizon 2020 Framework Programme
of the European Union under grant number 654142



- Introduction to EGI
- Services for distributed computing, data management and AAI
- New requirements, new challenges
- Towards the European Open Science Cloud

The EGI Services are provided by the EGI Federation

- **EGI Council participants:** national e-infrastructure providers and international research organisations (CERN and EMBL)
- **Integrated e-infrastructure providers**



826,000
Cores of compute capacity



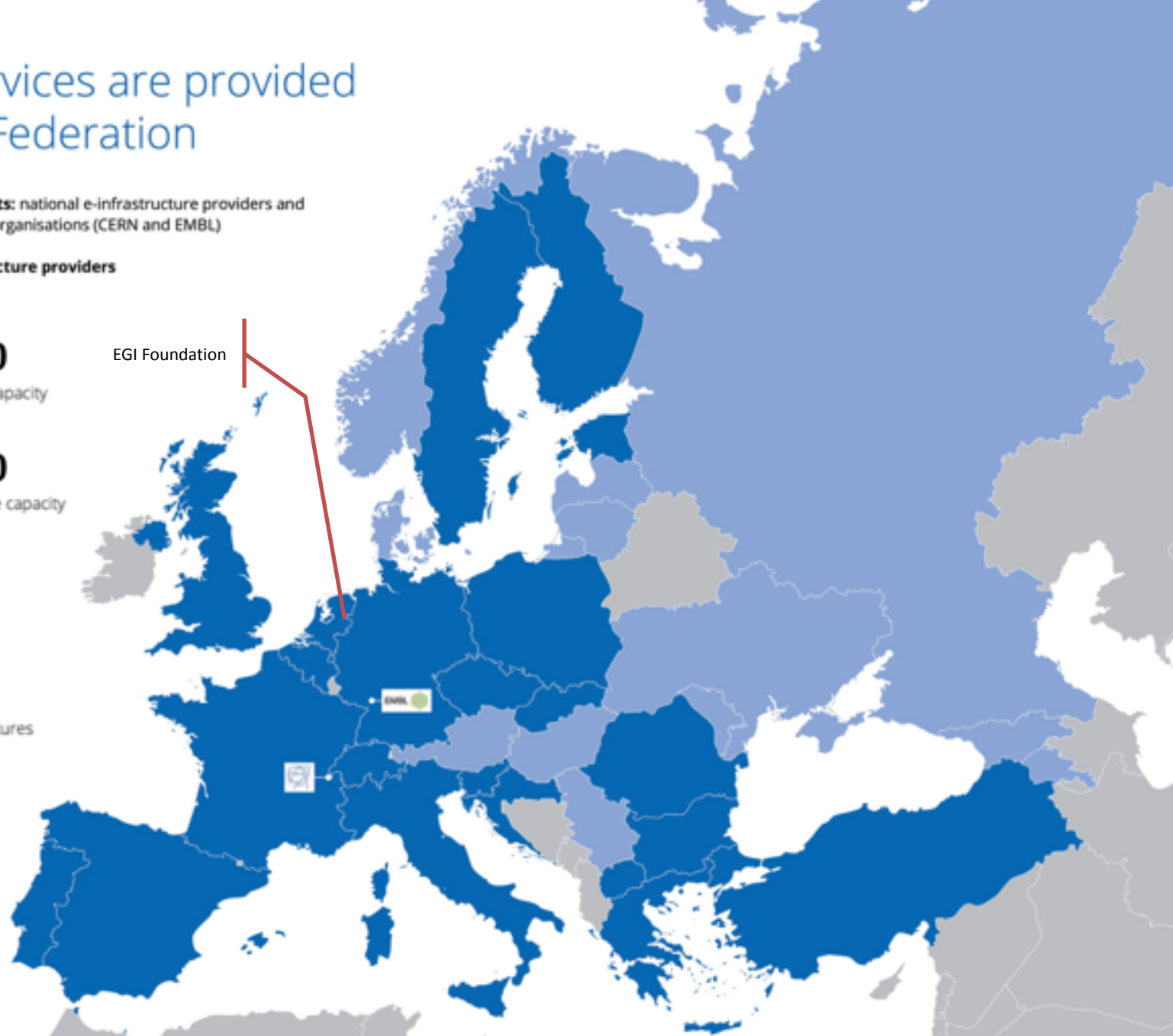
560,000
Terabytes of storage capacity



48,000
Users



15
Research Infrastructures
integrated with EGI



-
- A map of Europe with the European Union flag (a circle of twelve gold stars on a blue background) overlaid on various countries. The flag is shown in a larger size than in the previous slide, indicating a wider distribution or a different context.

