



EUDAT

Towards a pan-European Collaborative Data Infrastructure

Damien Lecarpentier
CSC-IT Center for Science, Finland
KE Research Data Working Group Meeting
Copenhagen, 14 August 2012



11101001 0101 11010010 10111 100100
011 0111 010011 010
01011000 0111

and now it's time for something
completely different





Topics

- Whats it all about?
- Whos Involved?
- What are we doing?
- Some High Level Tech Stuff...
- How does it relate to this workshop?
- The Future



European Data



EUDAT

- Start date: 1st October 2011
- Duration: 36 Months
- Budget: 16.3 M€ (9.3M€ EC)
- EC Call: INFRA-2011-1.2.2
- Consortium: 25 partners from 13 countries
 - National data centers, technology providers, research
- Objectives:
 - Cost-efficient and high-quality CDI
 - Meetings users' needs in flexible and sustainable way
 - Across geographical and disciplinary boundaries



<http://www.eudat.eu>

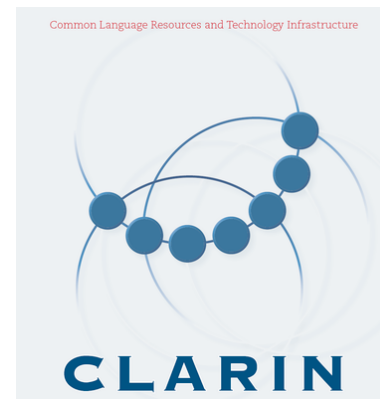
EUDAT Consortium



Data centers and Communities



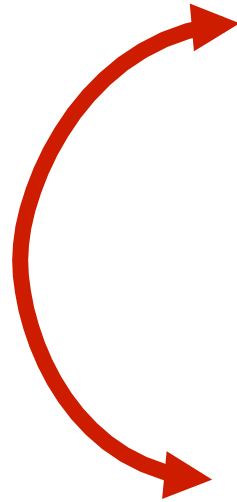
Communities



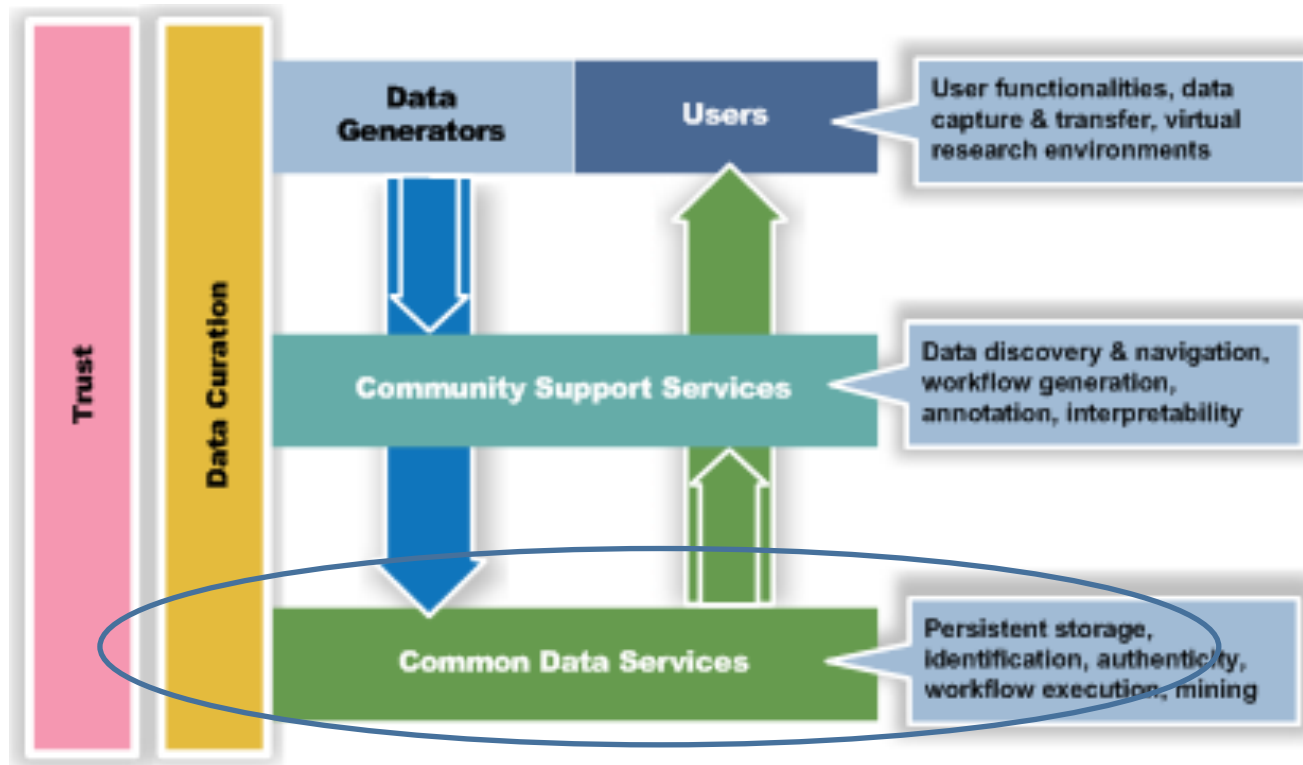
Communities and Data Centers

What are the basic requirements?

