



European Data Infrastructure - EUDAT

Data Services & Tools

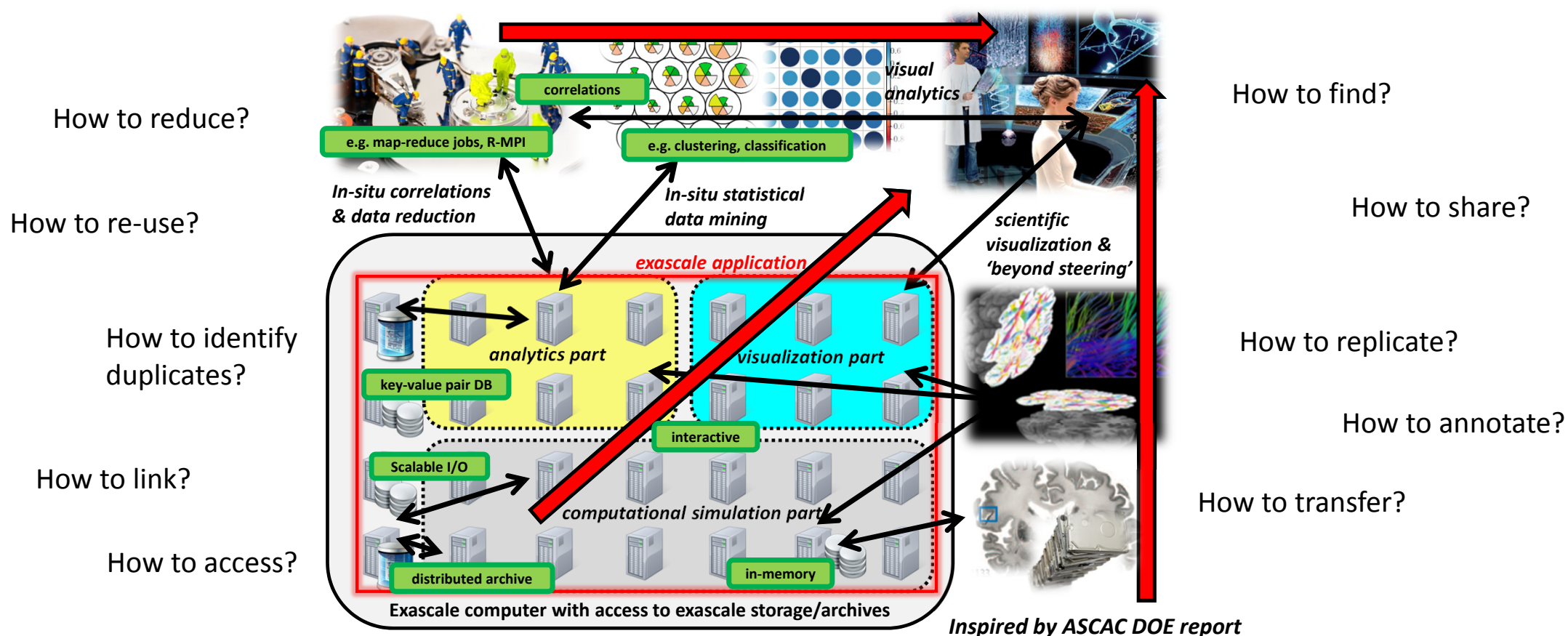


UNIVERSITY OF ICELAND
SCHOOL OF ENGINEERING AND NATURAL SCIENCES
FACULTY OF INDUSTRIAL ENGINEERING,
MECHANICAL ENGINEERING AND COMPUTER SCIENCE



Dr. – Ing. Morris Riedel
Research Group Leader, Juelich Supercomputing Centre
Adjunct Associated Professor, University of Iceland
BDEC2015, 2015-01-28

Relevance of Solving Big Data Challenges towards Exascale





Motivation - 'Need for Big Data Tools' in HPC & Exascale

Ever increasing volumes, varieties, velocities

- Shift from tape to active disks → **active processing**
- Data transfer-aware scheduling → **transfer takes time**
- Different copies of 'same data' → **sharing data necessary**
- Different copies of 'same data' in different representations → **delete some data** (e.g. tool-dependent data types, e.g. libsvm format vs. Original image, etc.)

Publication process changes

- Open referencable data is required for journals → **data publicly available**
- Long-lasting copies years after HPC users finished projects → **archiving**
- Technology changes, links need to persist in papers → **handle systems**

New toolsets

- Data replication, **in-memory & data sharing tools**, different filesystems, etc.
- Statistical data mining codes for classification, clustering, applied statistics, etc. (potential to validate, e.g. inverse problems, or reduce datasets, e.g. PCA)