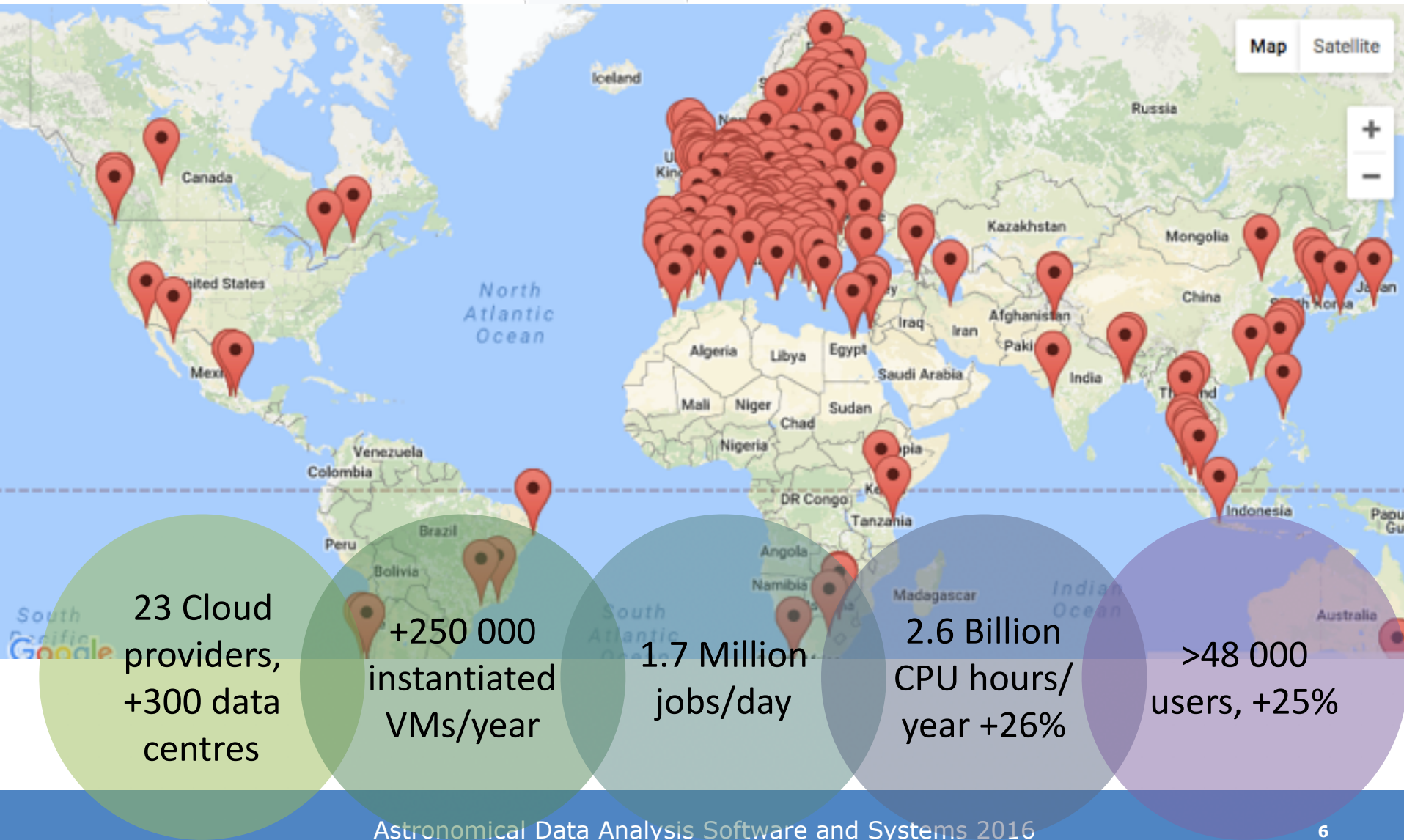


International Partnerships

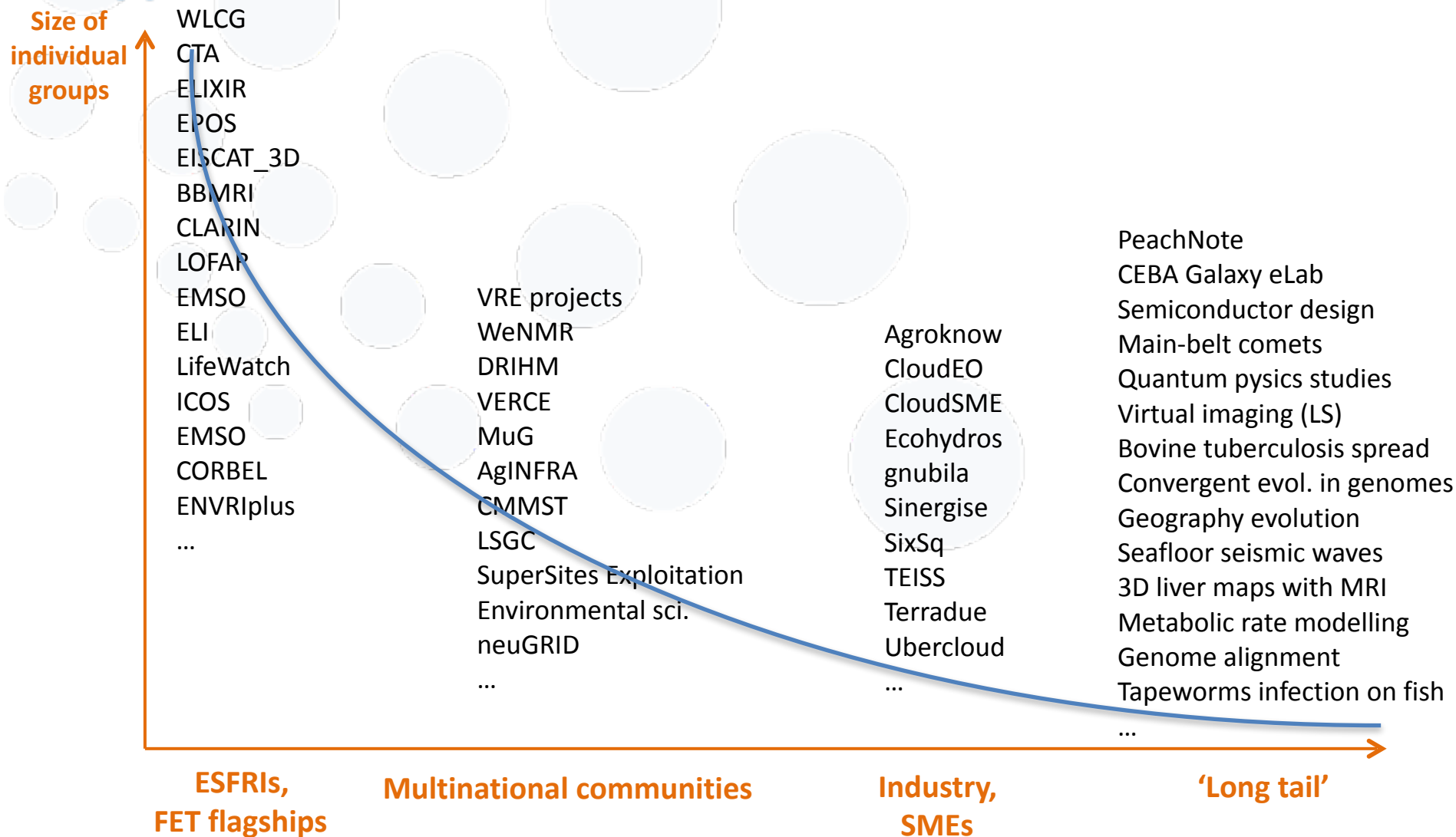


EGI Federation, 2016 QR3

The largest distributed compute e-Infra worldwide

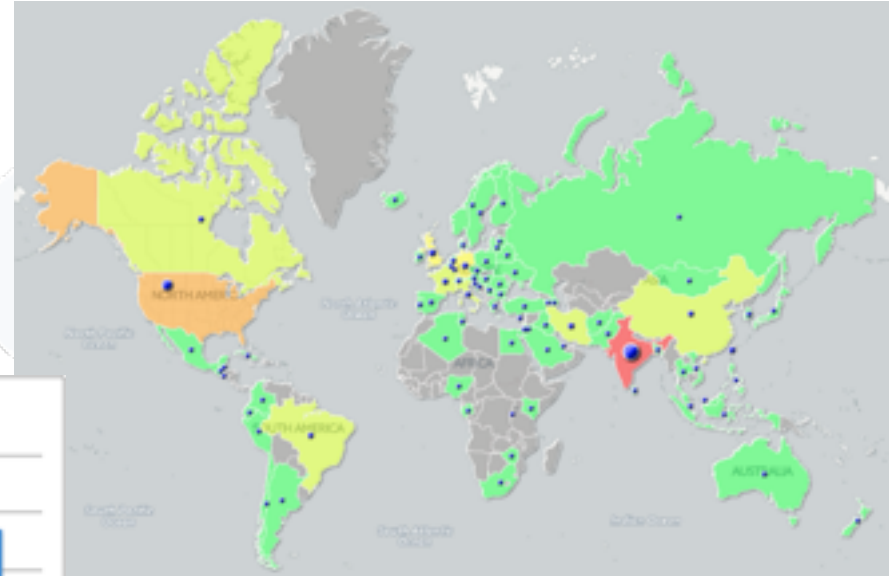
