# EUDAT Training
## 2ⁿᵈ EUDAT Conference, Rome
## October 28ᵗʰ

# Introduction, Vision and Architecture

Giuseppe Fiameni – CINECA
Rob Baxter – EPCC
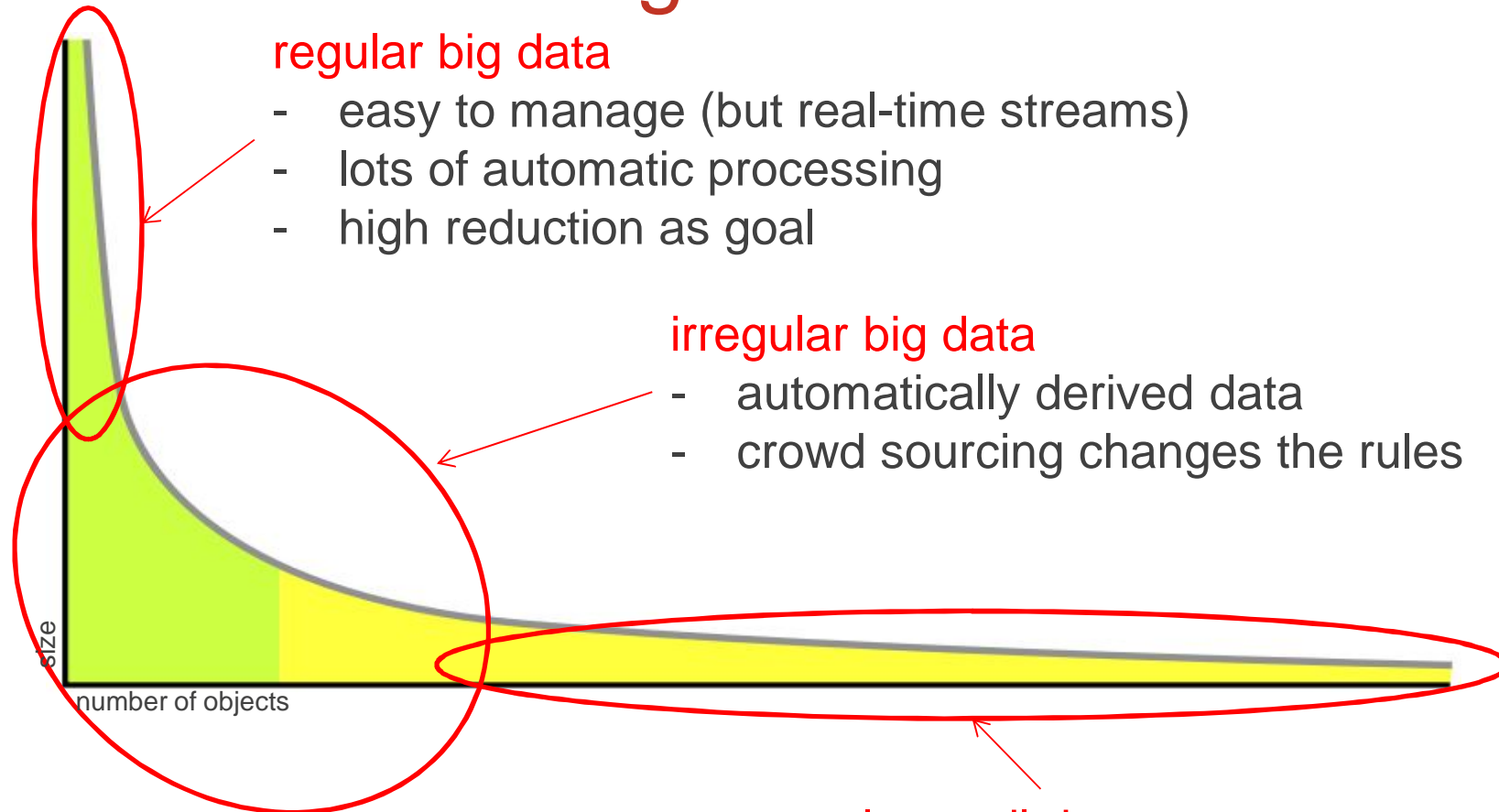EUDAT members

SEVENTH FRAMEWORK
PROGRAMME

EUDAT

# Agenda

- Background information
- Services
- Common Data Infrastructure

EUDAT

# Setting the scene



regular big data
- easy to manage (but real-time streams)
- lots of automatic processing
- high reduction as goal