Data Centers and Communities

26 European
partners



EBERHARD KARLS
UNIVERSITÄT
TÜBINGEN



KIT
Karlsruhe Institute of Technology



BSC
Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

| **epcc** |

INGV



CINES
Centre Informatique National
de l'Enseignement Supérieur

DKRZ
DEUTSCHES
KLIMARECHENZENTRUM



Red IRIS

**Max-Planck-Institut
für Meteorologie**

JÜLICH
FORSCHUNGSZENTRUM

rzg
RECHEN-
ZENTRUM
GARCHING

**Science & Technology
Facilities Council**

maatG



umweltbundesamt
ENVIRONMENT AGENCY AUSTRIA



SNIC

Trust-IT Trust-IT Services Ltd
Communicating ICT to markets

UNINETT
Jigma

SURF SARA

UCL

User Forums + 30 communities

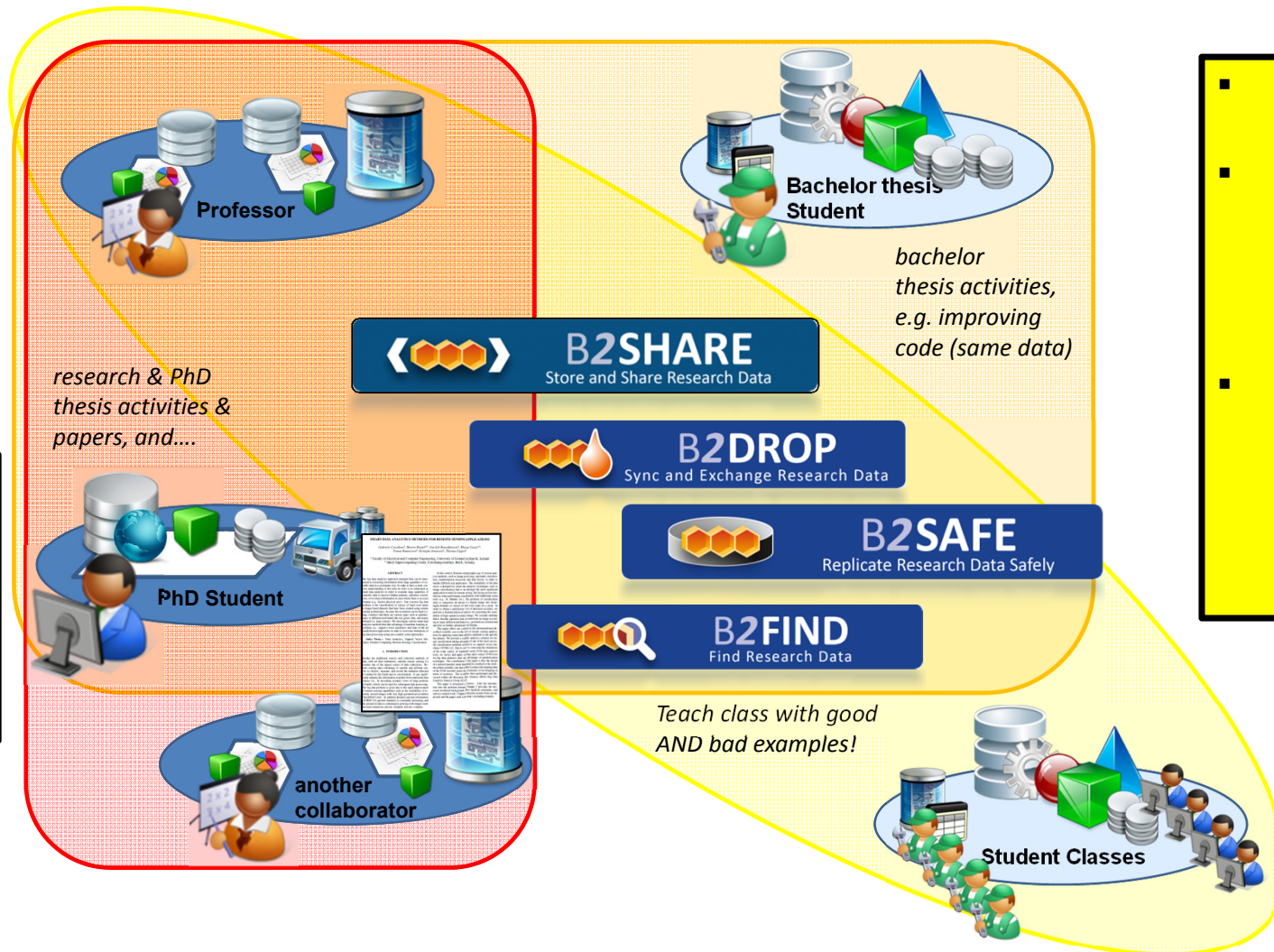


1st User Forum
7-8 March 2012, Barcelona



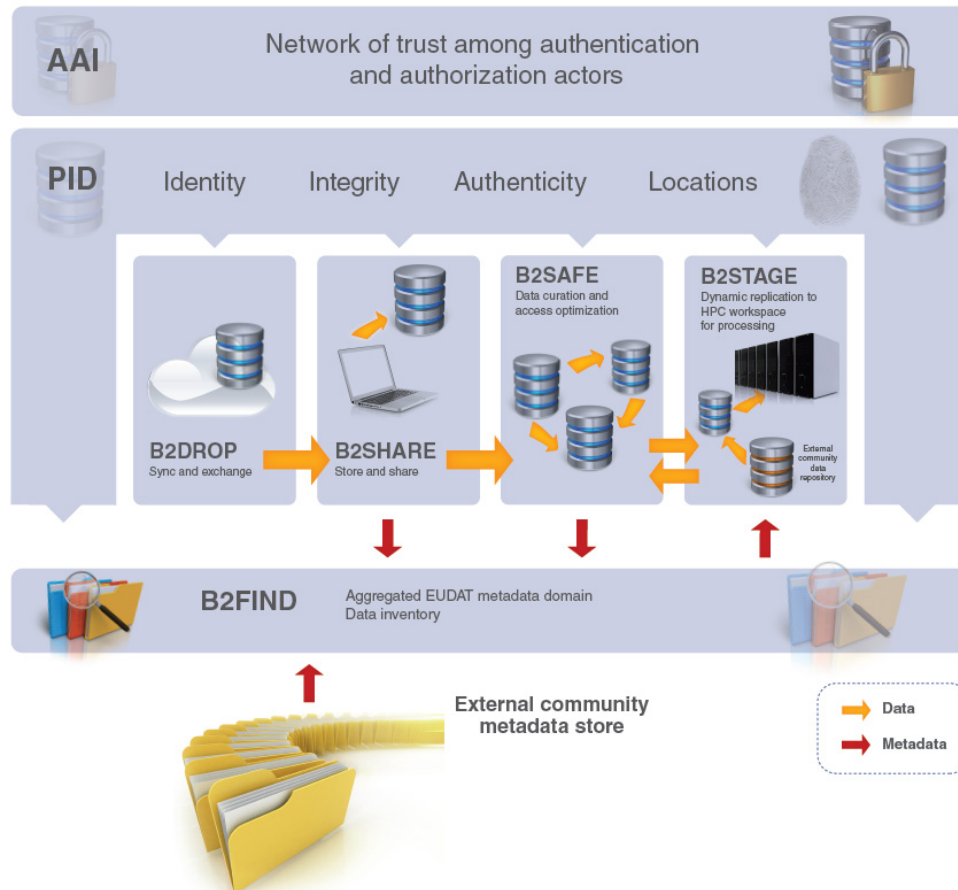
Work situation
in scientific
computing...

- Simple tools are important
- Avoid overheads in data management
- Realistic use within HPC environments



- Sharing different datasets is key
- One tend to loose the overview of which data is stored on which platform
- How do we gain trust to delete data when duplicates on different systems exist

Toolset Overview



- **Access and deposit,**
- **Informal data sharing**
- **Long-term archiving,**
- **Addressing identification, discoverability and computability**
- **long-tail and ‘big’ data**