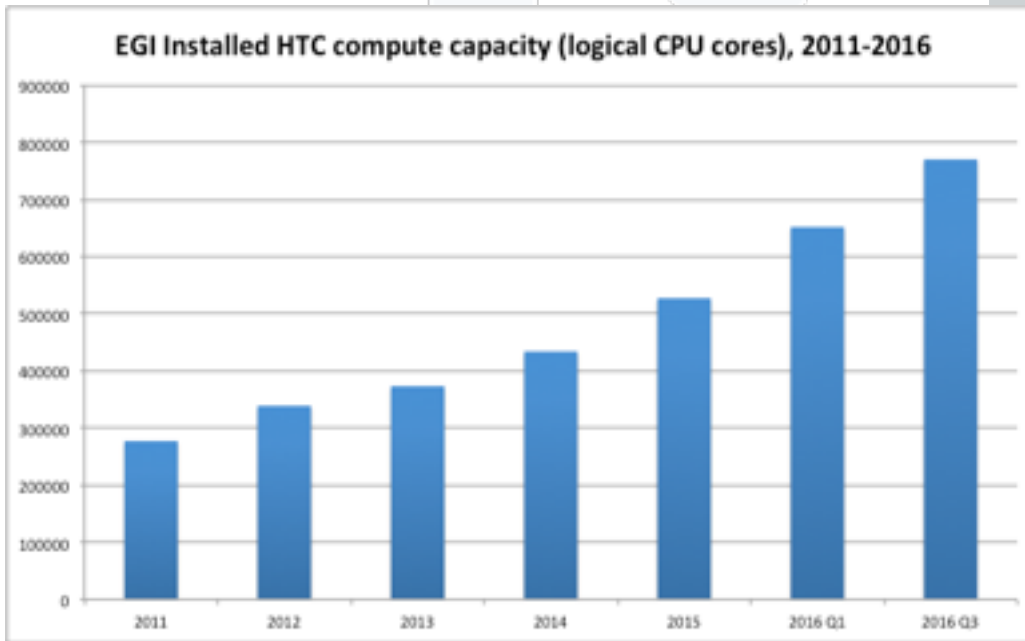


Serving researchers and innovators



Supporting international research communities and thematic services

Installed compute capacity trends 2011-2016



Example: Structural Biology
Distribution of users (2016, QR3)