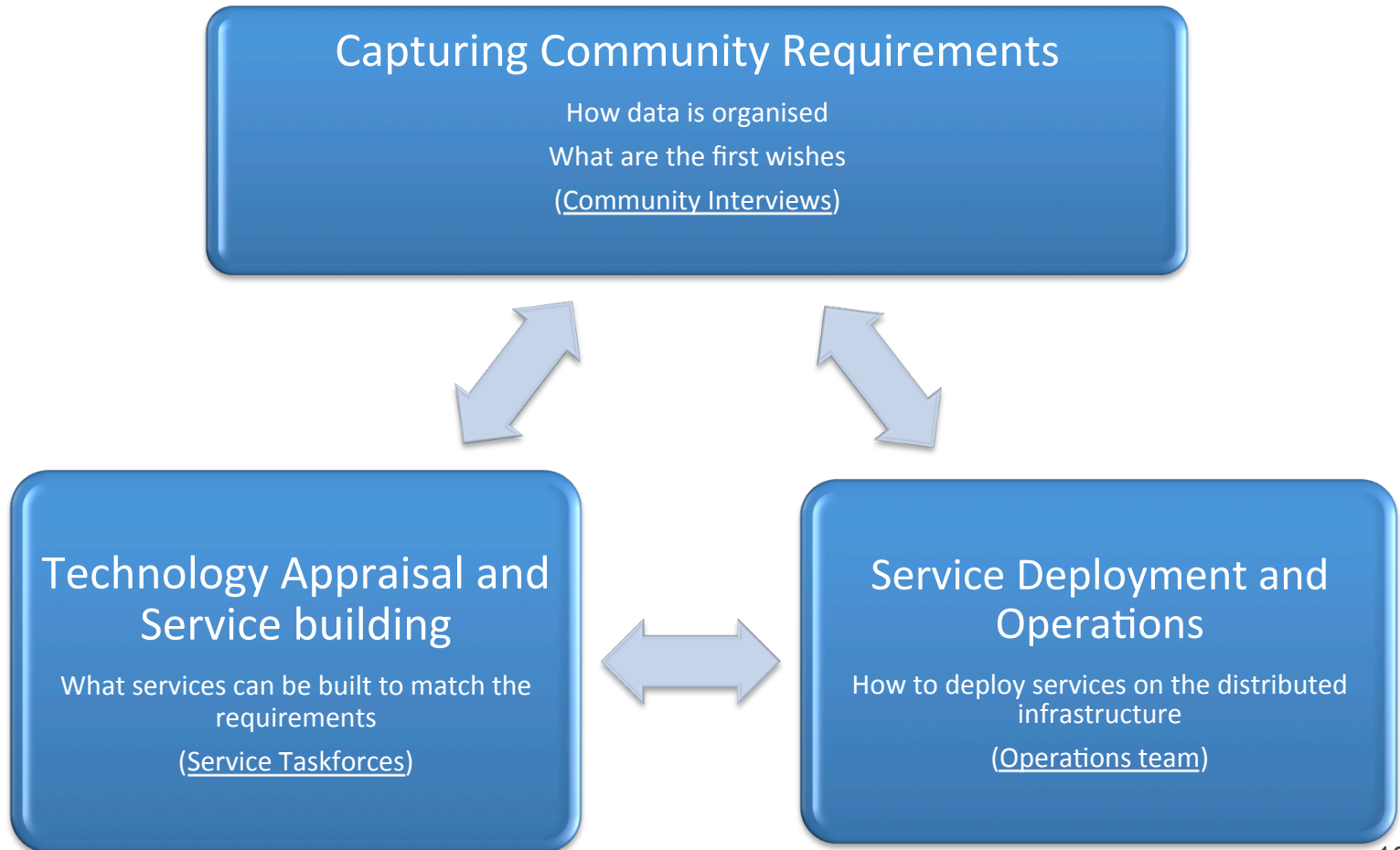
Which common services are needed?



The CDI concept



How Do We Achieve This?



What are the Requirements?

6 service/use cases identified

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

Dynamic replication: Enable communities to perform (HPC) computations on the replicated data.

Metadata: Create a common metadata domain for all data stored by EUDAT data centres and a searchable catalogue covering all the data stored within EUDAT, allowing data searches

Research data store: create an easy-to-use service that will enable researchers and scientists to upload, store and share data that are not part of the officially-managed data sets of the research communities

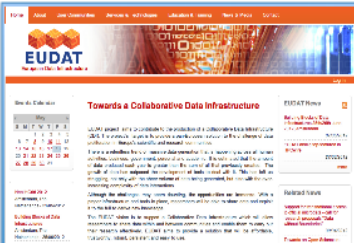
PID: a robust, highly available and effective PID system that can be used within the communities and by EUDAT.

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

EUDAT Core Services

Community-oriented

Enabling Services



EUDAT Portal

Integrated APIs and harmonized access to EUDAT facilities



Metadata Catalog

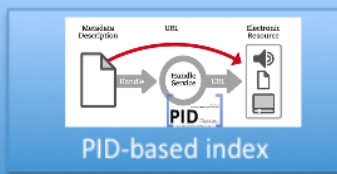


**Aggregated EUDAT metadata domain.
Data inventory**

Requirement
Provide an inventory of metadata across disciplines

Function
• Joint metadata domain
• Catalog indexing stored data

eudat-metadata@postit.csc.fi



AAI



Network of trust among authentication and authorization actors

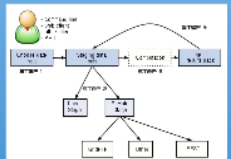
Data Staging



Dynamic replication to HPC workspace for processing

Requirement
Provide a service to stage data between EUDAT infrastructure and HPC/HTC resources

Function
Dynamically replicate subset of data stored in EUDAT to HPC workspace



eudat-datastaging@postit.csc.fi

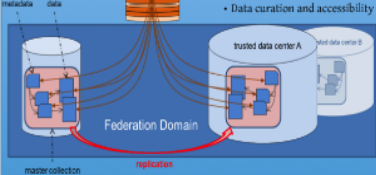
Safe Replication



Data curation and access optimization

Requirement
Provide a service to replicate and curate data to selected data center(s)

Function
• Safe replication from one data center to another
• Data curation and accessibility



eudat-safereplication@postit.csc.fi

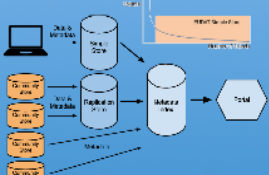
Simple Store



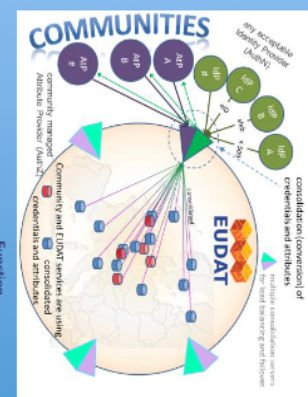
Researcher data store (simple upload, share and access)

Requirement
Provide a simple service to store user data temporarily

Function
Simple upload
Store data



eudat-simplestore@postit.csc.fi



Requirement
Provide a working AAI system in a federated scenario

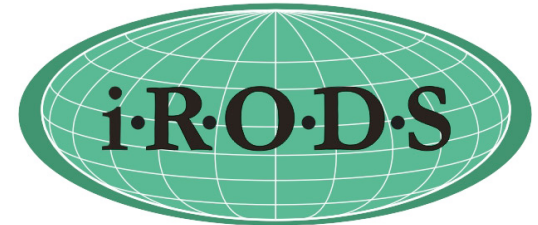
Function
• Integrate existing identification systems
• Establish a network of trust among AAI IdP and SP providers, attribute authorities and federations attribute harmonization