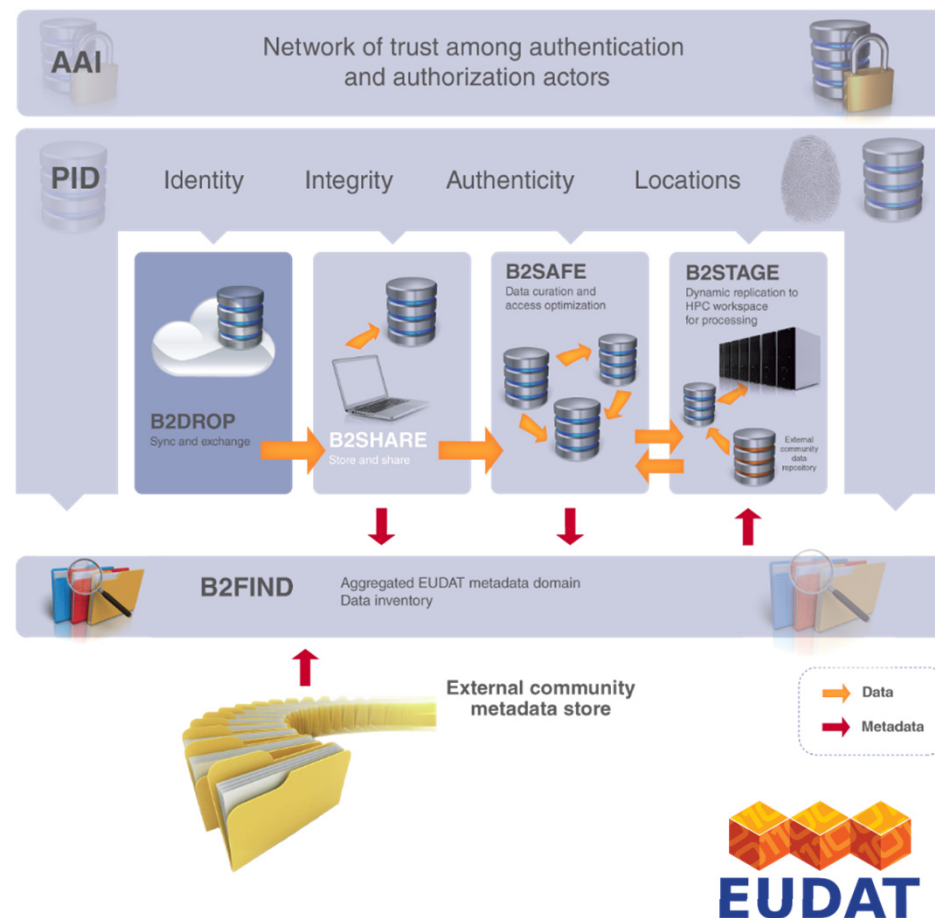
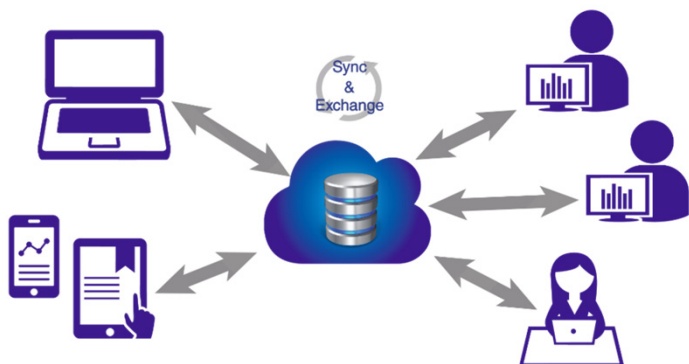
→ address full lifecycle of research data

→ adopt only what is needed

B2DROP is a **secure and trusted** data exchange service for researchers and scientists to keep their research data **synchronized** and up-to-date and to **exchange** with other researchers.

An ideal solution to:

- Store and exchange data with colleagues and team
- Synchronize multiple versions of data
- Ensure automatic desktop synchronization of large files



Features

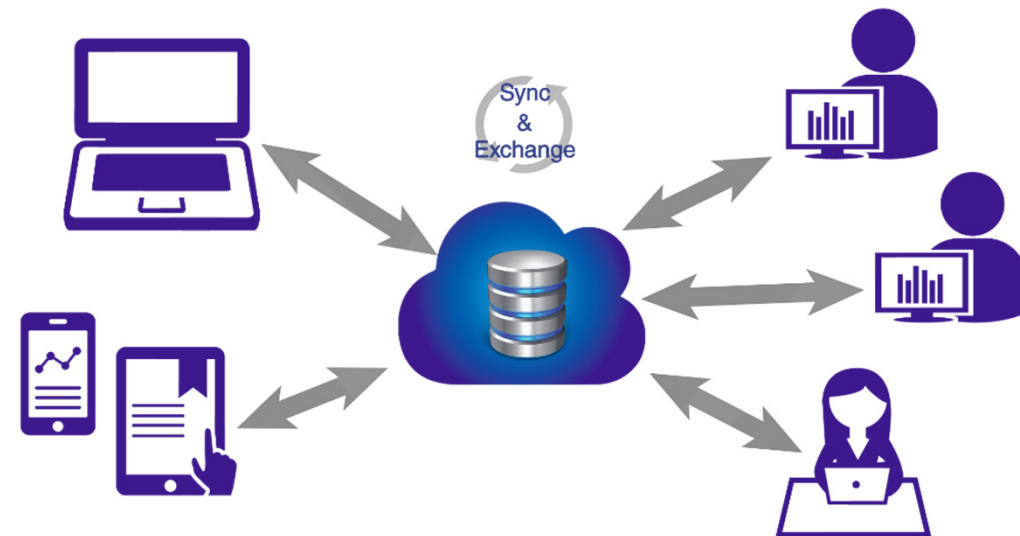
- future **integration with the B2 suite of services** to allow user-friendly data sharing
- users decide **with whom to exchange data, for how long and how**
- **up to 20GB of storage space** for research data
- **access and manage permissions** to files from any device and any location
- **simple to use and open to all** researchers, scientists, communities alike to **synchronize and exchange data with one or multiple users**



B2DROP



Sync and Exchange Research Data

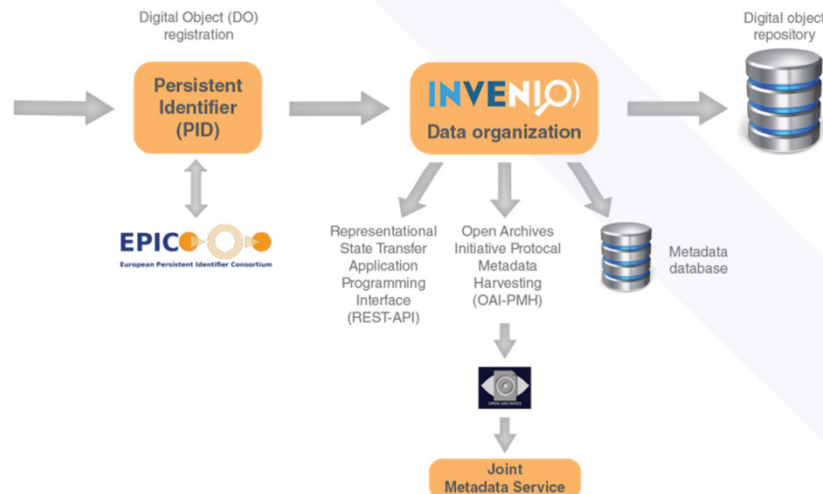
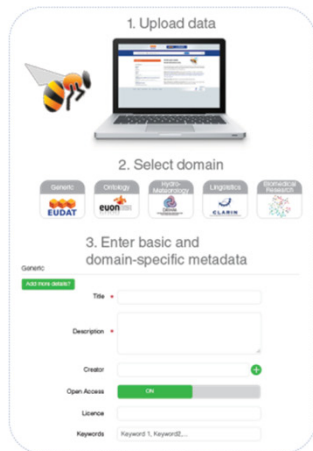
b2drop.eudat.eu



B2SHARE is a **user-friendly, reliable** and **trustworthy** way for researchers, scientific communities and citizen scientists to **store** and **share** small-scale research data from diverse contexts.