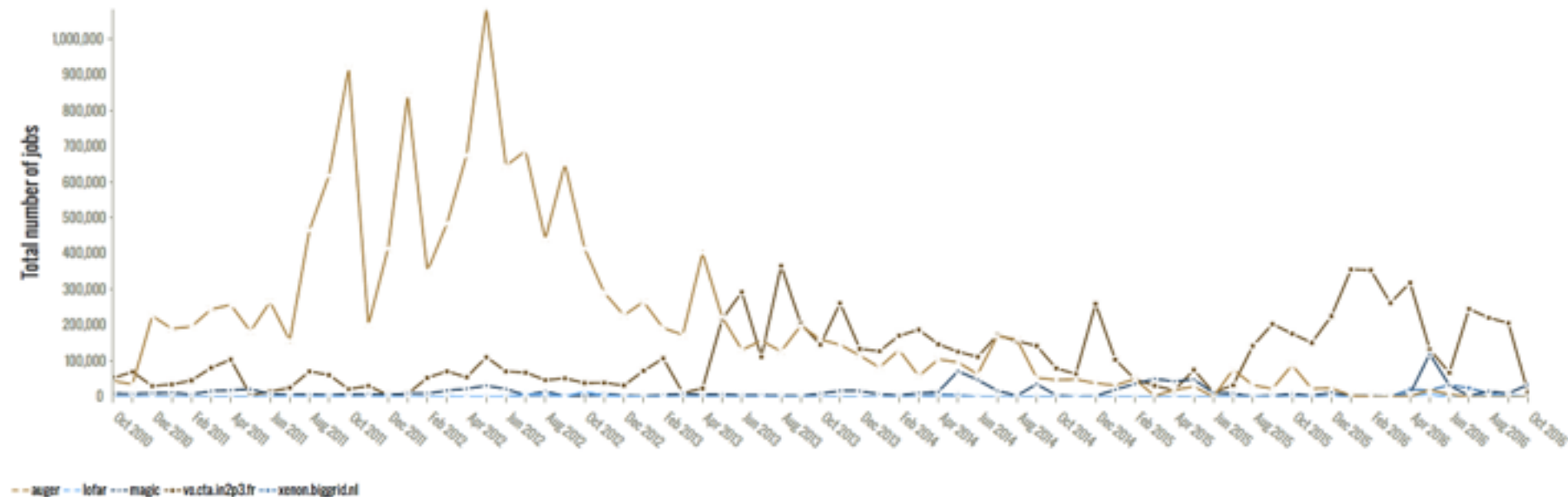
- 2700 users
- 81 countries

(credits: A. Bonvin, WeNMR)

Astronomy/Astrophysics/Astro-particle physics projects and RIs in EGI

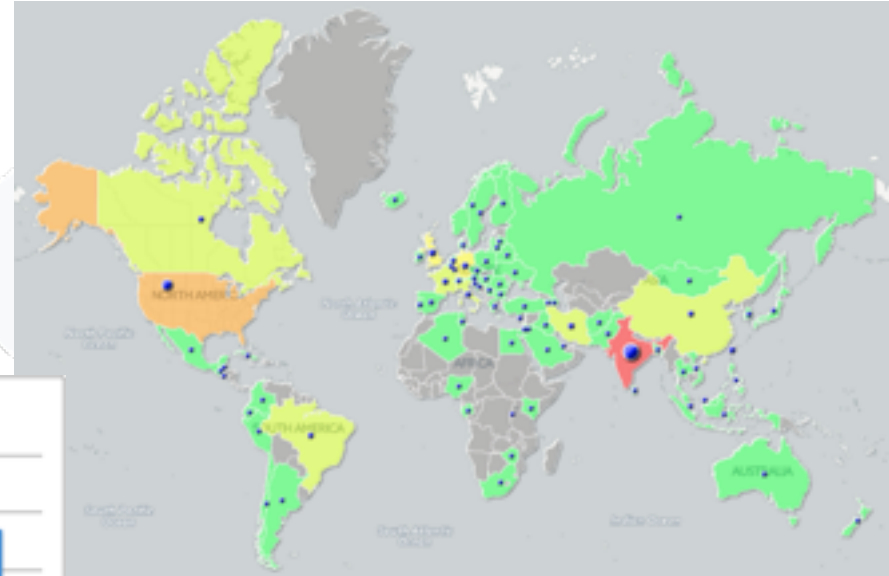
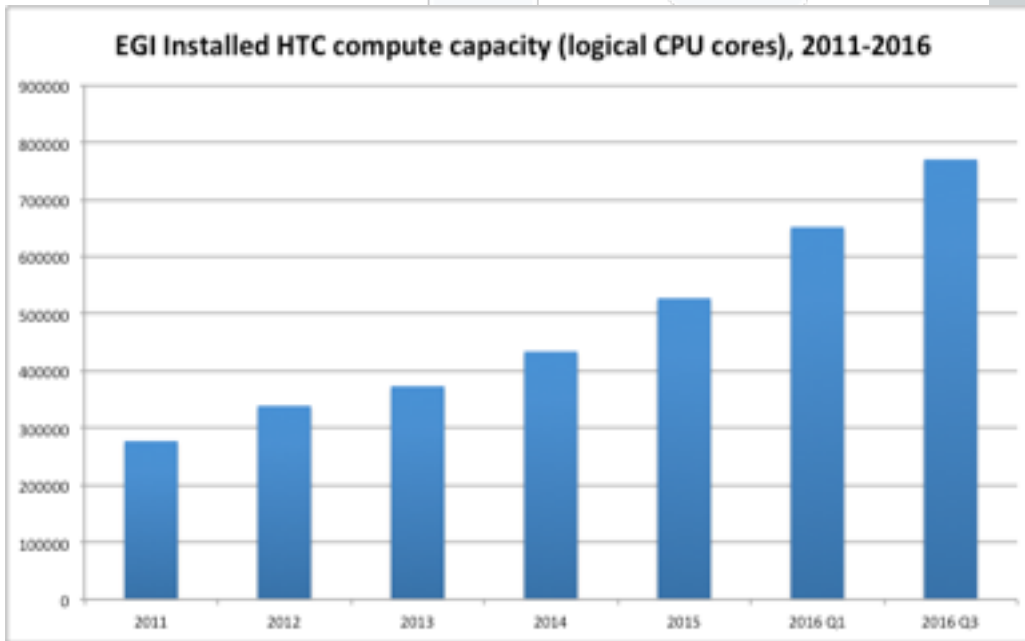
ARGO, AUGER, CTA, KM3NeT, LHCb, LOFAR, Large Synoptic Survey Telescope/LSST, PAMELA, ESA Planck Mission, XENON etc.