eudat-AAI@postit.csc.fi

Building Blocks of the Collaborative Data Infrastructure

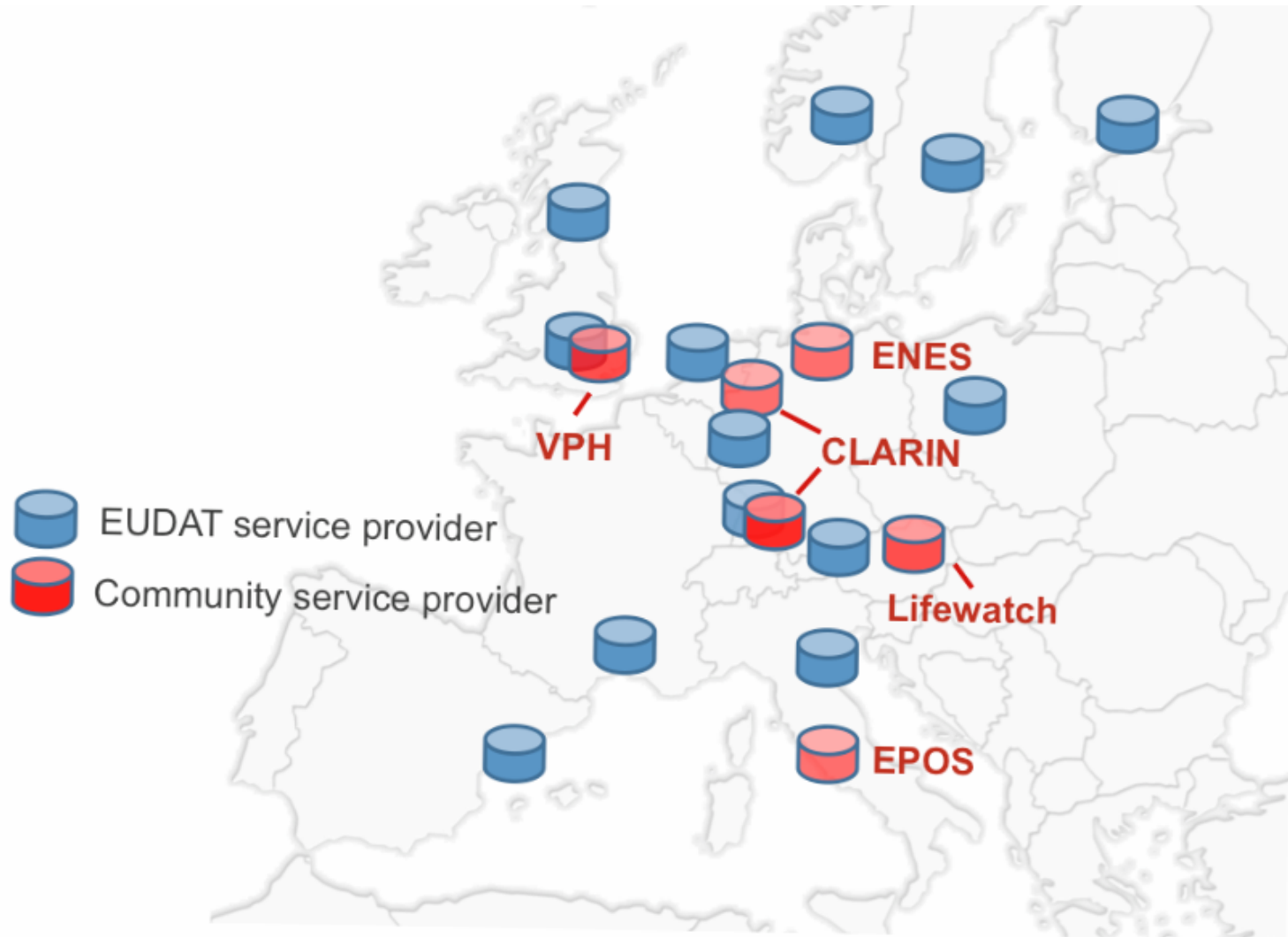
Its All About The Rules

- Federation Based on IRODS
 - Wide support and use 😊
 - Micro services allow new rules 😊
 - Federation built in 😊
 - Allows C and python plugins 😊
 - High Transfer Overhead 😞
 - Supports two interfaces to storage that don't work together 😞
 - Uncertain Scalability of database 😞



Integrated Rule-Oriented Data System

INFRASTRUCTURE



The Same, but Different

- WLCG

- Built from Scratch
- Homogenous Data
- Single Discipline
- Any Data, Any Time, Any Where
- Security
- Data Management
- Metadata
- No Services

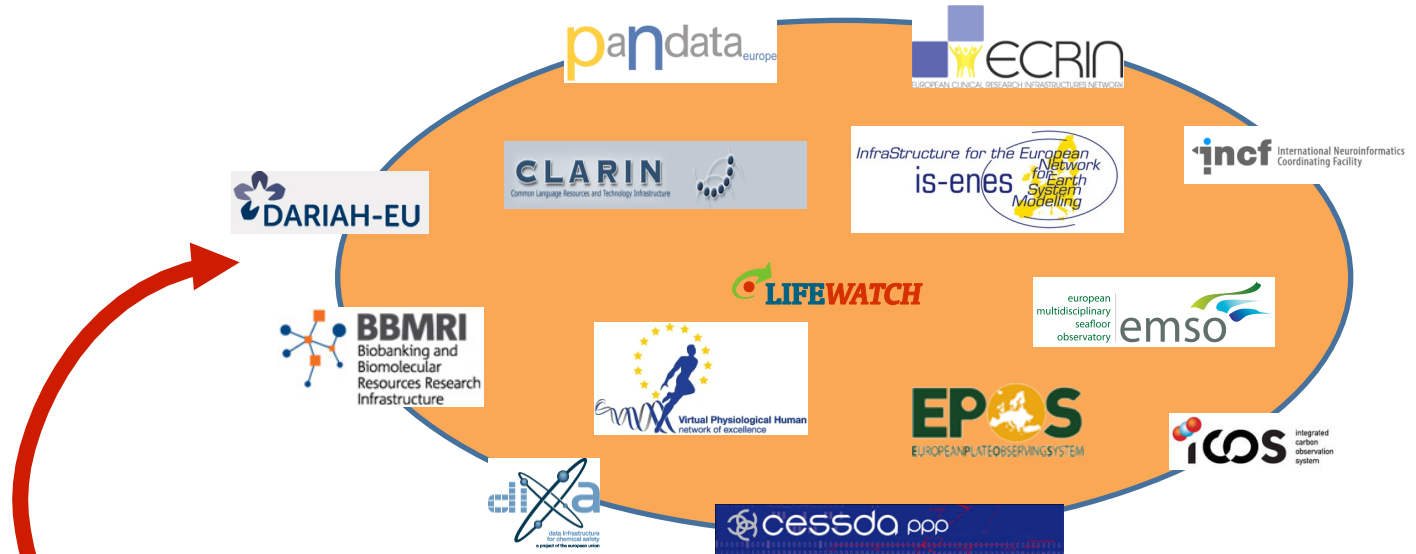
- EUDAT

- Data already exists
- Heterogeneous Data!
- Multi Discipline
- Any Data, Any Time, Any Where
- Security (Big Time)
- Data Management
- Metadata
- Services

Communities and Data Centers

What are the basic requirements?