irregular big data
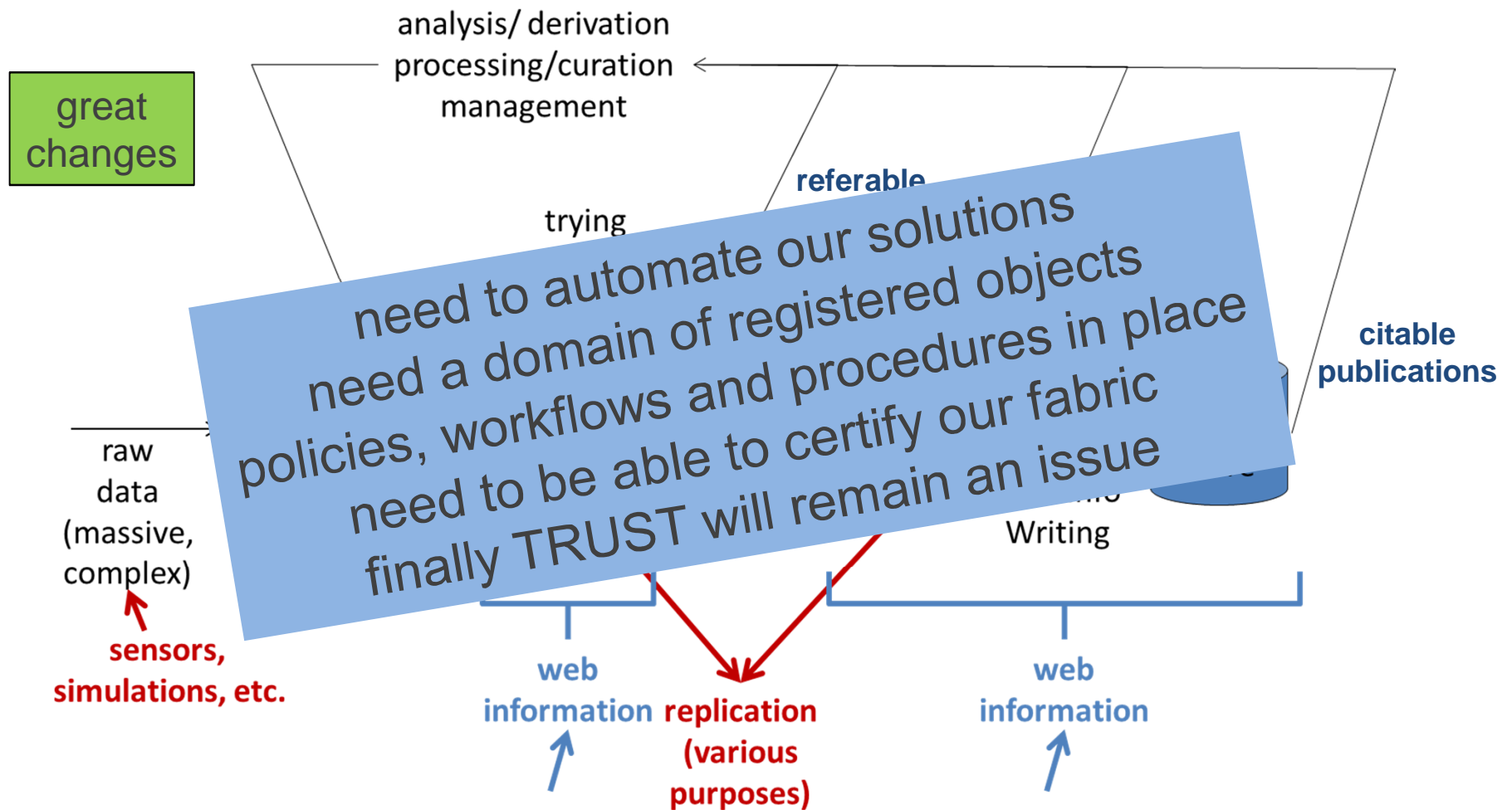- automatically derived data
- crowd sourcing changes the rules

all the same for industry, government, public services, citizens, etc.

long tail data
- difficult to manage
- lots of relations

*size*

*number of objects*

EUDAT

# big scientific data –> the data fabric

great changes

analysis/ derivation
processing/curation
management

referable

trying

raw
data
(massive,
complex)

**sensors,
simulations, etc.**

citable
publications

Writing

need to automate our solutions
need a domain of registered objects
policies, workflows and procedures in place
need to be able to certify our fabric
finally TRUST will remain an issue

web
information **replication
(various
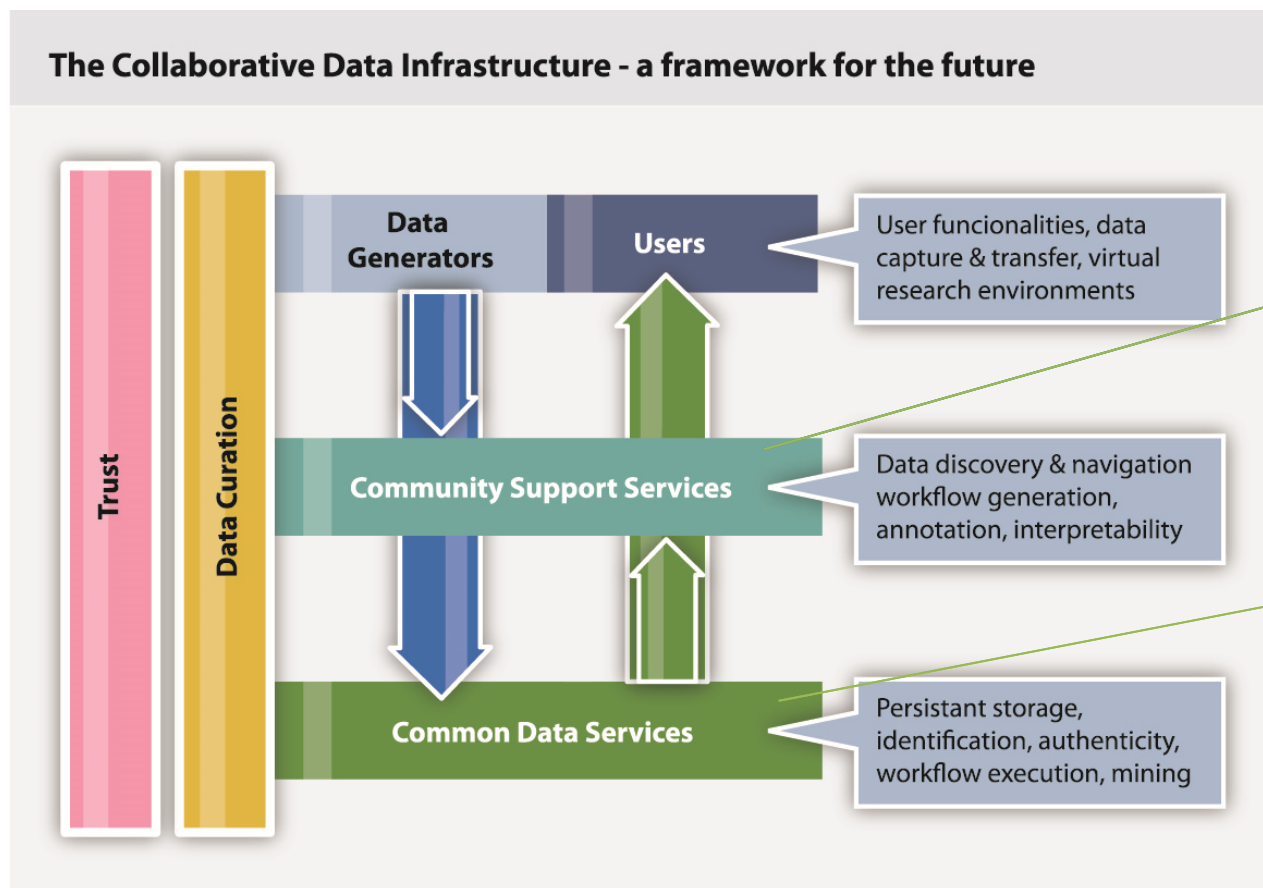purposes)**