Total number of jobs by VO and Date



Supporting international research communities and thematic services

Installed compute capacity trends 2011-2016



Structural Biology
Distribution of users (2016, QR3)

- 2700 users
- 81 countries

(credits: A. Bonvin, WeNMR)

Services Catalogue

Compute



Cloud Compute >

Run virtual machines on demand with complete control over computing resources



Cloud Container Compute >

Run Docker containers in a lightweight virtualised environment



High-Throughput Compute >

Execute thousands of computational tasks to analyse large datasets

Training



FitSM training >

Learn how to manage IT services with a pragmatic and lightweight standard



Training infrastructure >

Dedicated computing and storage for training and education

Storage and Data



Online Storage >

Store, share and access your files and their metadata on a global scale



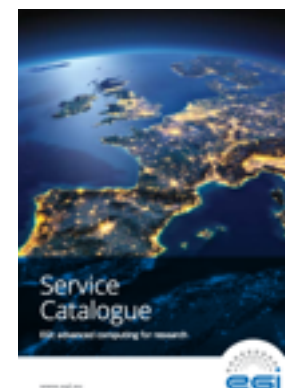
Archive Storage >

Back-up your data for the long term and future use in a secure environment



Data Transfer >

Transfer large sets of data from one place to another



<http://go.egi.eu/ServiceCatalogue>

Run virtual machines on-demand with complete control over the computing resources

- On-demand provisioning

Benefits

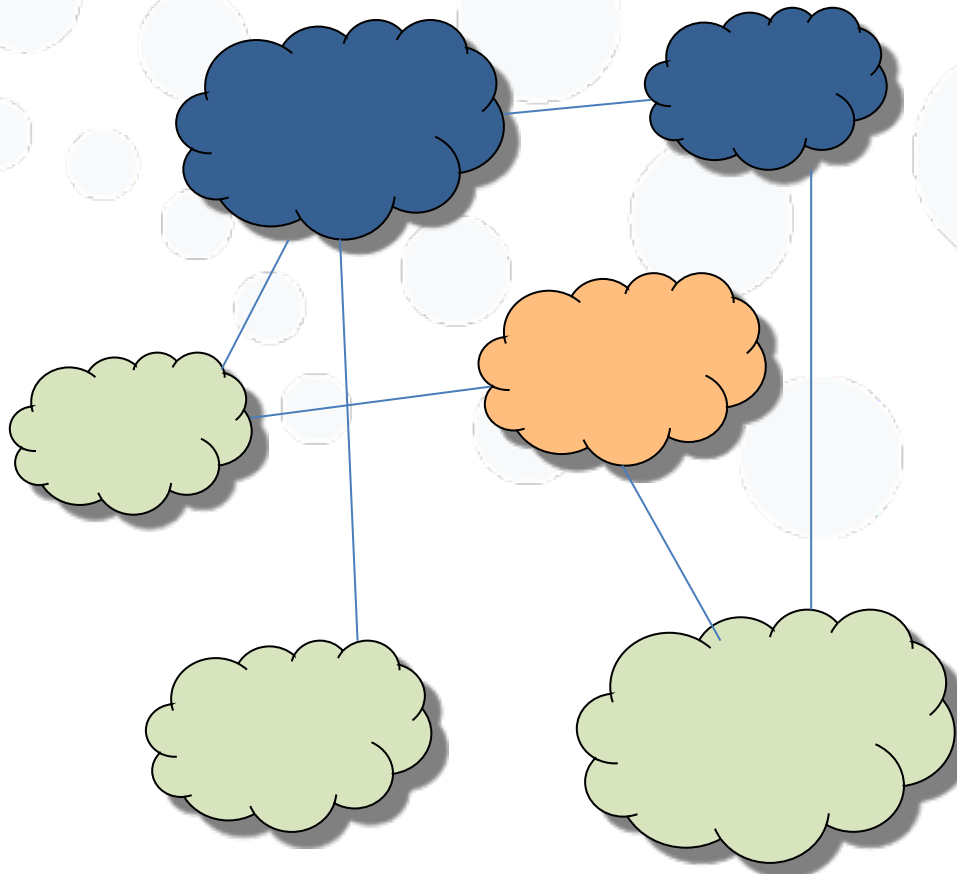
- Execute compute- and data-intensive workloads, including GPGPU computation



in the cloud

EGI Federated Cloud

- System of cloud infrastructures
- Standard user interfaces
 - Clouds and their interconnections are based on open standards, open technologies
 - Based on OCCI/OGF and OpenStack
- Harmonised operational behaviour
- Value proposition: distributed cloud computing for analysis of distributed large datasets



VM and block storage management:



Occi - On every site



OpenStack Nova - On OS sites

Object storage management (optional):