Which common services are needed?



Other Technologies

- XrootD
 - Scales well
 - Extensible API
 - But can it do what we want?
 - Fast Protocol allowing direct and streaming access
- HTTP
 - Scales well
 - Need separate rule engine
 - No direct access?

How Do We Sustain This?

▪ Organisational Model

- How do we move for a project collaboration to a federated infrastructure?
- Which are the actors of this infrastructure and what is/are their role(s)?
- How do we integrate new members?
- How will the infrastructure interact with other infrastructures and projects?

▪ Costs and Funding Models

- Who will pay for the infrastructure and the shared services?
- What are the costs of the services?
- How to define a business model that best support the interest of research communities, data centers and funders?

Contact Us



<http://www.eudat.eu>



eudat-info@postit.csc.fi

Project Coordinator: Kimmo Koski
kimmo.koski@csc.fi

Scientific Coordinator: Peter Wittenburg
peter.wittenburg@mpi.nl

Project Manager: Damien Lecarpentier
damien.lecarpentier@csc.fi

Dissemination Manager: Nagham Salman
nagham.salman@bsc.es

Industry Task Force: David Manset
dmanset@maatg.fr

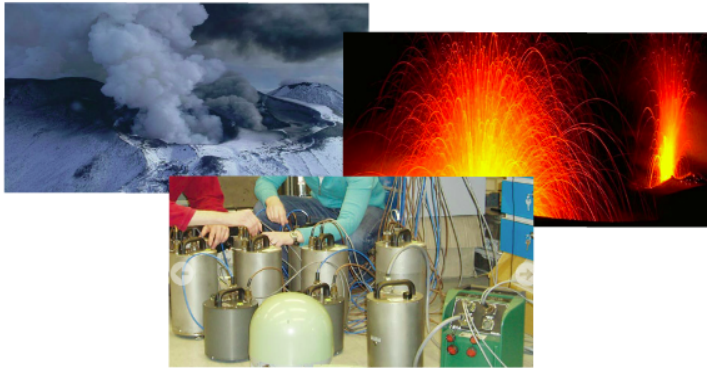




BACKUP SLIDES...

EPOS - European Plate Observatory System

- Distributed data sensors
- Large scale statistics
- Metadata schema
- Reference architecture



Research Infrastructure and E-Science for Data and Observatories on Earthquakes, Volcanoes, Surface Dynamics and Tectonics

CLARIN - Common Language Resources and Technology Infrastructure

- About 200 centers in EU
- Require PIDs, CMDI
- ISOcat, SCHEMcat
- Virtual Language Obs.

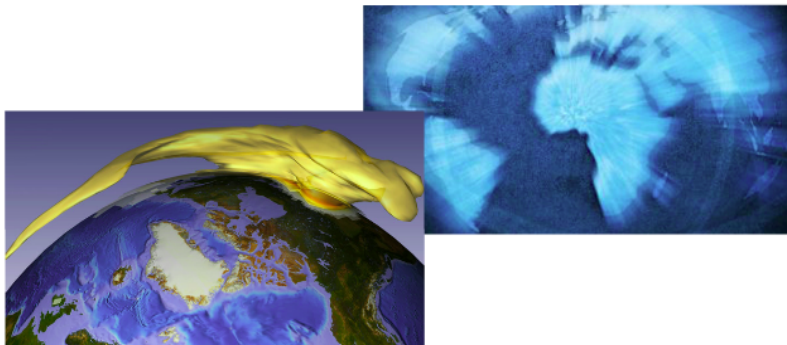
<http://www.clarin.eu/vlo/>



The CLARIN project is a large-scale pan-European collaborative effort to create, coordinate and make language resources and technology available and readily usable

ENES - Service for Climate Modeling in Europe

- About 20 centers in EU
- CIM data model
- Using CDI @ German Climate Center
- Using DOIs and EPIC
- Metadata based on ISO 11179



enes European Network for Earth System Modelling

Welcome

ENES Townhall Meeting at EGU 2010: Here is the [announcement!](#)

For latest news on IS-ENES click [here!](#)

A major challenge for the climate research community is the development of comprehensive Earth system models capable of simulating natural climate variability and human-induced climate changes. Such models need to account for detailed processes occurring in the atmosphere, the ocean and on the continents including physical, chemical and biological processes on a variety of spatial and temporal scales. They have also to capture complex nonlinear interactions between the different components of the Earth system and assess, how these interactions can be perturbed as a result of human activities.

Accurate scientific information is required by government and industry to make appropriate decisions regarding our global environment, with direct consequences on the economy and lifestyles. It is therefore the responsibility of the scientific community to accelerate progress towards a better understanding of the processes governing the Earth system and towards the development of an improved predictive capability. An important task is to develop an advanced software and hardware environment in Europe, under which the most advanced high resolution climate models can be developed, improved, and integrated.

ENES provides information and services to foster intricate simulations of the climate system using high performance computers as well as the distributions and dissemination of data produced by such simulations

VPH - The Virtual Physiological Human

- Pilot project with 5 hospitals
- Centralized data center
- Metadata aggregation
- DICOM, JPEG headers

<http://www.vph-share.eu/>



A screenshot of the VPH NoE website homepage. The page features a navigation menu with links for Home, WP1, WP2, WP3, WP4, WP5, VPH-I, MIP, and Login. The main content area includes a welcome message, a central diagram showing the VPH NoE structure with sub-components like VPH Initiative, VPH for Researchers, and VPH for Clinicians. There are also sections for Project, Activities, and Latest VPH Events. The page is designed with a clean, professional layout and includes a search bar and a sidebar with highlights.

VPH aims to help support and progress European research in biomedical modeling and simulation of the human body. This will improve our ability to predict, diagnose and treat disease, and have a dramatic impact on the future of healthcare, the pharmaceutical and medical device industries

LifeWatch - Biodiversity Data and Observatories

- Distributed data sensors
- Metadata standardisation
- Interoperability reqs
- Involving most nature infrastructures
- Common reference model

<http://envri.eu/>

<http://creative-b.eu/>



LIFEWATCH e-science and technology infrastructure for biodiversity data and observatories

Home Contact About Participants Get involved News Cases Events Press Documents

LIFEWATCH COUNTRIES
Austria Belgium Denmark Finland France Greece Hungary Italy Netherlands Norway Poland
Portugal Romania Slovak Republic Slovenia Spain Sweden Turkey United Kingdom

LIFEWATCH NEWS
2011-02-16 **LIFEWATCH RESEARCH INFRASTRUCTURE STARTS CONSTRUCTION IN 2011** - The initial country consortium establishing the LifeWatch research infrastructure agreed to finance... [Read more](#)
2011-01-19 **LIFEWATCH CLOSING EVENT** - On this page you can download all the slides presented at the closing event of the LifeWatch preparatory project... [Read more](#)
2011-01-17 **LIFEWATCH CONSTRUCTION KICKS OFF ON JANUARY 19TH** - On 19 January 2011, at the closing conference of the LifeWatch preparatory project a first group of... [Read more](#)

LIFEWATCH FOCUS
LifeWatch research infrastructure starts construction in 2011
The initial country consortium establishing the LifeWatch research infrastructure agreed to finance the start-up activities for the infrastructure construction. These countries will host the Common Facilities of LifeWatch.
On 19th January 2011 representatives from organisations in Hungary, Italy, the Netherlands, Romania and Spain signed a Memorandum of Understanding to cooperate for an early start of the LifeWatch infrastructure for biodiversity and ecosystem research. The LifeWatch Stakeholders Board, representing the ten countries aiming at establishing the LifeWatch ERIC, welcomed the initiative to start early construction.