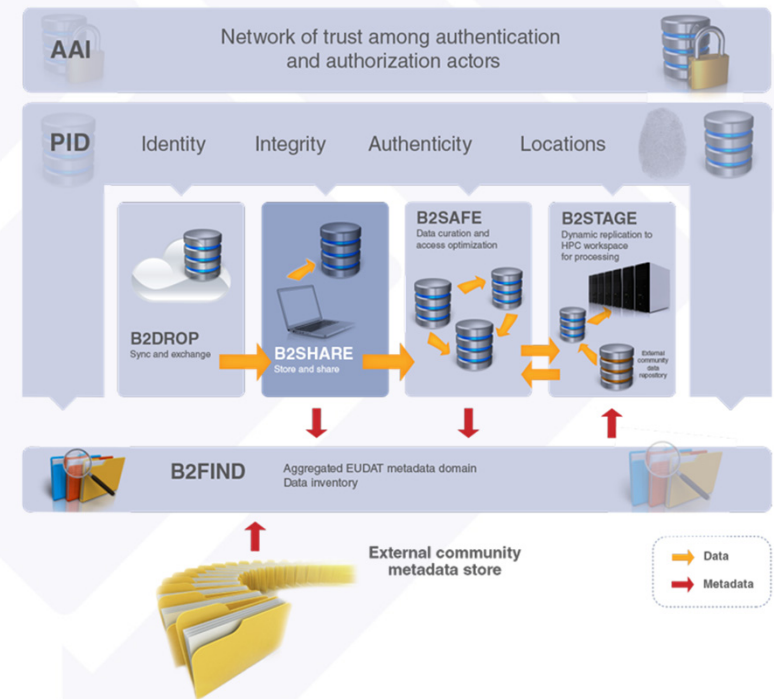
A winning solution to:

-  **Store:** facilitates research data storage
-  **Preserve:** guarantees long-term persistence of data
-  **Share:** allows data, results or ideas to be shared worldwide



B2SHARE
Store and Share Research Data

b2share.eudat.eu



EUDAT

Features



B2SHARE
Store and Share Research Data

b2share.eudat.eu

- Targets **small-scale research data collected as part of international collaboration and looking for a central repository**
- integrated with the **EUDAT collaborative data infrastructure**
- **free upload and registration** of stable research data
- data assigned a **permanent identifier**, which can be retraced to the data owner
- **community-specific metadata extensions and user interfaces**
- **openly accessible and harvestable metadata**
- **representational state transfer application programming interface (REST API)** for integration with community sites
- **data integrity ensured by checksum** during data ingest
- **professionally managed storage service** – no need to worry about hardware or network
- **monitoring of availability and use**

'A four-click service'



B2SHARE
Store and Share Research Data

b2share.eudat.eu



Search 139 records for

+ - ⚙️ 🔍 Search

You are logged in as Damien.

Deposit >>>

Latest Deposits

2014-10-23
PARADE, Strategy for a European Data Infrastructure White Paper

by Kimmo Koski[...] Strategy for a European Data Infrastructure White Paper by PARADE, Partn ...

2014-10-21
Knowledge Exchange Sustainability Index

Output from the Knowledge Exchange workshop: Sustainable Business Models ...

2014-10-10
MHD run for a single cluster simulated with 640^3 cells with ENZO-MHD at z=0.

by Franco Vazza
Hdf5 files for 3D dataset for a cluster run using a uniform grid box wit ...

2014-10-09
ENZO-AMR data for cluster E1

by Franco Vazza
Hdf5 monolithic 256^3 files for a simulated galaxy cluster at z=0, using ...

STORE AND SHARE YOUR RESEARCH DATA



A user-friendly, secure, robust, reliable and trusted service to share and store your research data *adding value to your research data by assigning Persistent Identifiers to ensure long-lasting access and reference.*

Deposit **and release** your data via the generic interface or select a community extension including specific metadata fields. **Releasing your data implies that your deposited data can be referred to, therefore any changes should be reflected in new data uploads.**

Share your data with others in a safe and trusted environment.

Do you belong to a scientific community? Brand and create your own community collection with specific metadata fields customized for your field.

Step 01

Drag and drop files here

Select files

Stop upload

Add basic details

Filename

EUDAT-DEL-WPS-D5 2 2-EUDAT Eas
Services.pdf

Generic

Add more details?

Title

Description

Creator

Open Access

Licence

Publisher

Publication Date

Tags

Linguistics

Language Code

Country/Region

Resource Type

Project Name

Quality

* indicates required field

Step 02

Select a domain

Generic



Ontology




Step 04

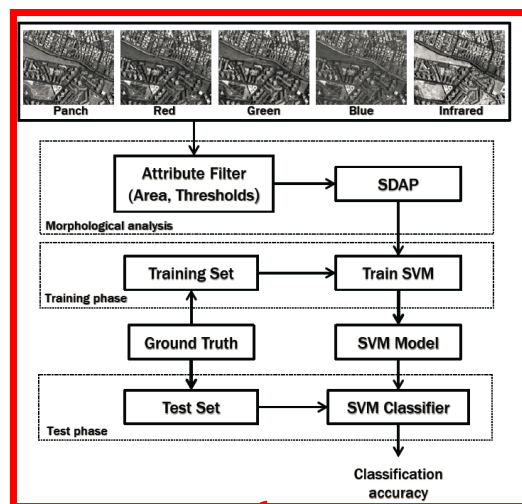
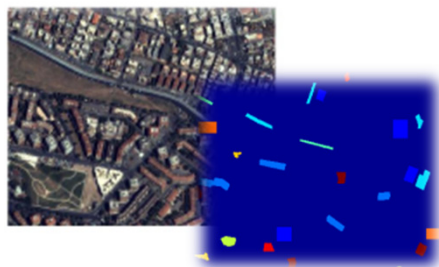
Deposit

HPC Usage Example



b2share.eudat.eu


Classification
(one field in data mining)



Class	Training	Test
Buildings	18126	163129
Blocks	10982	98834
Roads	16353	147176
Light Train	1606	14454
Vegetation	6962	62655
Trees	9088	81792
Bare Soil	8127	73144
Soil	1506	13551
Tower	4792	43124
Total	77542	697859

```
#!/bin/bash
#MSUB -N Train-tune-rec86-4-16-32
#MSUB -l nodes=4:ppn=16:performance
#MSUB -l walltime=03:00:00
#MSUB -M m.riedel@fz-juelich.de
#MSUB -m abe
#MSUB -W x=naccesspolicy:singlejob
#MSUB -v tpt=2
#MSUB -q devel

### jobscript

cd $PBS_O_WORKDIR
echo "workdir: $PBS_O_WORKDIR"

NSLOTS=32

echo "running on $NSLOTS cpus..."

### location
PISVM=/homeb/zam/mriedel/pisvm-1.2/pisvm-1.2/pisvm-train
TRAINDATA=/homeb/zam/mriedel/bigdata/86-romeok/sdap_area_all_training.el

### submit
mpirun -np $NSLOTS $PISVM -o 1024 -q 512 -c 10000 -g 16 -t 2 -m 1024 -s 0 $TRAINDATA
```

„Reference Data Analytics“
for reusability & learning

CRISP-DM
Report



Openly
Shared
Datasets



Running
Analytics
Code



-  **HPC JobScripts**
-  **HPC run in-/outputs**
-  **Input data & metadata**
-  **Output data & metadata**
-  **PIDs for Trust to Delete**
-  **Handle for Publications**

Classification
Study of
Land Cover
Types



Search for PiSVM Big Data Analytics in B2SHARE

Sattelite Data (Quickbird)