web
information

EUDAT

# EUDAT: Vision & Architecture

- EUDAT began with the concept of the Collaborative Data Infrastructure
  - See "Riding the Wave" (High Level Expert Group on Scientific Data, Final Report, 2010)

- This identified a handful of core Service Cases
  - See http://www.eudat.eu/services-and-technologies

- And the implementation of the Service Cases led to our current distributed Architecture
  - See later ☺

# EUDAT's mission: common services in CDI



The Collaborative Data Infrastructure - a framework for the future

Trust | Data Curation

**Data Generators** | **Users** — User funcionalities, data capture & transfer, virtual research environments

**Community Support Services** — Data discovery & navigation workflow generation, annotation, interpretability

**Common Data Services** — Persistant storage, identification, authenticity, workflow execution, mining

CLARIN, LifeWatch, ENES, EPOS, VPH, INFC etc.
6 Core Infrastructures about 20 infrastructures

⇒ 12 EUDAT data centers
⇒ and/or cross-disciplinary initiatives

diagram taken from EC's HLEG report "Riding the Wave"

**EUDAT**

# What is the EUDAT CDI?

- The EUDAT Collaborative Data Infrastructure is

  - a **pan-European**, **cross-disciplinary** domain of research data for both **big community** researchers and **"long tail"** scientists

  - where data are **registered, preserved**, **accessible** and made **re-usable**

# What does this mean?

- ***Pan-European***
  - Fundamentally, a wide-area distributed architecture

- ***Cross-disciplinary***
  - Five core stakeholder communities, many other interested; many sources of conflicting requirements!
  - Including simplified services to encourage the "long tail" to participate
  - All implies a significant systems integration challenge!

**EUDAT**

# What does this mean? (2)

- ***Registered*** means EUDAT data are
  - Globally identified and discoverable (the PID Service)

- ***Preserved*** means EUDAT data are
  - Stored at big European HPC and data centres
  - Replicated for safety (the Safe Replication Service)
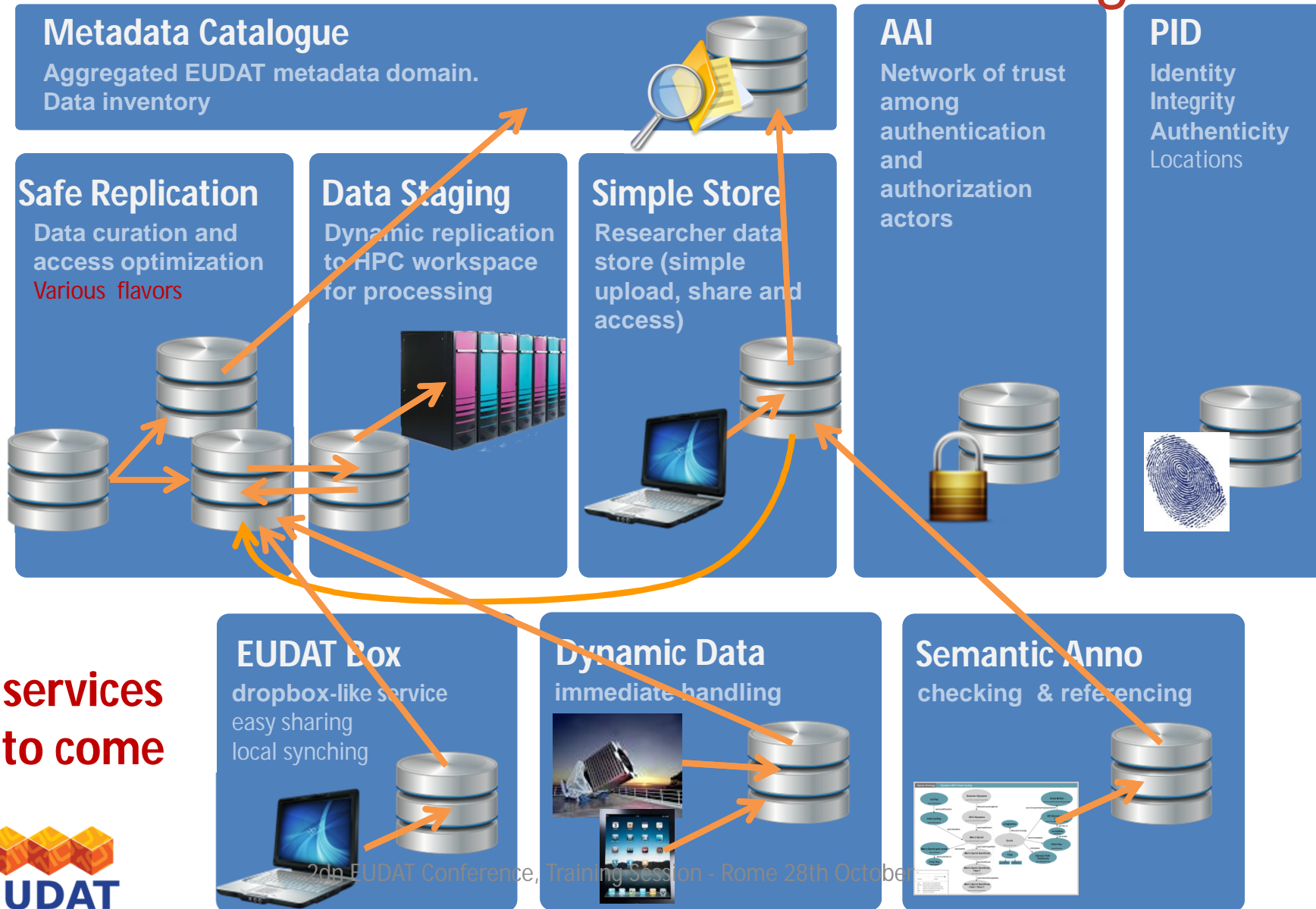  - Governed by policy rules (the Policy Management Service)