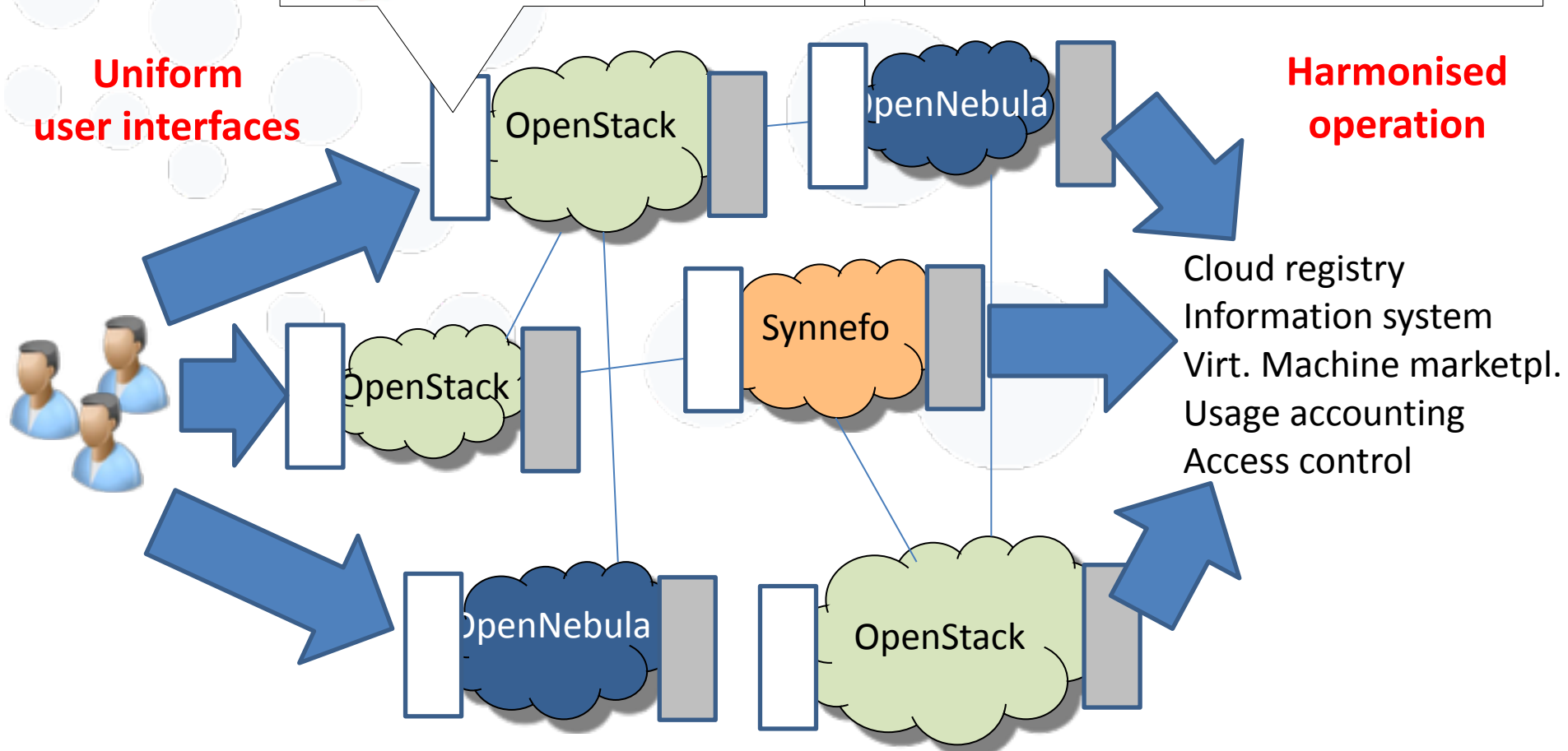


CDMI - on any site

OpenStack SWIFT – on OS sites

**Uniform
user interfaces**

**Harmonised
operation**



EGI Federated Cloud

EGI Federated Cloud is a collaboration of communities developing, innovating, operating and using cloud federations for research and education.

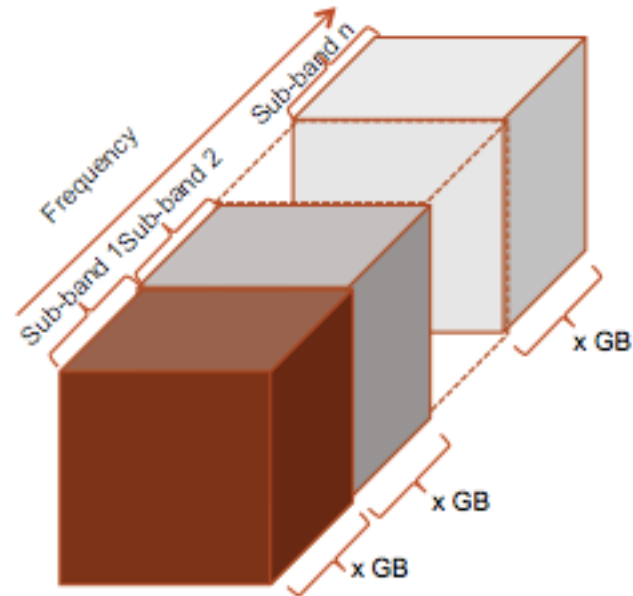
Today:

- 23 providers from 14 NGIs
 - 15 OpenStack
 - 6 OpenNebula
 - 1 Synnefo
- ~ 7.000 cores in total



LOFAR use case details

- Measurement sets: **datacubes** (3D data): two Fourier spatial coordinate axes plus a spectral axis.
- A datacube can reach several **terabytes**.
- LOFAR telescope allows **up to 488 subbands**, which can reach several GBs.
- Each subband processed independently.



Credits: Susana Sánchez Expósito - CSIC

Porting of LOFAR calibration pipeline

2/2

Implementing the Use Case: COMPSs application

- **COMPSs:**

- a data-driven programming model
 - it exploits the inherent parallelism of the applications
 - It executes the application tasks as soon as their input data are ready
- a VM orchestrator
 - It starts and contextualize the VM instances needed to execute the application tasks
 - It also checks the status, gathers the outputs and deletes the VM instances

- **Our COMPSs application:**

- A python script
- It iterates over the subbands executing for each one a COMPSs task
- They calls the LOFAR software (= executes a script) to process the subband.

```
import subprocess
import sys
import os
from pycompss.api.task import task
from pycompss.api.parameter import *

@task(script_name = FILE)
def iter_calib(script_name):
    os.chmod(script_name,0744)
    subprocess.call(script_name)
    print "end executiong"

if __name__ == "__main__":
    args = sys.argv[1:]
    DATA_PATH=args[0]
    TEMPLATE_FILE=args[1]
    f=open(TEMPLATE_FILE,'r')
    content=f.read()
    f.close()
    list_f=os.listdir(DATA_PATH)
    for directory in list_f: # Iterate over the data inputs
        if os.path.isdir(DATA_PATH+"/"+directory):
            new_content=content.replace("INPUTDATAPATH",directory)
            script_name="job"+directory+".sh"
            f=open(script_name,"w")
            f.write(new_content)
            f.close()
            iter_calib(script_name)
```

