# Safe Data Replication Service

## A Simple and Reliable EUDAT Service for Scientific Communities

Pre-Production State

Morris Riedel et al.
Juelich Supercomputing Centre
EUDAT Conference, Barcelona

Date: 23th October 2012

# Safe Replication Service in a Nutshell



metadata    data

PID Domain

trusted data center A

sted data center B

Federation Domain

master collection

**replication**

**Safe Replication**
Data curation and access optimization

**Make data referencable**

**Better accessibility of scientific data**

**More optimal data curation**

**High degrees of reliability and trust**

EUDAT

# Federated Approach for Use Cases



Create M replications at different data centers for N years,
exclude data centers X to data centers Z from the replication scheme
and make them all accessible by maintaining the given access permissions.

# Forming Strong EUDAT Collaborations

EPOS - European Plate Observatory System

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

**EPOS**
EUROPEAN PLATE OBSERVING SYSTEM

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

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 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

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
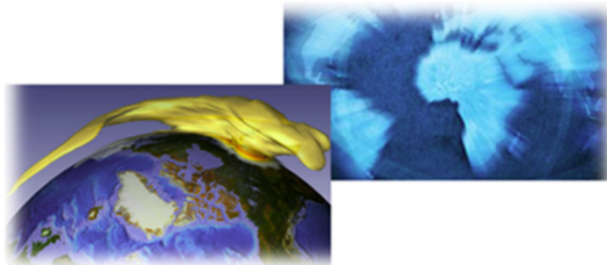
| Scientific Community | Community Centers | Data Centers |
|---|---|---|
| CLARIN | MPI-PL | RZG, SARA |
| ENES | DKRZ | JSC, CSC |
| EPOS | INGV | CINECA, SARA |

**EUDAT**

# Use Case Example: Climate Science Data

- ## ENES: Service for Climate Modelling in Europe
  - Provides services to foster intricate simulations of the climate system using high-performance computers
  - **Enables the distribution and dissemination of data produced by such simulations**
  - Other Facts: *about 20 EU centres; CIM data model; uses DOIs and EPIC handles; metadata in ISO 11179;*



**Slide content kindly provided by Hannes Thielmann, DKRZ, a climate scientist in earth system modelling**

# Concrete Replicated Climate Scientific Data

- Complexity: ENES & CMIP5 & IPCC AR5
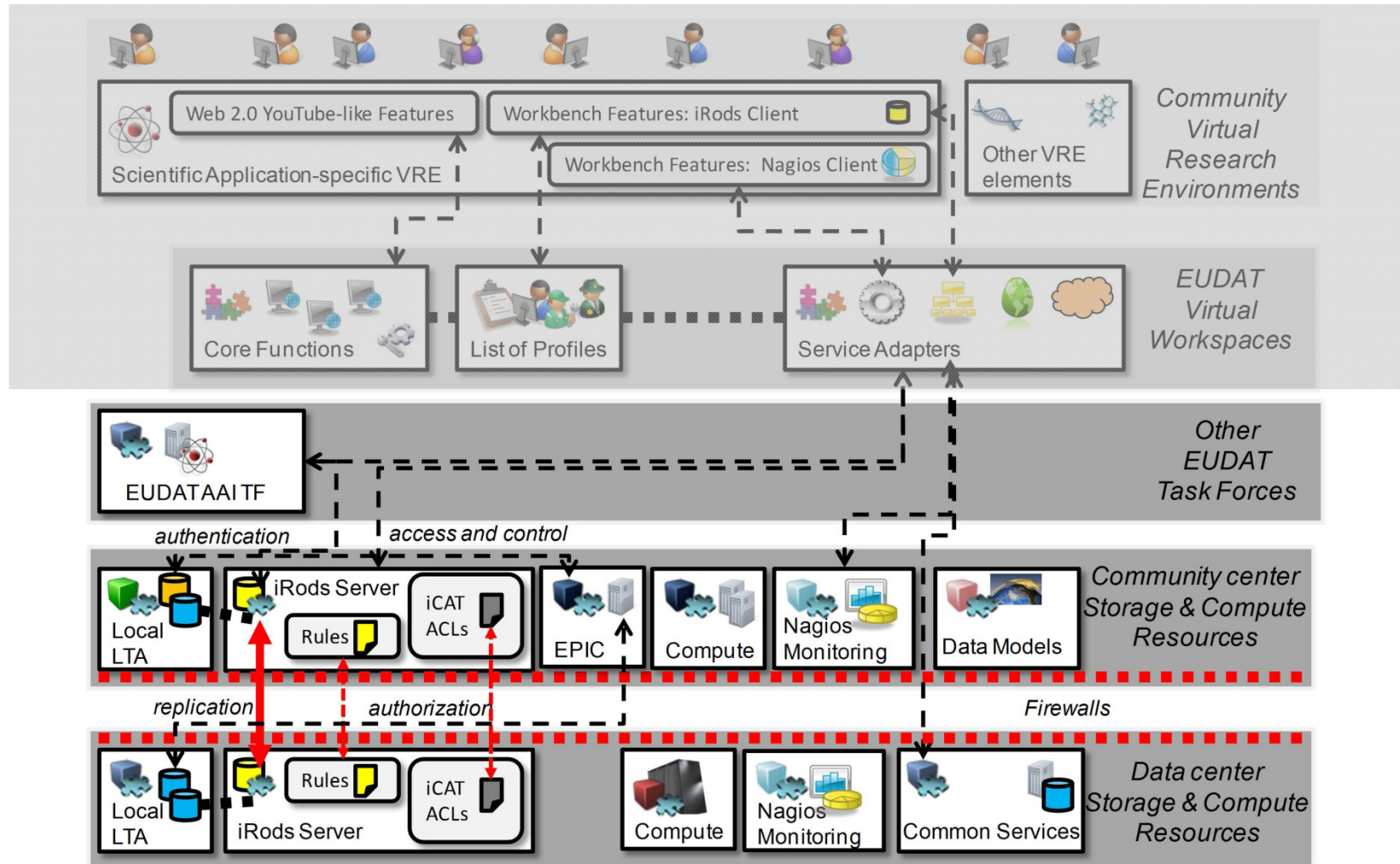  - ENES contributes to the Intergovernmental Panel on Climate Change (IPCC) Fifth Assessment Report (AR5).
  - Coupled Model Intercomparison Project Phase 5 (CMIP5) model data that will serve as the basis for IPCC AR5.
  - This data prepared will be made available to the international climate community.
  - The Earth System Grid Federation (ESGF) is a partnership of climate modeling centers created to provide secure, web-based, distributed access to CMIP5 model data.

**EUDAT**

# Use Cases Derived Reference Architecture

# Selected References

- Documentation on iRODS and EPIC/Handle system available on the Web

- 1st EUDAT Conference Training Day - Many training sessions yesterday!
  - PID handling & services, iRODS policies, rules, micro-services, etc.

- EUDAT Newsletter April 2012
  - Check the EUDAT WebSite

- M. Riedel and P. Wittenburg et al.
  *'A Data Infrastructure Reference Model
  with Applications - Towards
  Realization of a ScienceTube
  Vision with a Data Replication Service'*,
  Journal of Internet Applications,
  to be published early 2013

- Contact to specialists:
  eudat-safereplication@postit.csc.fi