# What does this mean? (3)

- ***Accessible*** means EUDAT data are
  - Identifiable and findable (the PID Service)
  - Retrievable efficiently (the Data Staging Service)
  - Governed by suitable access control (the AAI Service)

- ***Re-usable*** means EUDAT data are
  - Findable (the PID Service)
  - Comprehensible (the Joint Metadata Service)
  - Composable and combinable (future workflow and computational services)

**EUDAT**

# What does this mean? (4)

- For both **big communities** and ***"long tail"*** means
  - Stable, web-service APIs for existing tool-stacks to use (the Common Service Layer Interface)
  - Low barriers to use (the Simple Store Service)

- Hence the core EUDAT service cases

- Identifying solutions for these cases *that work with our stakeholder communities' existing solutions* led us to the current CDI architecture
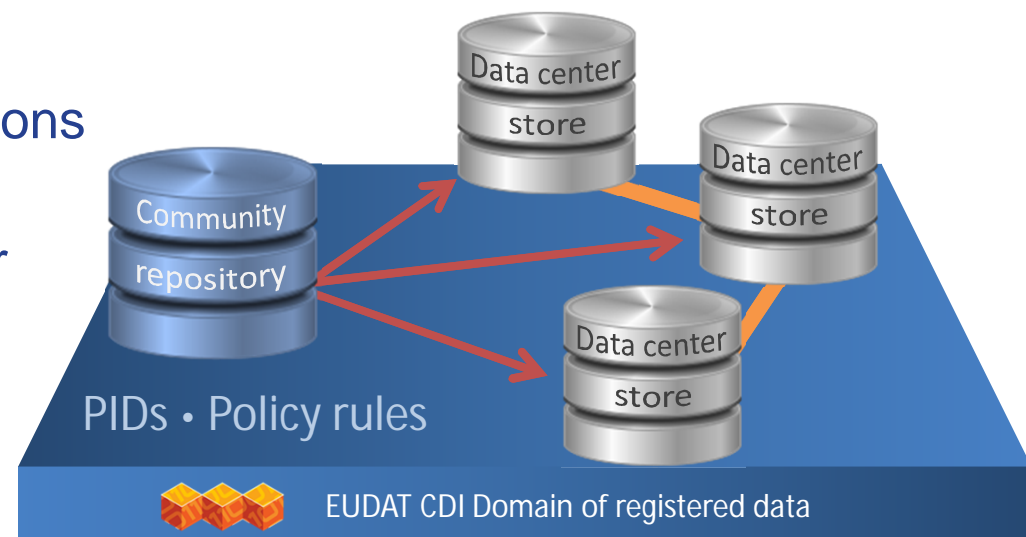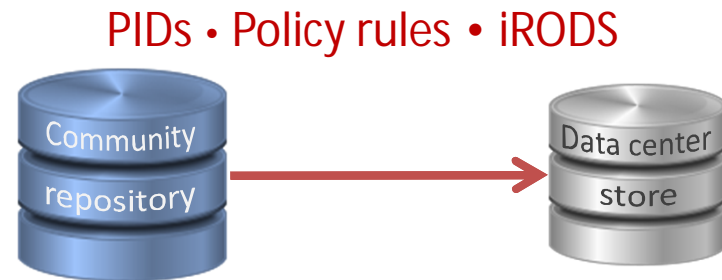
# common services EUDAT is working on

## Metadata Catalogue
Aggregated EUDAT metadata domain.
Data inventory

## Safe Replication
Data curation and access optimization
Various flavors

## Data Staging
Dynamic replication to HPC workspace for processing

## Simple Store
Researcher data store (simple upload, share and access)

## AAI
Network of trust among authentication and authorization actors

## PID
Identity
Integrity
Authenticity
Locations

## services to come

## EUDAT Box
dropbox-like service
easy sharing
local synching

## Dynamic Data
immediate handling

## Semantic Anno
checking & referencing

EUDAT

# Safe Replication Service

- Robust, safe and highly available data replication service for small- and medium- sized repositories
  - To guard against data loss in long-term archiving and preservation
  - To optimize access for user from different regions
  - To bring data closer to powerful computers for compute-intensive analysis

# Safe Replication Service (B2SAFE)

PIDs • Policy rules • iRODS



communities/departments
- do not have IT people
- can't install iRODS
- can't adapt their repository solution
- are partly using ready-made solutions (Fedora, D-SPACE, WikiMedia, etc.)
- don't know how to register PIDs
- etc.

- not as easy as thought

- so what to do?