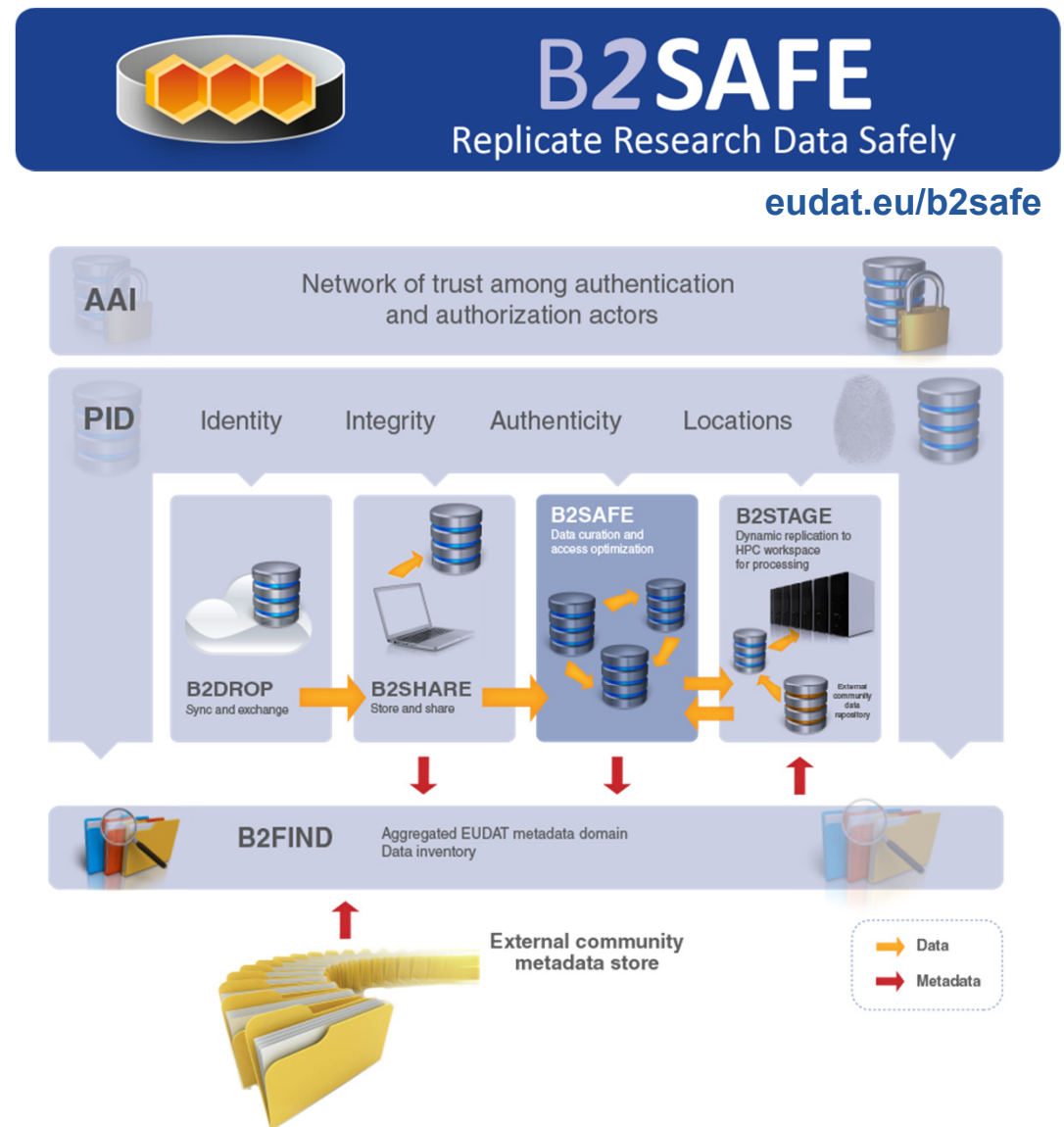
Parallel
Support Vector
Machines (SVM)

HPC / MPI code

B2SAFE is a **robust, safe and highly available service** which allows community and departmental repositories to **implement data management policies on their research data** across multiple administrative domains in a trustworthy manner.

A solution to:

- Provide an **abstraction layer which virtualizes large-scale data resources**
- Guard against data loss in long-term **archiving and preservation**
- Optimize access** for users from different regions
- Bring data **closer to powerful computers** for compute-intensive analysis



Features

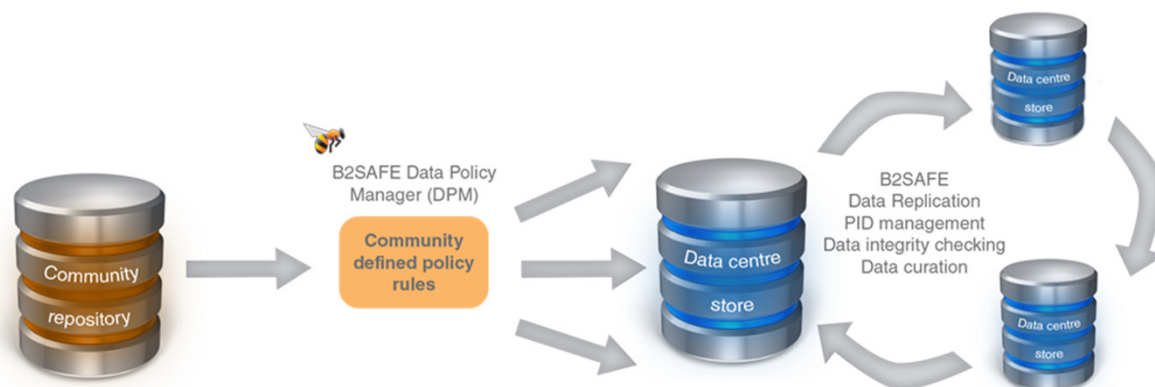


B2SAFE

Replicate Research Data Safely

eudat.eu/b2safe

- based on the execution of **auditable data policy rules** and the use of **persistent identifiers (PIDs)**
- respects the **rights of the data owners** to define the **access rights for their data** and to decide how and when it is made **publicly referenceable**
- data policies are **centrally managed via a Data Policy Manager**, and the policy rules are implemented and enforced by **site-local rule engines**
- able to **aggregate data from different disciplines** into a storage system of **trustworthy and capable data service providers**
- support for **repository packages** (e.g. DSPACE, FEDORA) and a **lightweight HTTP-based solution**



B2STAGE is a **reliable, efficient, light-weight and easy-to-use** service to **transfer research data sets** between EUDAT storage resources and high-performance computing (HPC) workspaces.