Newsletter
Subscribe to our newsletter. Send an email to info@lifewatch.eu

Quote
"Through our Memorandum of Understanding GBIF and LifeWatch, based on our respective complementary mandates, now have a formal framework for co-operation and collaboration on infrastructural developments, building on GBIF's 10 years of investment to date."
Dr. Nick King
Director Global Biodiversity Information Facility (GBIF)

LifeWatch will construct and bring into operation the facilities, hardware, software and governance structures for all aspects of biodiversity research. Facilities for data generation and processing, data integration and interoperability. A network of observatories, virtual laboratories. A Service Center supporting scientific and policy users.

How Requirements are Shared?

Service	SR	DR	MD	SS	PID	AAI
Community						
CLARIN	X	+	X	X	+	X
ENES	X	X	X		+	X
EPOS	X	X			X	X
VPH	X	X			X	X
LifeWatch	X	+	X	+	+	X

NB: “X”= this service is relevant to this community, “+” = this community has interest in this service but at a later stage or has a similar service already running in production.

SAFE_REPLICATION@EUDAT

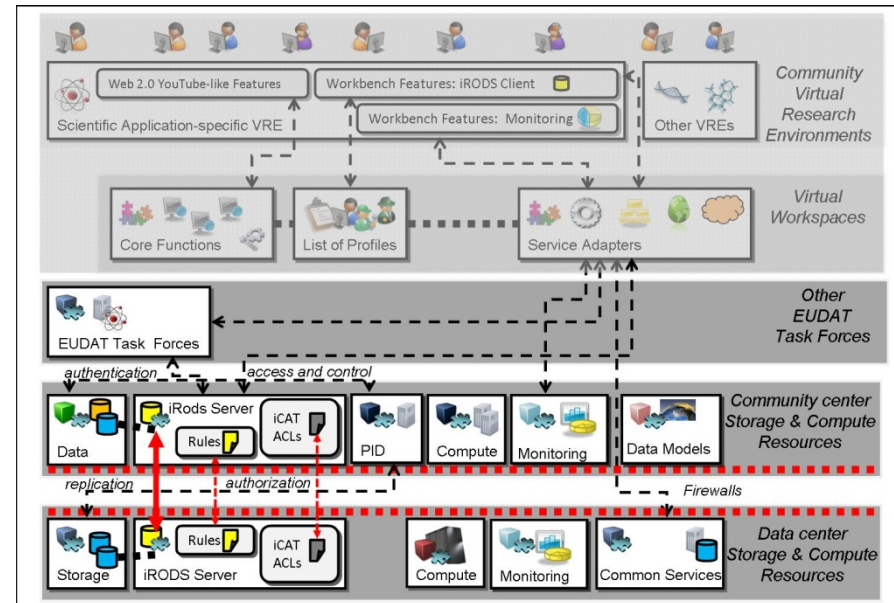
Objective: Enable communities to easily replicate data to selected data centers for storage in a robust and reliable manner.

Key benefits: data bit stream preservation, more optimal data curation, better accessibility

Description: Data replication management based on users' requirements and constraints; data replication solutions and services embedded into critical security policies, including firewall setups and user accounting procedures.

Technology: iRODS to be used as an initial replication middleware, implemented across the community centers and data centers; as more user communities join the task force, other storage technologies may be added, depending on user needs.

➤ Production setup expected by 2013, such that users will be able to safely replicate data across different user community centres and data centres.



Integrated Rule-Oriented Data System

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

DATA_STAGING@EUDAT

Objective: Enable communities to perform (HPC) computations on the replicated data

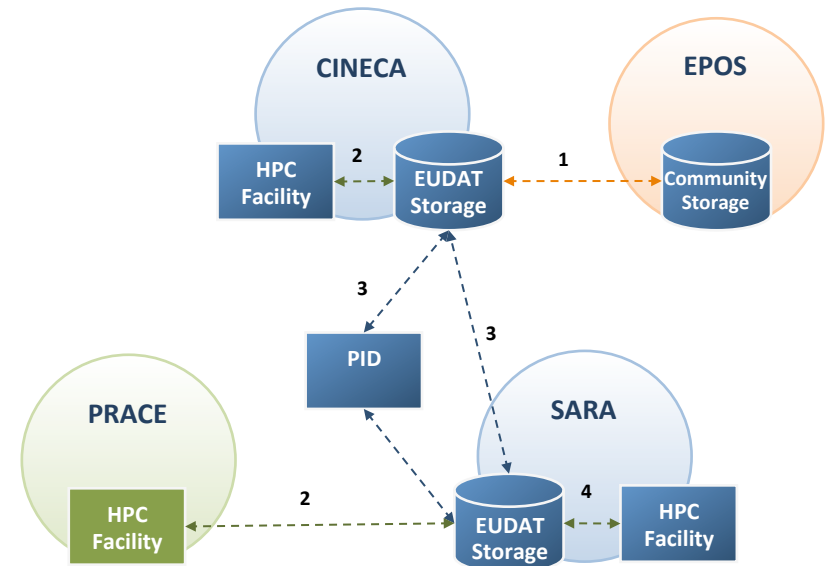
Key benefits: Access to large computing facilities

Description: This service will allow the EUDAT communities to dynamically replicate subsets of their data stored in EUDAT to HPC machine workspaces for processing.

Differences with the safe replication scenario:

- replicated data are discarded when the analysis application ends;
- Persistent Identifier (PID) references are not applied to replicated data into HPC workspaces;
- Users initiate the process of replicating data while in the safe replication scenario data are replicated automatically on a policy basis.