- created/are creating various flavors
  - FULL SR via iRODS policies    ready
  - Light SR via GridFTP client    ready
  - Lighter SR via HTTP client    in p
  - Packages for
    - Fedora    ready
    - D-SPACE    in p
    - WikiMedia    tbd
    - ?

http://eudat.eu/safe-replication | eudat-safereplication@postit.csc.fi
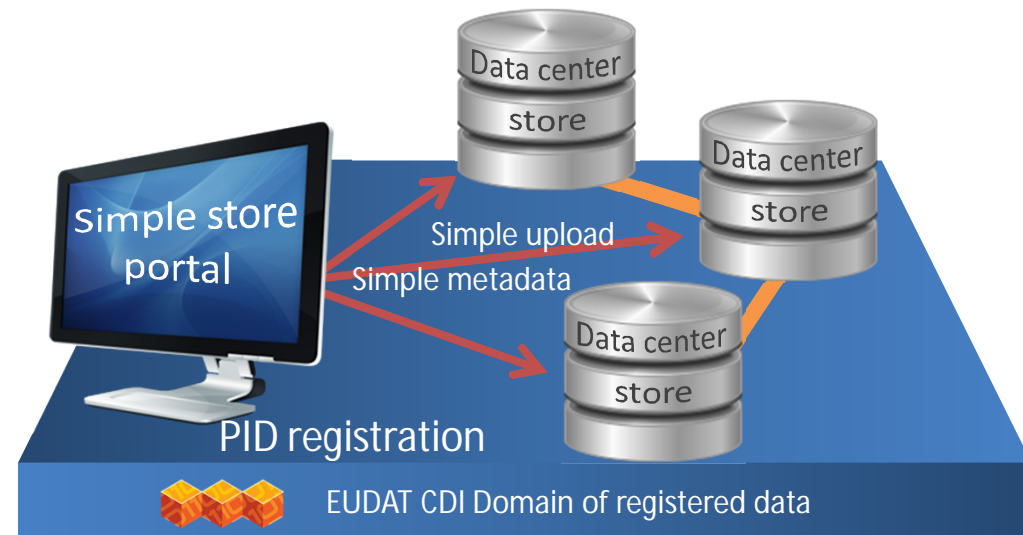
# Data Staging Service

- Support researchers in transferring large data collections from EUDAT storage to HPC facilities

- Reliable, efficient, and easy-to-use tools to manage data transfers

- Provide the means to re-ingest computational results back into the EUDAT infrastructure



- not a simple service!
- politics involved (access to HPC)

# Simple Store Service

- Allow registered users to upload "long tail" data into the EUDAT store

- Enable sharing objects and collections with other researchers

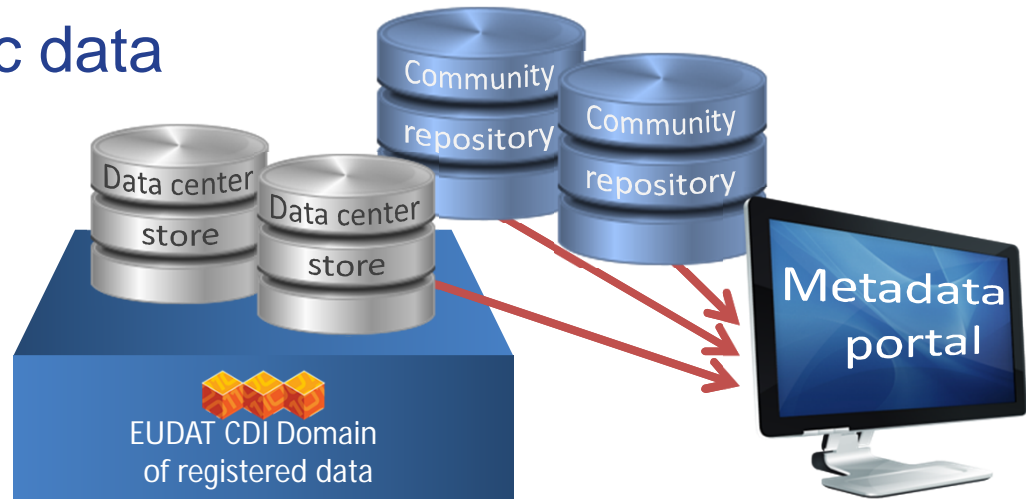- Utilise other EUDAT services to provide reliability



simple store portal

Simple upload
Simple metadata

Data center
store

Data center
store

Data center
store

PID registration

EUDAT CDI Domain of registered data

- much competition

- see it as complementary – finally it is about trust

EUDAT

# Metadata Service

- Easily find collections of scientific data – generated either by various communities or via EUDAT services

- Access those data collections through the given references in the metadata to the relevant data stores

- Europeana of scientific data

- how to offer metadata in a cross-disciplinary space?

- scalability issue?



http://eudat.eu/metadata | eudat-metadata@postit.csc.fi

# EUDAT Box Service

- some similarity to SimpleStore of course

- just similar to Dropbox incl. load balancing and replication
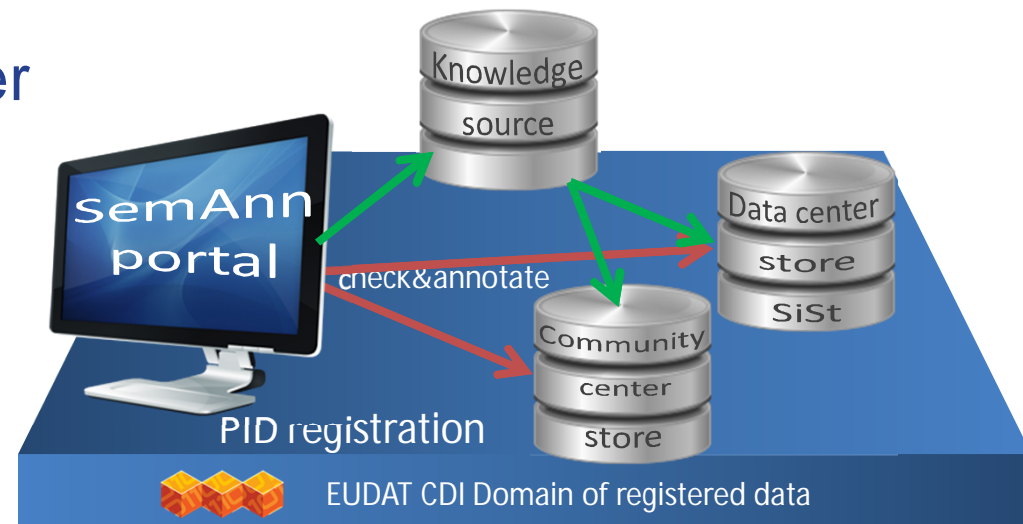
- there is no metadata – just data

- how to integrate into registered domain of data?