- The computing capabilities fulfil the requirements from the use case
 - The memory and cpu needs depends on the specific pipeline step, and the **EGI federated cloud allows to configure virtual machines with different capabilities.**
- A better storage solution is needed
 - **The user data are too large to be stored in the VM images.** They should be stored in volumes easily mountable from several VMs and synchronized across different cloud providers.
- COMPSs facilitates the porting and deployment of the application

Cloud Container Compute

Run Docker containers within isolated user-space with no overhead

- On-demand provisioning
- Lightweight environment for

Benefits

- Reduce time to production by removing friction between development and operations



High-Throughput Compute

Analyze large datasets by executing large numbers (thousands) of computational tasks

- Access to high-quality computing

Benefits

- Large amounts of processing capacity over long periods of time
- Faster results for your research



- ## Query result

Online storage

Store, share and access your files and their metadata on a global scale

- Assign global identifiers to files
- Access highly-scalable storage from

Benefits

- Highly scalable storage system accessible from anywhere
- Easily share data



Archive storage

Back-up your data for the long term and future use in a secure environment

- Store data for long-term retention
- Store large amount of data

Benefits

- Stores large amounts of data
- Long-term retention
- Reliable and interoperable



Data Transfer

Transfer large sets of data from one place to another

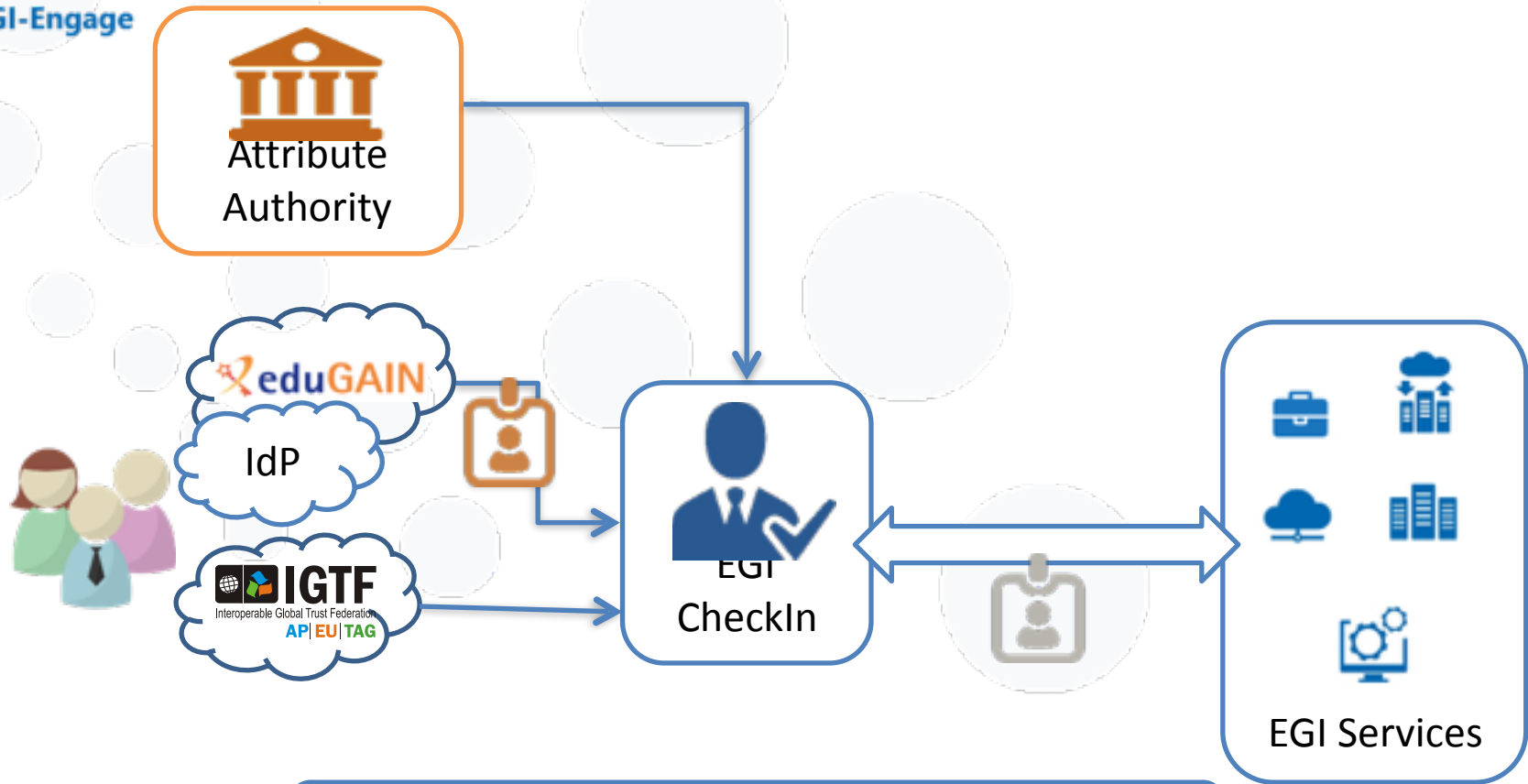
- Move research data fast
- Specialized analytics of on-going transfers

Benefits

- Ideal for very large files
- Able to handle large amounts of files
- Transfer process with automatic retry



EGI AAI New Architecture



Mandatory
Attributes

EGI UID

First name, last
name

email

affiliation

Why a IdP/SP Proxy?