Technologies: GridFTP, Griffin, gTransfer, FTS (under appraisal)



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

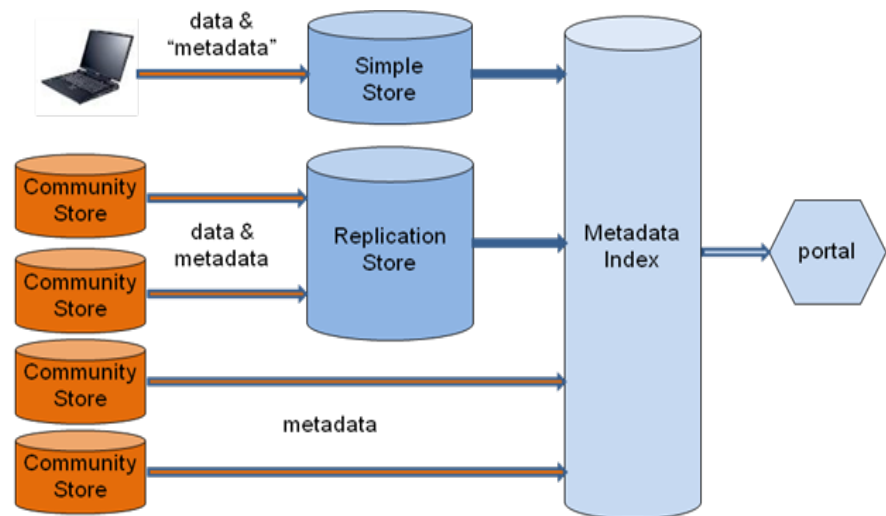
METADATA@EUDAT

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.

Key benefits: Advertising platform for data sets, metadata service for less mature communities

Description: EUDAT will handle metadata for more resources than just those deposited within the EUDAT CDI. In the initial phase we will target mainly resources contributed by the participating communities augmented with those of interested well-organized communities that are ready to contribute. Then, later, other interested communities can be approached depending on the respective community capabilities.

Technology: OAI-PMH and embeds domain specific metadata, as XML, within the OAI-PMH record



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

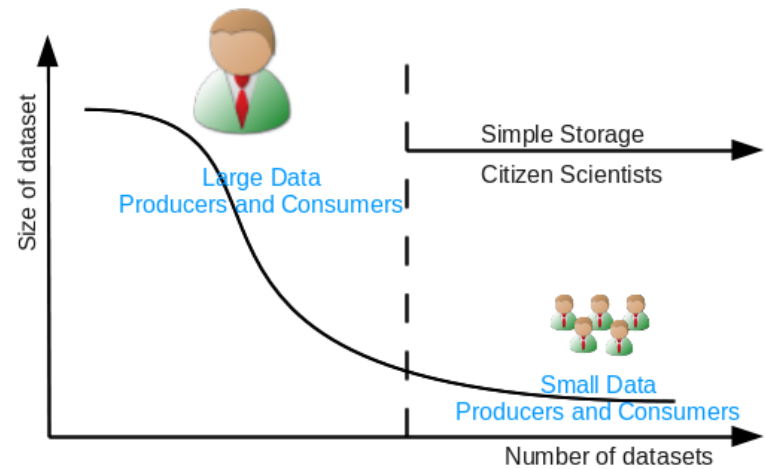
SIMPLE_STORE@EUDAT

Objective: create an easy-to-use service that will enable researchers and scientists to upload, store and share data that are not part of the officially-managed data sets of the research communities.

Key benefits: Store, share, and retrieve smaller sets of data not officialt handled.

Description: This service will address the long tail of "small" data, and the researchers/citizen scientists creating and manipulating it. Typically this type of data comes in a wide range of formats including text, spreadsheets, number series, audio and video files, photographs and other images. The Research Data Store is complementary to the other EUDAT services that manage the large volumes of official community data.

Technologies: Invenio, figshare, beehub and MyExperiment.



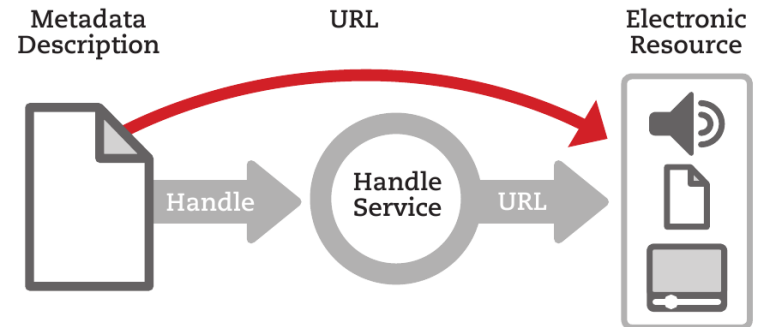
More info: eudat-simplestore@postit.csc.fi

PIDS@EUDAT

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.

Technologies: 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

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:

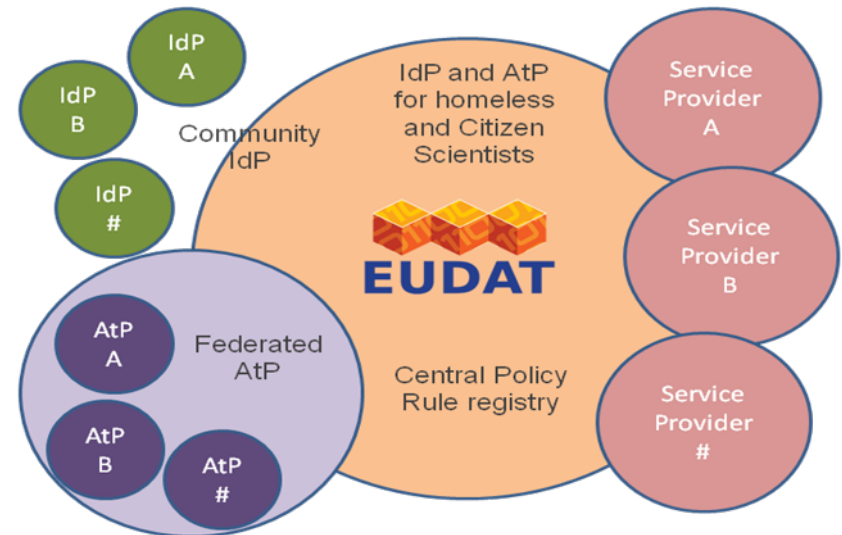
Leveraging existing identification systems within communities and/or data centers

Establishing a network of trust among the AA actors:
Identify Providers (IdPs), Service Providers (SPs), Attribute Authorities and Federations