- much competition

- see it as complementary – finally it is about trust

# Semantic Annotation Service

- acts as a plugin component to be executed before uploading a resource with tags (crowd sourcing etc.)
- check tags against Knowledge Source & correct/refer/etc.

- could be used as trigger in Simple Store
- plugin available to everyone

- not center dependent

# Service Targeting

- **Replication**: targeted at data managers/archivists/projects/departments without facilities

- **Data Staging**: same plus "easy" access to HPC

- **SimpleStore**: place for individuals/projects/groups to store & exchange data

- **EUBox**: share data via synchronization

- **Metadata**: EUDAT data & everyone interested

- **SemAnn**: individual/projects working with massive amounts of human created data

data stored in domain of registered data is not EUDAT's data!
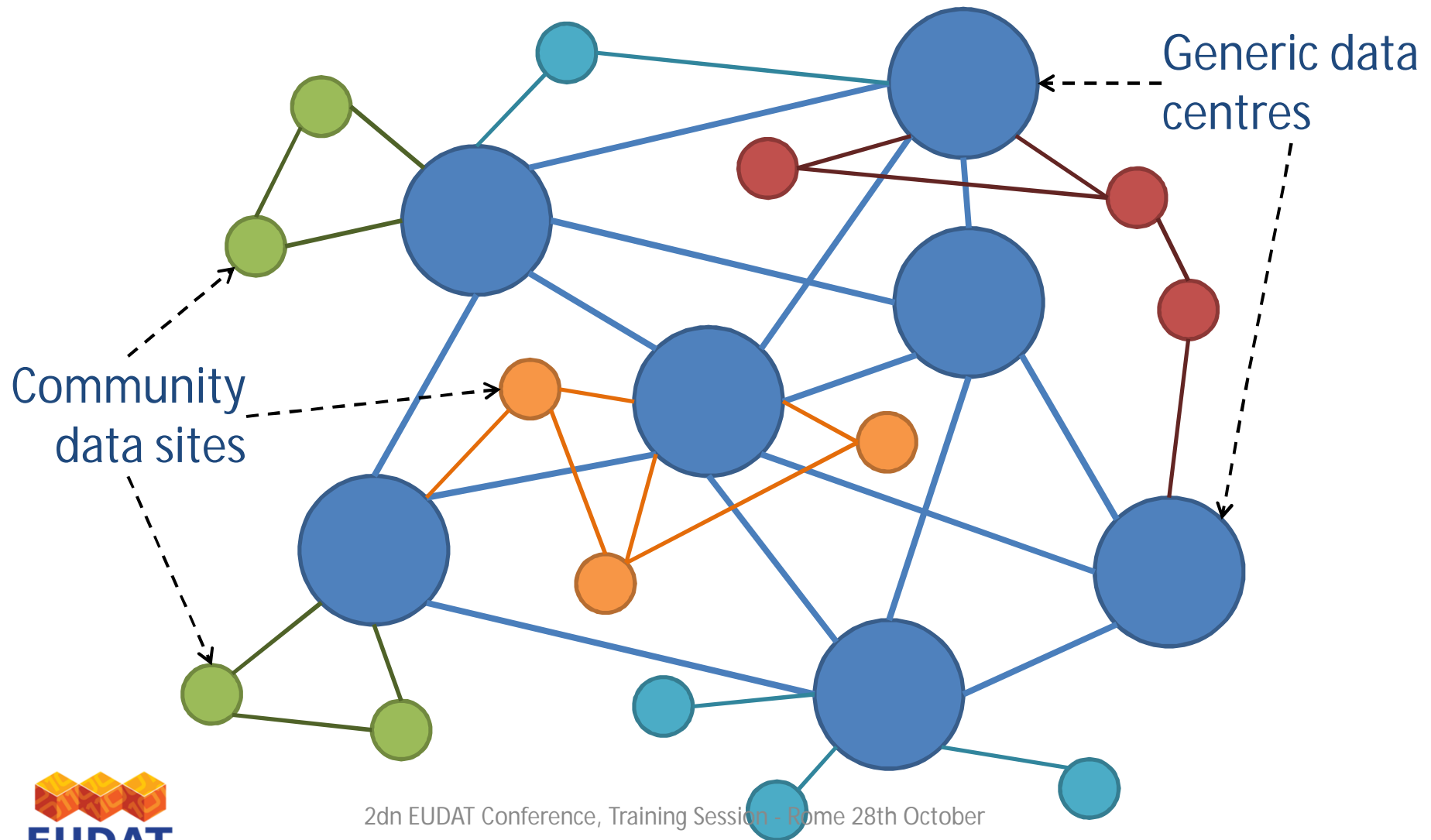
**EUDAT**

# COMMON DATA INFRASTRUCTURE

EUDAT

# The CDI network architecture

- The CDI is a connected network of European research institutions and data centres (collectively *Nodes*) each offering one or more common EUDAT data services to both participating research communities and independent researchers

