# EUDAT

## Towards a pan-European Collaborative Data Infrastructure

Giuseppe Fiameni (g.fiameni@cineca.it)  – Claudio Cacciari
*SuperComputing, Application and Innovation – CINECA*

Johannes Reatz
*RZG, Germany*

Damien Lecarpentier
CSC-IT Center for Science, Finland

Napoli, 16 May 2012

# EUDAT Key facts

| Project Name | EUDAT – European Data |
|---|---|
| Start date | 1st October 2011 |
| Duration | 36 months |
| Budget | 16,3 M€ (including 9,3 M€ from the EC) |
| EC call | Call 9 (INFRA-2011-1.2.2): Data infrastructure for e-Science (11.2010) |
| Participants | 25 partners from 13 countries (national data centers, technology providers, research communities, and funding agencies) |
| Objectives | "To deliver cost-efficient and high quality Collaborative Data Infrastructure (CDI) with the capacity and capability for meeting researchers' needs in a flexible and sustainable way, across geographical and disciplinary boundaries." |

# The current data infrastructure landscape: challenges and opportunities

- Long history of data management in Europe: several existing data infrastructures dealing with established and growing user communities (e.g., ESO, ESA, EBI, CERN)

- New Research Infrastructures are emerging and are also trying to build data infrastructure solutions to meet their needs (CLARIN, EPOS, ELIXIR, ESS, etc.)
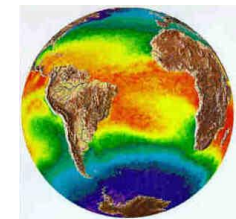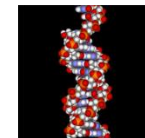
➢ **However, most of these infrastructures and initiatives address primarily the needs of a specific discipline and user community**

## Challenges

- Compatibility, interoperability, and cross-disciplinary research
    ➢ how to re-use and recombine data in new scientific contexts (i.e. across disciplinary domains)

- Data growth in volume and complexity (the so-called "data tsunami")
    ➢ strong impact on costs threatening the sustainability of the infrastructure
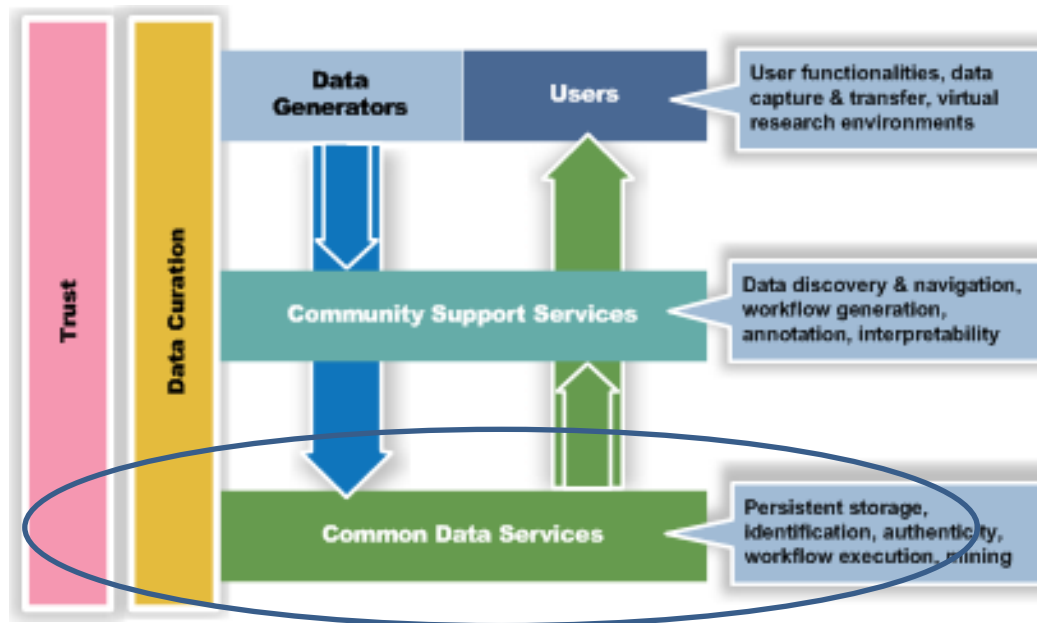
## Opportunities

- Potential synergies do exist: although disciplines have different ambitions, they have common basic needs and service requirements that can be matched with generic pan-European services supporting multiple communities, thus ensuring at the same time greater interoperability.