- Service Providers (SPs) can have **one statically configured IdP**
- **No need to run an IdP Discovery Service** on each EGI SP
- Connected SPs get **harmonised user identifiers and accompanying attribute sets** from one or more AAs that can be interpreted in a uniform way for authZ purposes
- External IdPs only deal with a **single EGI SP** proxy

EGI services will not have to deal with the complexity of multiple IdPs/Federations/Attribute Authorities/technologies.

- Manage entire data life cycle from raw data to preservation
- Combine efficient computation services with open data managed by federated infrastructures
 - No local staging of data for processing
- Share public datasets for download or reuse
- Make public datasets discoverable

Open Data Platform interfaces

GUI

- Web based
- Easy data management and sharing, access control
- Publication of data

REST

- Advanced data and collection management API for integration with community tools and portals

CDMI

- Standard data management operations
- Advanced metadata queries
- Integration with

POSIX

- Enable direct mounting of spaces in the local filesystem without full data transfer

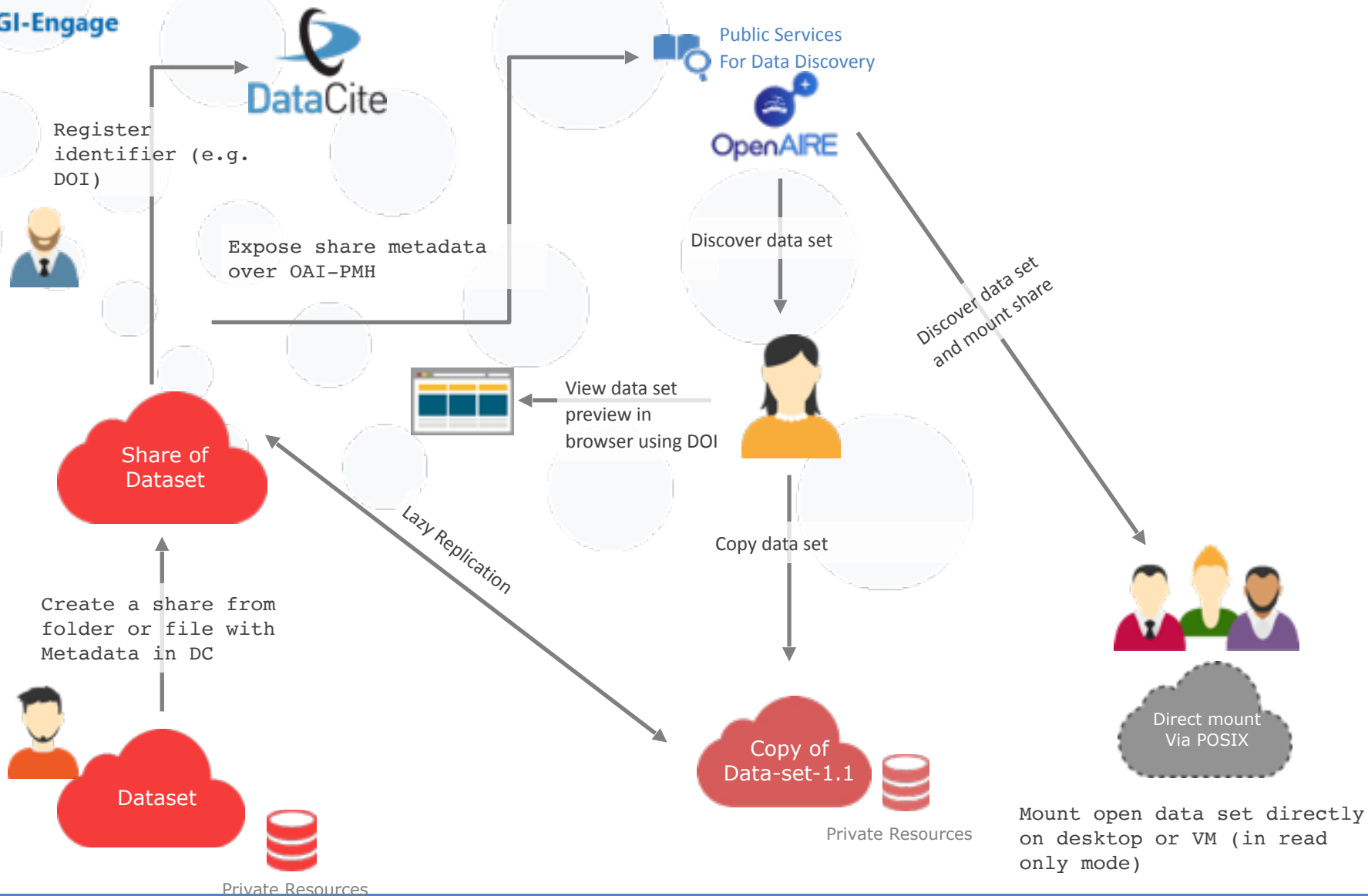
OAI-PMH

- OAI Data Provider interface
- Dublin Core metadata by default
- More complex metadata

HTTP

- Direct download of open data from URL's

Open Data Platform workflow



EGI role towards the European Open Science Cloud

- An **Open Science Service Exchange** as partnership of public/commercial organizations and initiatives responsible for
 - Provisioning of **wide set of services** to researcher and innovators
→ consolidation of national e-Infrastructures, open standards, technical and business process integration among the suppliers (e-Infrastructures, Research Infrastructures etc.)
 - **Platform integration** for community-specific capabilities with coordinated outreach
 - **Aggregation of demand** for economies of scale, technical requirements translations, cross-border access via **brokering and procurement, end-to-end operations**
 - Development of **human capacity**
 - A “Digital innovation hub” to support innovation with industry/SMEs

Thank you for your attention.

Questions?



Acknowledgements

This presentation used icons made by Freepik from www.flaticon.com

www.egi.eu

This work by Parties of the EGI-Engage Consortium is licensed under a [Creative Commons Attribution 4.0 International License](https://creativecommons.org/licenses/by/4.0/).