- Data centre Nodes have lots of connections

- Research community Nodes need only one

- Connections have both technical & policy agreement aspects

**EUDAT**

# The CDI network architecture



Generic data centres

Community data sites
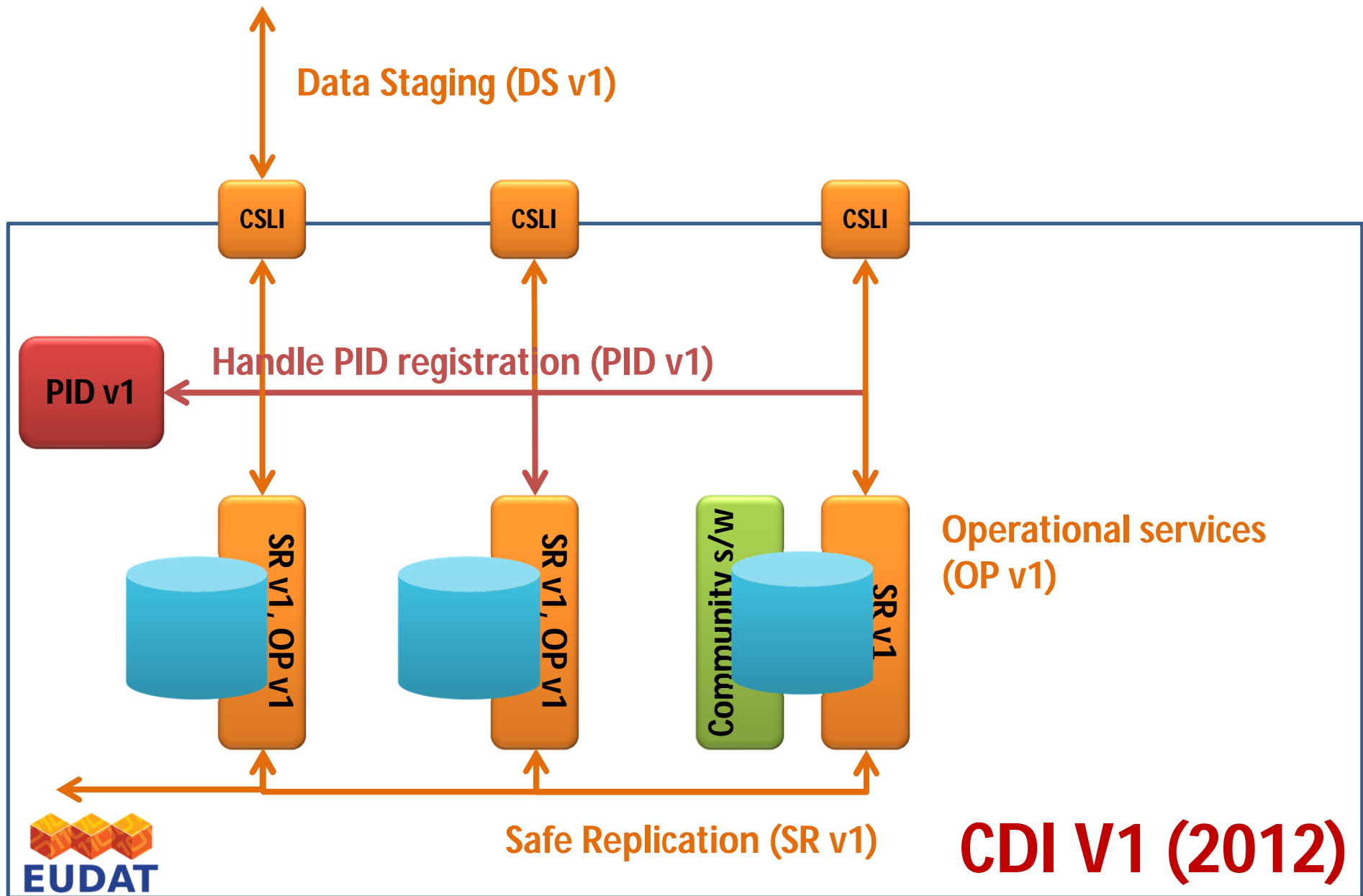
EUDAT

# CDI Node architecture

- Nodes run parts of the CDI Node software suite, depending on which services they want to offer

- All Nodes should offer Safe Replication and PID
  - This is really what being in the CDI is all about