➢ **Strategy needed at pan-European level**

EUDAT

# The CDI concept

# EUDAT Core Service Areas

**Community-oriented services**

- Simple Data Acces and upload
- Long term preservation
- Shared workspaces
- Execution and workflow (data mining, etc.)
- Joint metadata and data visibility

**Enabling services (making use of existing services where possible**

- Persistent identifier service (EPIC, DataCite)
- Federated AAI service
- Network Services
- Monitoring and accounting

**Core services are building blocks of EUDAT's Common Data Infrastructure**
mainly included on bottom layer of data services

# Data centers and Communities

# First EUDAT Communities

# Building the services

6 service/use cases identified

Safe replication: Allow communities to safely replicate data to selected data centers for storage and do this in a robust, reliable and highly available way.

Dynamic replication: Perform (HPC) computations on the replicated data. Move (part of) the safely replicated data to a workspace close to powerful machines and move the results back into the archives.

Metadata: A joint metadata domain for all data that is stored by EUDAT data centers by harvesting metadata records for all data objects from the communities.

Simple store : A function that will help researchers mediated by the participating communities to upload and store data which is not part of the officially handled data sets of the community.

PID: a robust, highly available and effective PID system that can be used within the communities and by EUDAT.

AAI: A solution for a working AAI system in a federation scenario.

# SAFE_REPLICATION@EUDAT

**Safe Replication**

**Objective**: Allow communities to replicate data to selected data centers for storage and do this in a robust, reliable and highly available manner.

**Description** The ability to safely and simply replicate data from one data center to another is essential to EUDAT's task of improving data curation and accessibility.

Several EUDAT user communities (CLARIN, ENES, EPOS, and VPH) have identified safe replication as a common need, and are working to design a blueprint for managing data replication based on users' requirements and constraints

Data replication solutions and services are embedded into critical security policies, including firewall setups and user accounting procedures.

**More info:** eudat-safereplication@postit.csc.fi

# DATA_STAGING@EUDAT

**Data Staging**

**Objective**: Allow communities to stage data between EUDAT resources and HPC/HTC resources for computational purposes.

**Description***:* This service will allow the communities to dynamically replicate a subset of their data stored in EUDAT to an HPC machine workspace in order to be processed.



**More info: eudat-datastaging@postit.csc.fi**

# METADATA@EUDAT

**Metadata**

**Objective**: Create a joint metadata domain for all data stored by EUDAT data centers and a catalogue which exposes the data stored within EUDAT, allowing data searches.

**Description:** The EUDAT repository should provide an inventory of metadata from different communities

**More info: eudat-metadata@postit.csc.fi**

# SIMPLE_STORE@EUDAT

## Simple Store

**Objective**: Create an easy to use service that will help researchers mediated by the participating communities to upload and store data which is not part of the officially handled data sets of the community.

**Description:** This service will address the long tail of "small" data and the researchers/citizen scientists creating/manipulating them and NOT the short head of big data.

**More info: eudat-simplestore@postit.csc.fi**

# PIDS@EUDAT

**Persistent Identifiers**

**Objective:** Deploy a robust, highly available and effective PID service that can be used within the communities and by EUDAT.

**Description:** Keeping track of the "names" of data sets or other digital artefacts deposited with the CDI requires more robust mechanisms than "noting down the filename". The PID service will be required by many other CDI services, from Data Movement to Search and Query.

Currently considering use of both EPIC for data objects, and DataCite to register DOIs (Digital Object Identifiers for published collections.

**More info: eudat-persistentidentifiers@postit.csc.fi**

# AAI@EUDAT

**AAI – Distributed Authentication**

**Objective:** Provide a solution for a working AAI system in a federated scenario.

**Description:** Design the AA infrastructure to be used during the EUDAT project and beyond.

**Key tasks**:
Leveraging existing identification systems within communities and/or data centers
Establishing a network of trust among the AA actors:
Identifty Providers (IdPs), Service Providers (SPs), Attribute Authorities and Federations
Attribute harmonization

**More info: eudat-AAI@postit.csc.fi**

# Background

# AA process'actors

1. Federations
2. Multiple IdPs (e.g. home institute IdP)
   – Provision for supporting "homeless" users, cf SWITCH
   – Attributes from home institute
   – Technology – IdPs should use the same technology
3. Attribute authorities
   – Attributes relating to collaborations/communities (e.g. roles, memberships)
   – Each community should be prepared to manage and publish the user attributes
4. Multiple service providers
   – All consuming the *same* identities and attributes
   – Single Sign on: single IdP

consolidation (conversion) of credentials and attributes

any acceptable Identity Provider (AuthN)

IdP A

IdP B

IdP C

IdP #

shib

x.509

eID

COMMUNITIES

AtP A

AtP B

AtP #

community managed Attribute Provider (AuthZ)

multiple consolidation servers for load balancing and failover.

EUDAT

consolidated

Community and EUDAT services are using consolidated credentials and attributes

# Assumptions, Statements

- The IdP is an issuer of any kind of acceptable identity credential (x.509, shibboleth/saml2, card based eID, OpenID, credentials from social networks)

- Communities are assumed to manage their AtP (but they can offer a IdP services too if needed)

- AtPs can make use of the consolidated identity credentials to map their attributes (roles) to identities (green arrows indicate the usage of consolidated credentials).