The service allows users to:

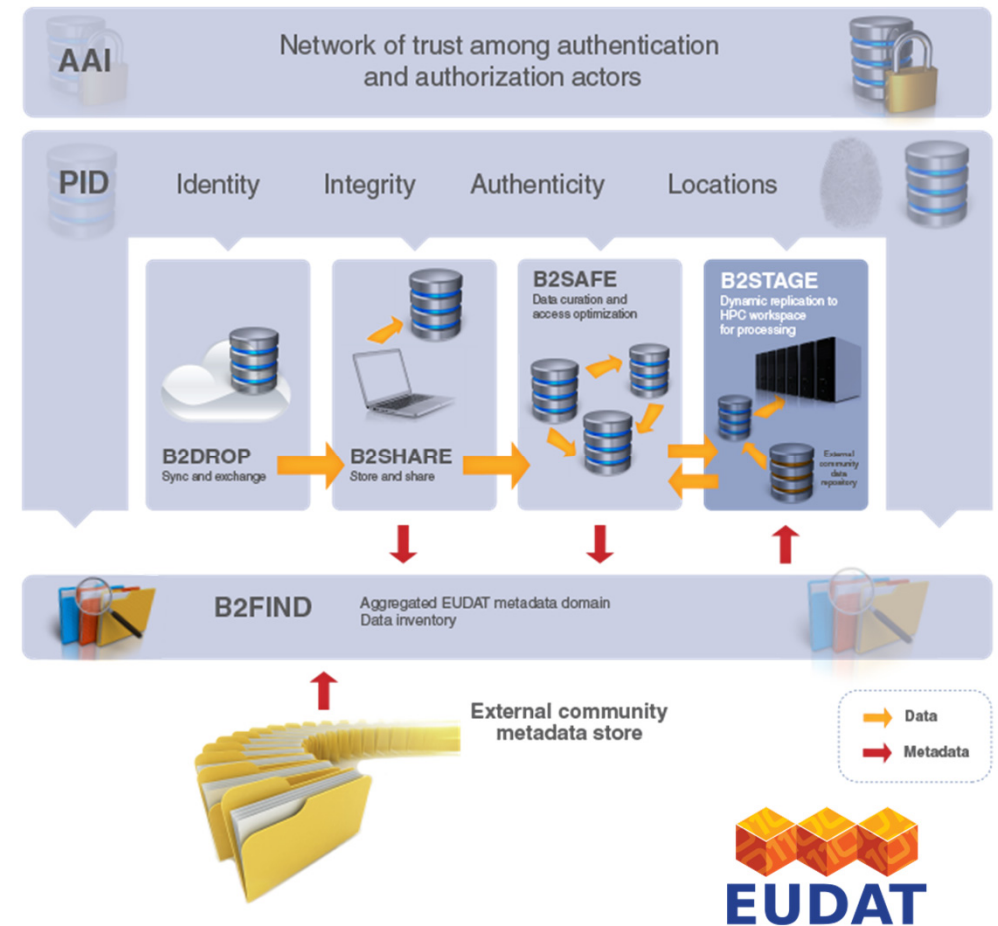
- Transfer large data collections from EUDAT storage facilities to **external HPC facilities for processing**
- In conjunction with B2SAFE, **replicate community data sets, ingesting them** onto EUDAT storage resources for long-term preservation
- Ingest computation results** into the EUDAT infrastructure
- Access data through a **RESTful HTTP interface** (in progress)



B2STAGE

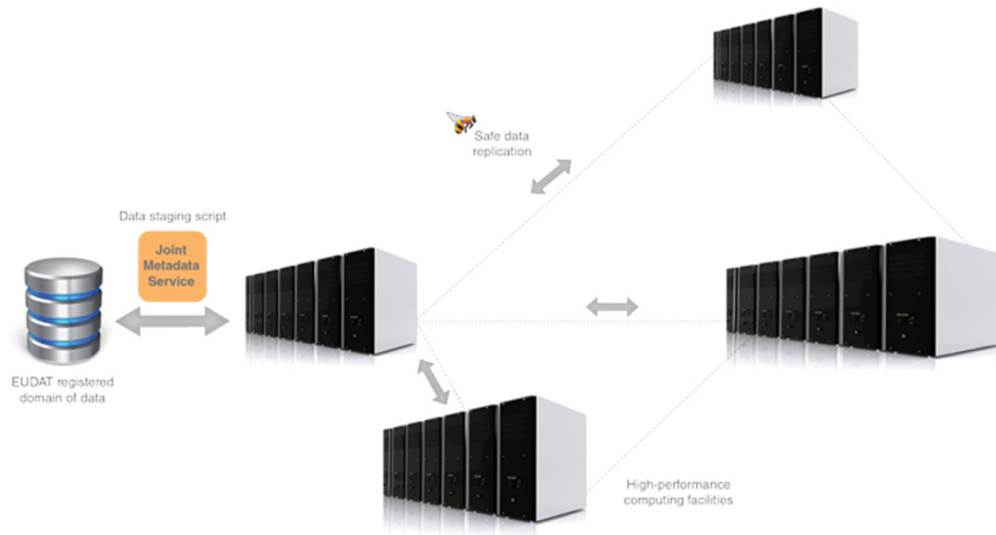
Get Data to Computation

eudat.eu/b2stage



Features

- **an extension of the B2SAFE and B2FIND services**, which allow users to store, preserve and find data
- **data-staging script** facilitates staging, ingestion and retrieval of persistent identifier (PID) information of transferred data
- service available **to all registered researchers and interested communities**