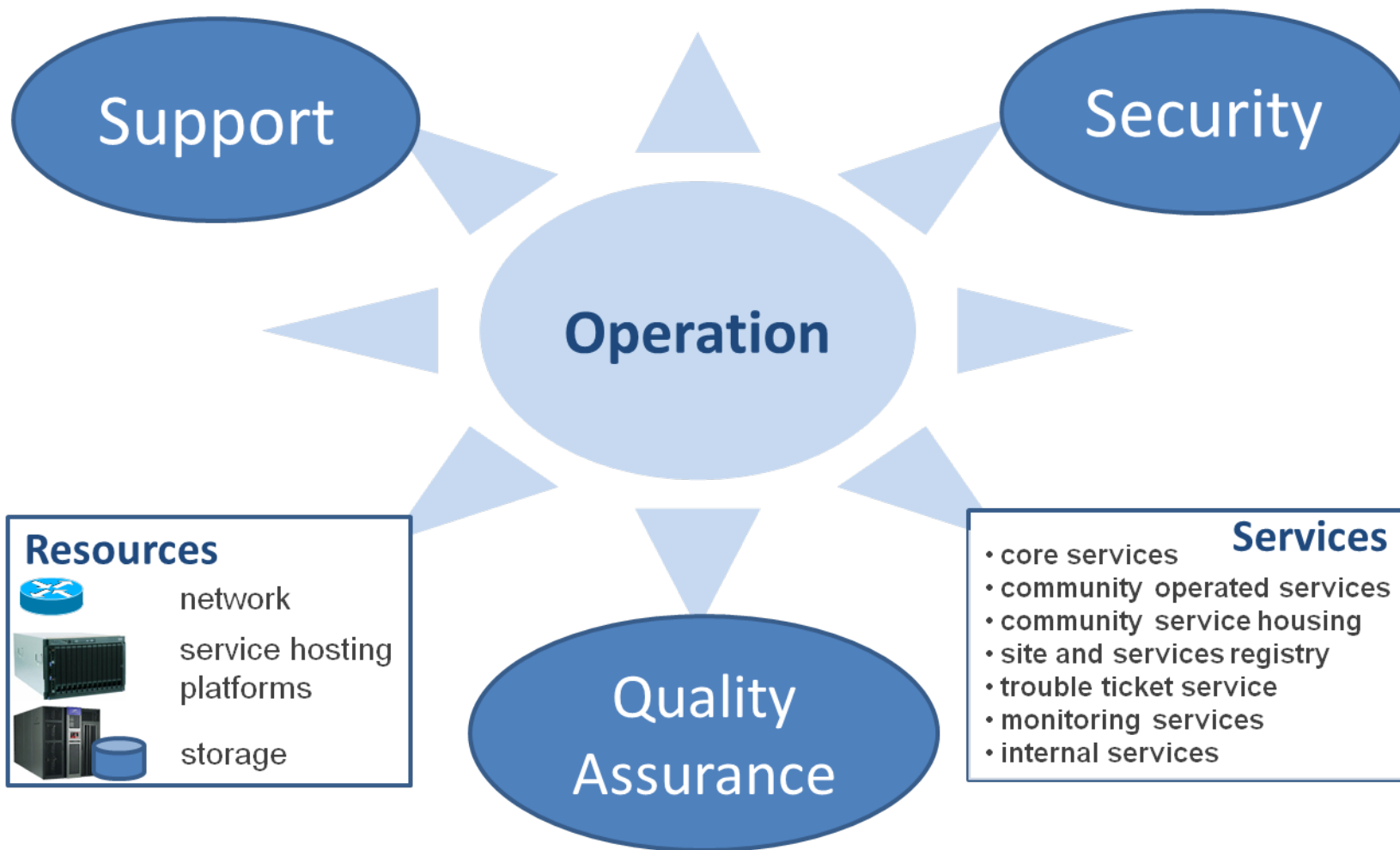
Attribute harmonization

Technologies: Oauth2, OpenID, RADIUS, SAML2, X.509, XACML, etc.