- The credential/attribute conversion service is a gateway to EUDAT services which must be high-available. Therefore this service should be distributed over more than one server (load balancing, failover).

- The conversion service must be safe and trustworthy. Domains of trust can be fragmented (although they are encourages to collaborate). As a possible solution, specific centres could offer their conversion service for „their" affiliated communities (and service providers).
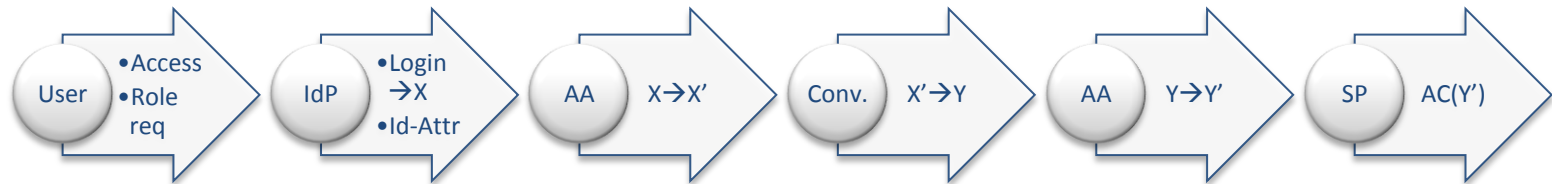
- The AtP of the „community" of (homeless) citizen scientists can be managed by EUDAT.

- Citizen scientists should be able to use any acceptable means for AuthN (including the eID on their national ID card)

- The EUDAT services need to build a trusted connection only to these credential consolidation gateway. No need to maintain large distributions of (e.g. IGTF) CA certificates etc. at the SP side.

EUDAT

# AA process: general overview

Some steps are of course optional



User • Access • Role req → IdP • Login →X • Id-Attr → AA X→X' → Conv. X'→Y → AA Y→Y' → SP AC(Y')

Technology
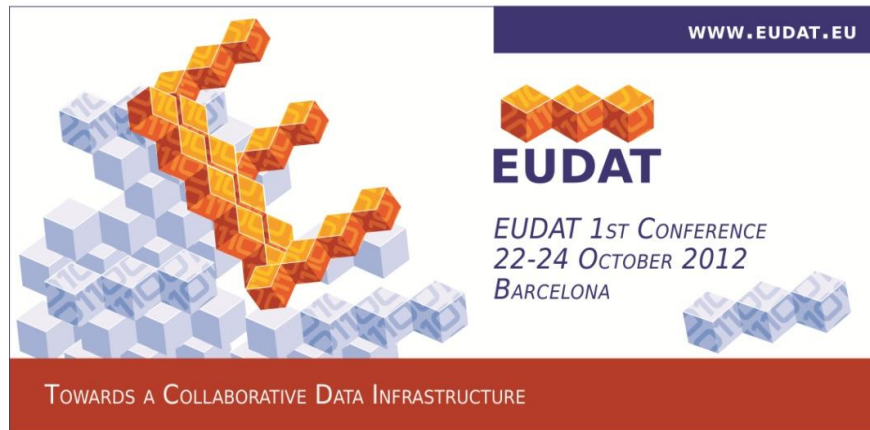- Within the federation
  - **Shibboleth** (Web), **Moonshot** (Non Web)
  - eduRoam (based on **RADIUS**: Remote Authentication Dial-In User Service)
  - **User certificates**, **OpenId**
  - **XACML** (eXtensible Access Control Markup Language)
  - **Oauth2** (Google, Facebook, Microsoft)
- Outside the federation
  - Credential conversion: special SP to create "external" credential

# Challenges

- Leveraging <u>existing identification systems</u>

- Establishing a <u>network of trust</u> among the AA actors: IdPs, SPs, Attribute Authorities, Federations

- <u>Attributes harmonization</u>: it is necessary to agree on a common way to interpret different set of attributes.

22

# Welcome to the 1st EUDAT Conference!

22-24 October 2012, Barcelona

• International event with keynotes from Europe and US

• A forum to discuss the future of data infrastructures

• Project presentations and poster sessions

• 2nd EUDAT User Forum

• Training tutorials

# Welcome to the 1st EUDAT Training Days

**Building Blocks of Data Infrastructures 1 , 25-26 June 2012. Amsterdam**

•25 June (12pm-6pm): Policy-Rule based Data Management

•26 June (9am-11am): Use of Handles (EPIC, DataCite) for Persistent Identification

•26 June (11:30am-3pm): Distributed Authentication and Authorization

**EUDAT**