B2STAGE
Get Data to Computation

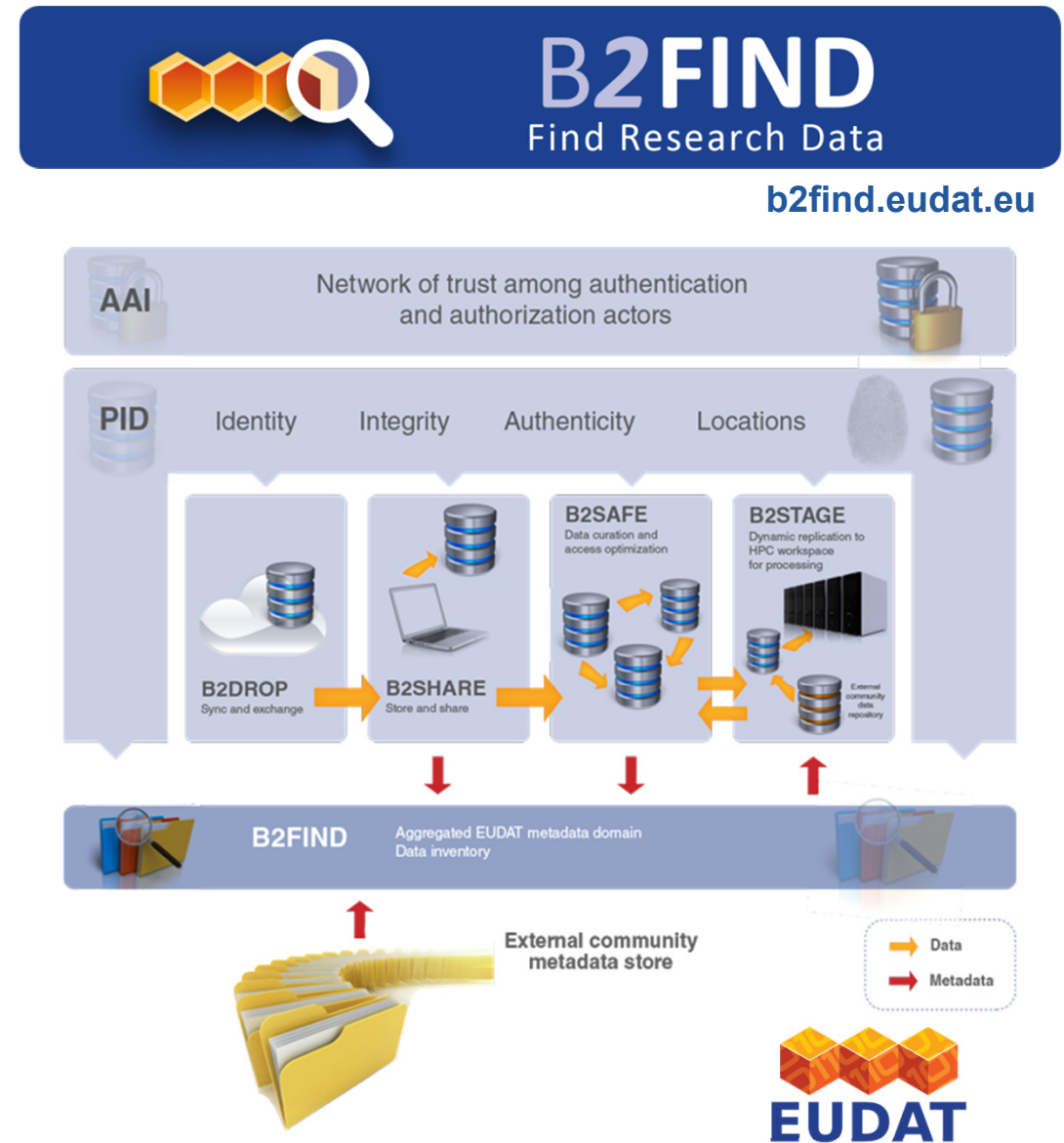
eudat.eu/b2stage

- **users negotiate access to remote HPC services in parallel**
- **collaboration with other infrastructures**, such as the European Grid Infrastructure (EGI) and Partnership for Advanced Computing in Europe (PRACE)
- **documentation, educational material and service helpdesk** available to support users

B2FIND is a simple, user-friendly **metadata catalogue of research data collections** stored in EUDAT data centres and other repositories.

A service which allows users to:

-  **Find** collections of scientific data quickly and easily, irrespective of their origin, discipline or community
-  Get quick **overviews** of available data
-  Browse through collections using **standardized facets**



Features

- supports faceted, geospatial and temporal metadata searches
- allows users to search and browse datasets via keyword searches
- initially available for communities in the EUDAT registered domain of data
- EUDAT will then extend the service to other interested and reliable data and metadata providers
- results displayed in user-friendly format and listed in order of relevance
- access to the scientific data objects is given through references provided in the metadata

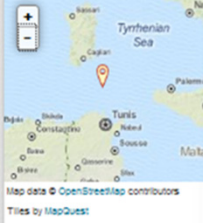


B2FIND
Find Research Data

b2find.eudat.eu

/ Datasets / Foraminifera abundance of sediment core MD81-LC07

Dataset extent



Dataset | Related

Foraminifera abundance of sediment core MD81-LC07

Hayes, Angela (2008): Foraminifera abundance of sediment core MD81-LC07, doi:10.1594/PANGAEA.407622

Data and Resources

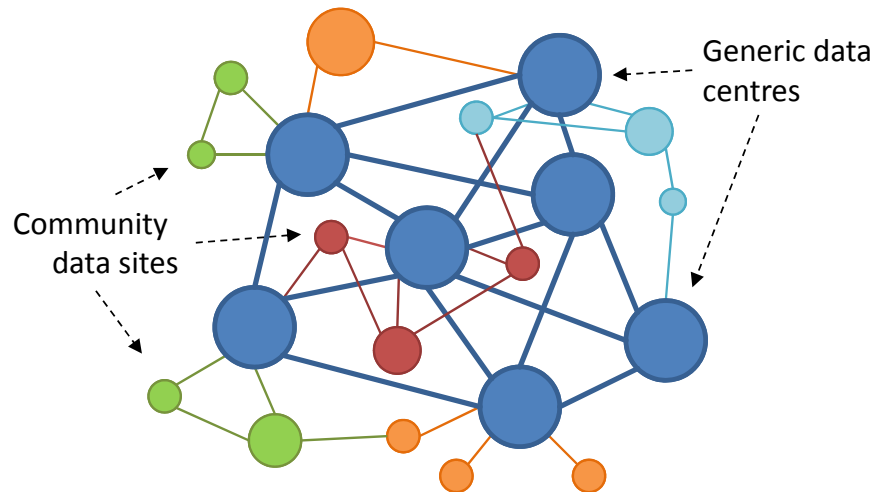
This dataset has no data

Globigerina bulloides

Additional Info

Field	Value
Source	http://dx.doi.org/10.1594/PANGAEA.407622
Author	Hayes, Angela (2008)
Version	68c7252d8134ebb8a19ce88a84e98ade
Discipline	Species
GeographicCoverage	Informationnotprovided
MetadataAccess	http://metadata.gbif.org/catalogue/OAIHandler?verb=GetRecord&metadataPrefix=eml&identifier=oai:metadata.gbif.org:eml/portal/2290.xml
Origin	PANGAEA - Publishing Network for Geoscientific and Environmental Data
PublicationYear	2008
TemporalCoverage:BeginDate	1995-01-01
TemporalCoverage:EndDate	1995-01-01