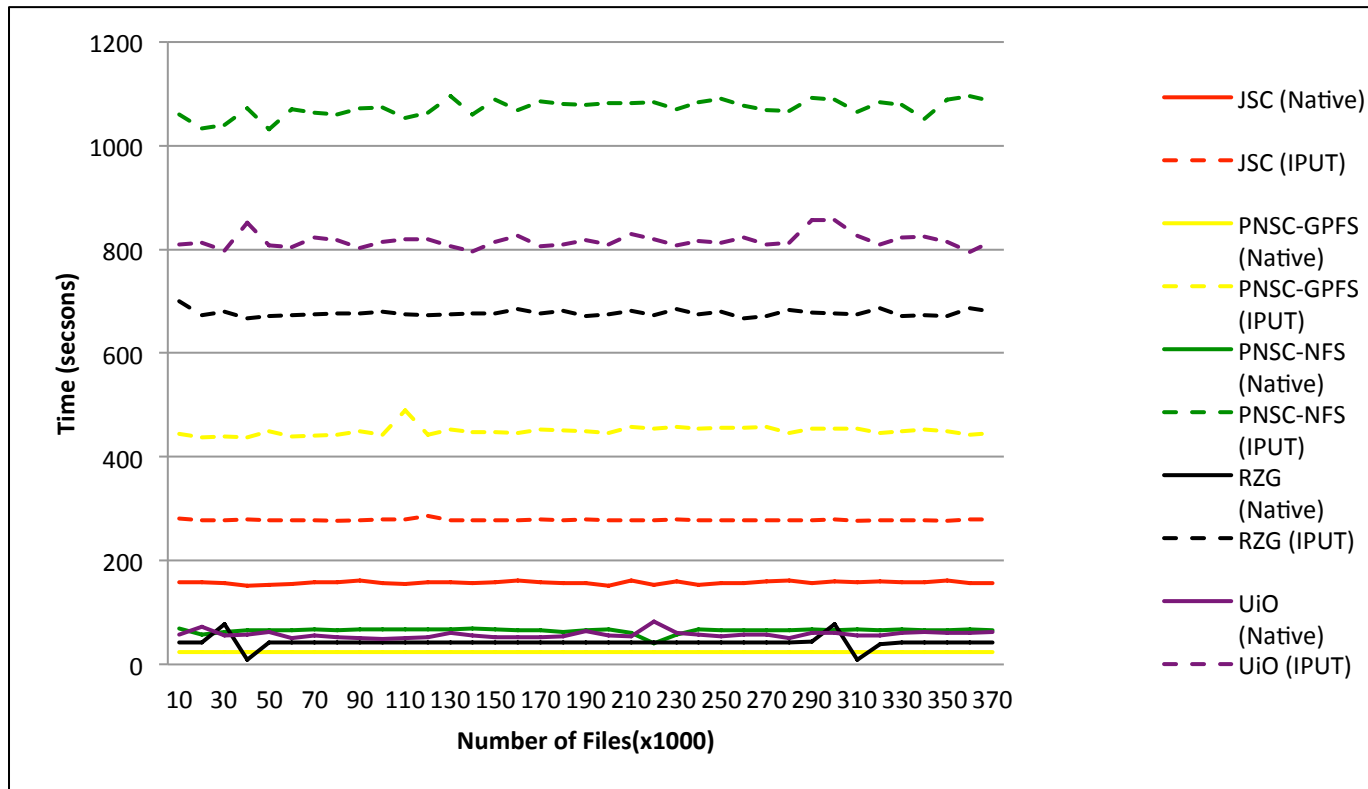
More info: eudat-AAI@postit.csc.fi



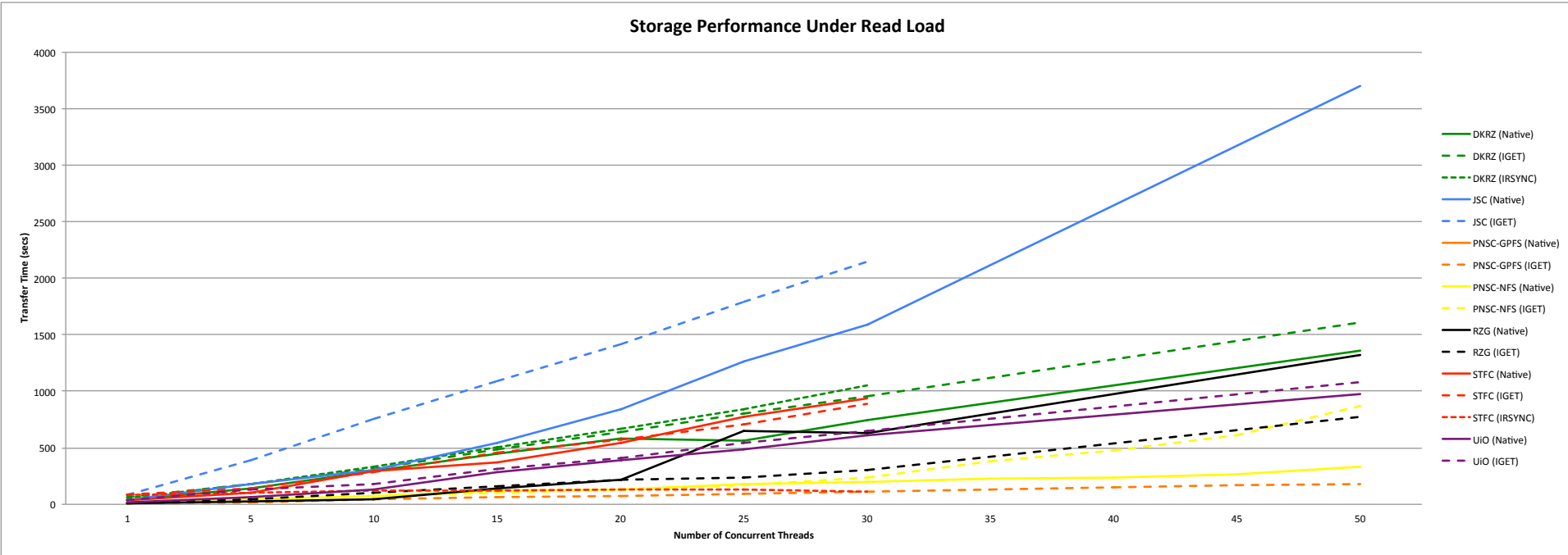
OPERATION TEAM



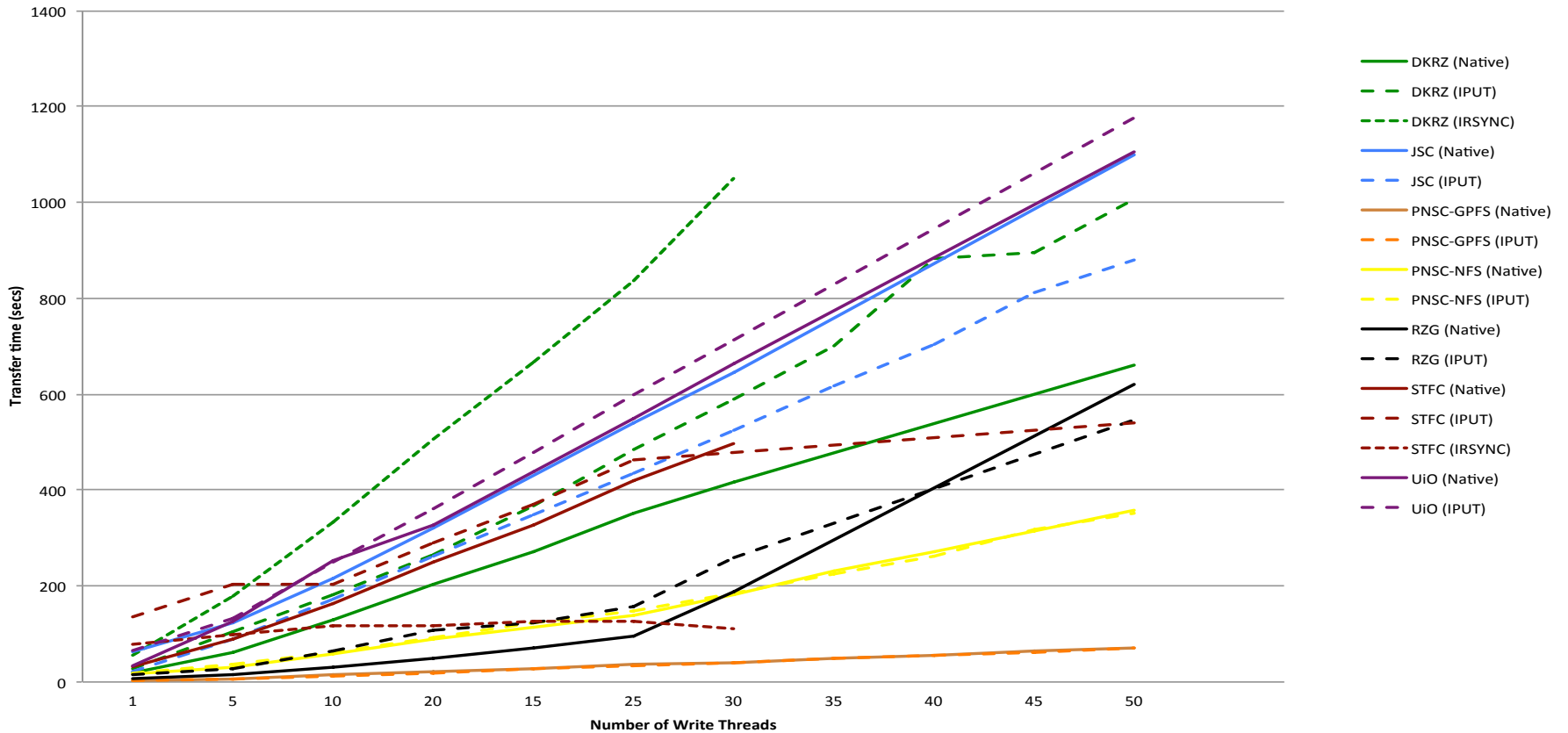
Preliminary results



And More



And more



And finally