A Federated and Distributed CDI

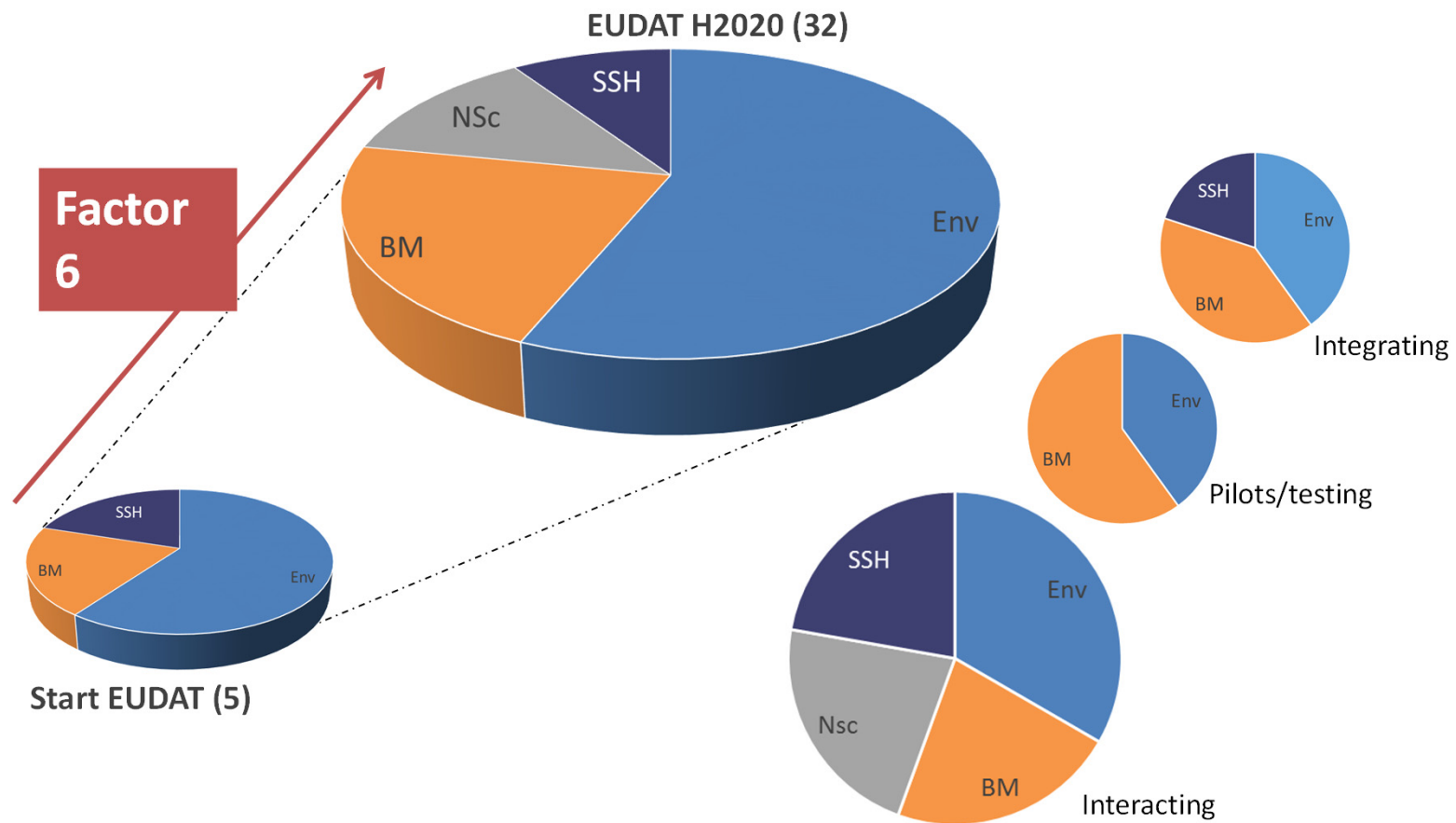


- **Using EUDAT services:** finding and accessing data, for instance, or storing smaller data sets by interacting with one of the CDI public front-end services

vs

- **Joining the CDI:** implies a tighter integration with at least one of the EUDAT centre → partnership between legal entities relying on OLAs and SLAs

Community Outreach & Service take Up



EUDAT Policies / Data Access and Reuse

- Open Access?
 - Funders: “Yes, absolutely!”
 - Researchers: “Yes, but...”
 - Some data is “sensitive”
 - What about credit and merit – others ‘harvesting’?
 - How to find one’s way in the legal minefield?
 - ‘Data-driven application-enabling’ activities
 - Providing tools and services to handle sensitive data
 - Licensing guidance, PIDs and usage statistics
 - Training & working on case studies (e.g. HPC simulation data demands)

OPEN DATA – WHAT DO EUDAT COMMUNITIES REALLY THINK ABOUT IT?

MARIE SANDBERG, PAWEŁ KAMOŃSKI, DAMIEN LECARPENTIER, ROB BAXTER

TO APPEAR IN ERCIM NEWS No. 100 (JANUARY 2015) [1]

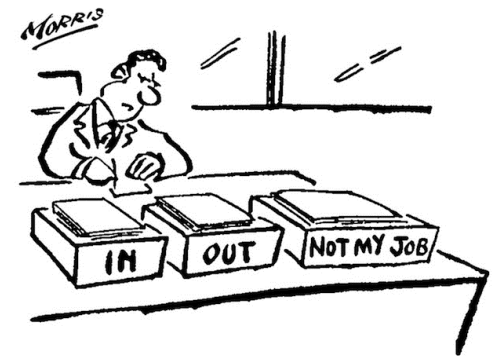
Facilitating open access to research data is a principle endorsed by an increasing number of countries and international organizations, and one of the priorities in the European Commission’s Horizon 2020 funding framework [2][3]. But what do researchers themselves think about it? How do they perceive the increasing demand for open access and what are they doing about it? What problems do they face, and what sort of help are they looking for?

As a pan-European research data infrastructure, these are questions that are of fundamental interest to EUDAT. To understand better what researchers think, EUDAT has conducted a programme of interviews with fourteen major research communities from the fields of life sciences, Earth and atmospheric science, astrophysics, climate science, biodiversity, agricultural science, social science and humanities – a broad cross-section of European research interests. While one cannot, of course, interpret the views of any given individual as the official position of a whole research community, they nevertheless provide useful information on the general attitude, requirements and challenges researchers face with regard to opening up their research data. In this article we report on our initial conclusions from this survey.

GROWING AWARENESS

Open access to research data is increasingly seen as a compelling principle in many research communities. There is a growing awareness of the global move towards open access, the potential benefits of it, and the necessity to implement open access policies within their disciplines. According to preliminary figures on the first wave of open data pilot projects in Horizon 2020, the opt-out rate among proposals submitted to the “open by default” categories was below 30%, and the opt-in rate among other proposals was around about the same. This underlines our findings in EUDAT – researchers are pretty happy about sharing their data.

CHALLENGES AHEAD



Need to Understand Computational Scientists

- Research Infrastructures → CDI users, partners & stakeholders
 - Uptake plans: work with computational scientists & HPC users to understand where data services make a difference
- It is not only about developing technical solutions, but also about defining the right partnership model
 - Take into account existing arrangements within pan-European research communities (organisational structure, funding schemes, business models, etc.)





E-Infrastructure Commons

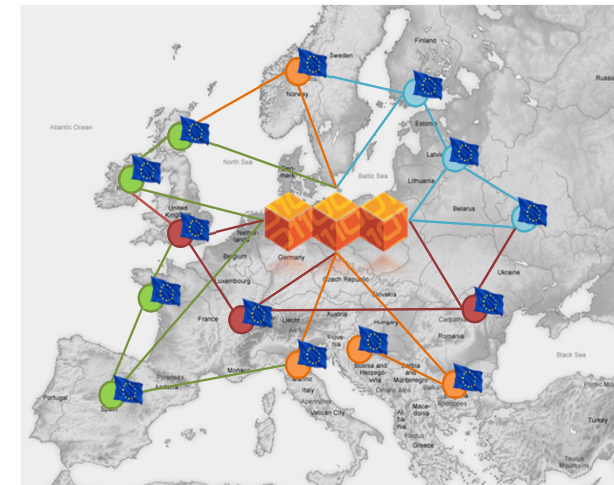
- Users have a “right” to a seamless access to network, data, and computing resources funded by public money
 - It is our role to make it as easy as possible for users → Users should not care about which e-Infrastructures they are using
- Cross-Infrastructure services
 - Based on pilots with interested communities
- E-Infrastructure Commons Roadmap





Bridging National and European Data Solutions

- Making national resources more available
 - Making visible valuable national collections through EUDAT
 - Access to European resources through national catalogues
- Enhancing cross-national collaborations
 - EUDAT provides a European extension to national solutions
 - True data sharing & archiving are pan-European challenges





Thank you

Talk available on
<http://www.morrisriedel.de/talks>