- Others are optional
  - Depends on what a Node's expected user base requires

- (Some data centre Nodes also need to run the Operational Services suite)
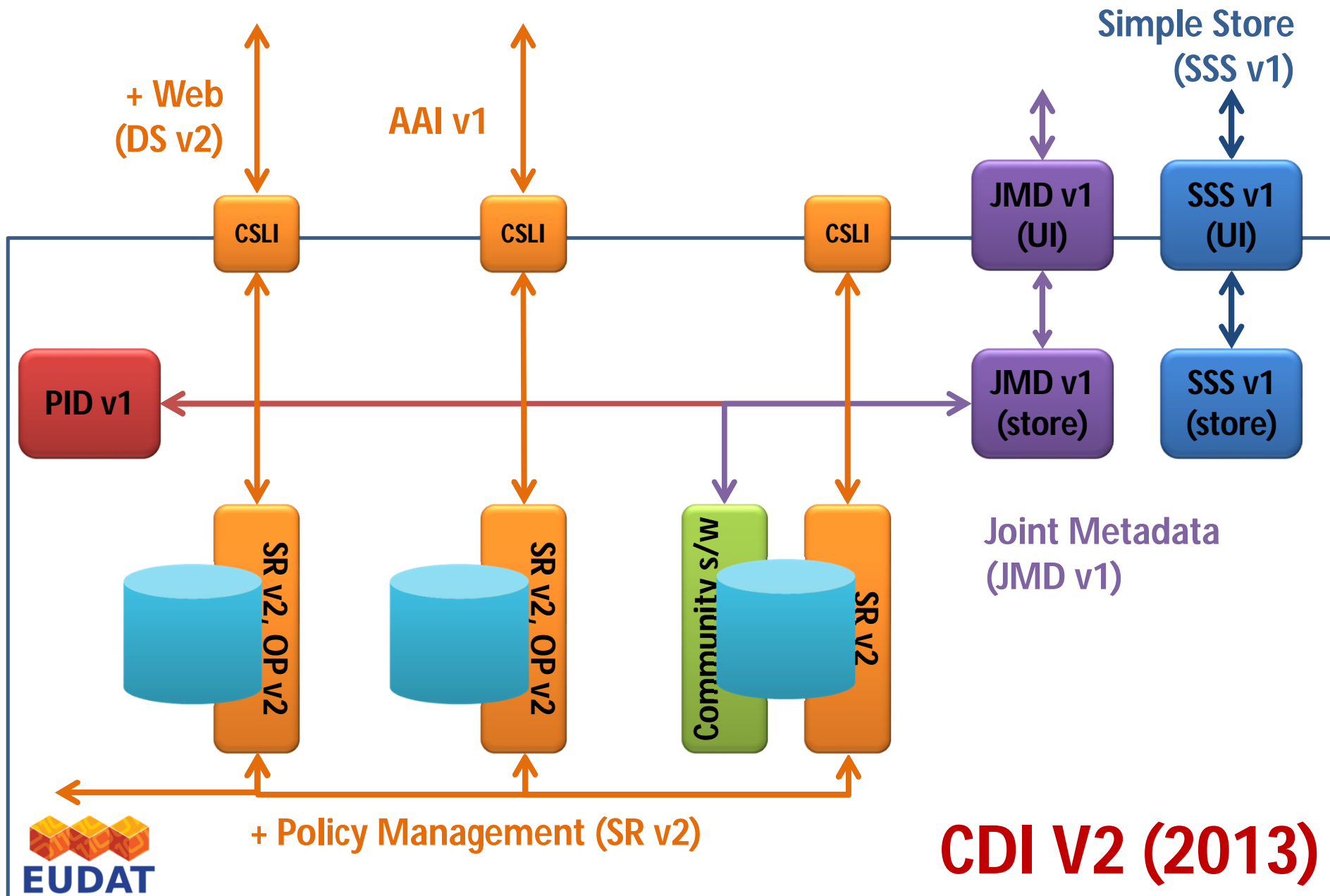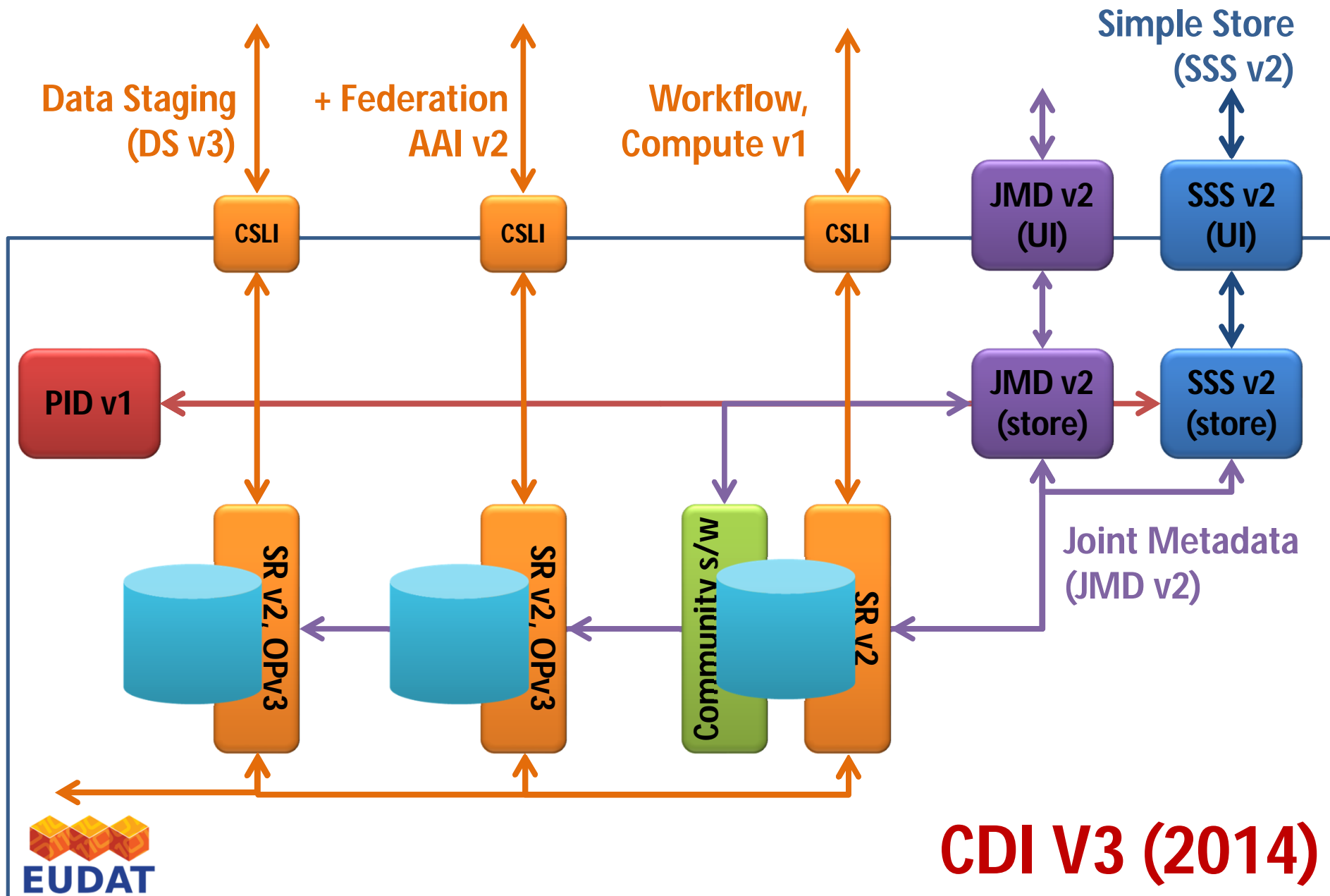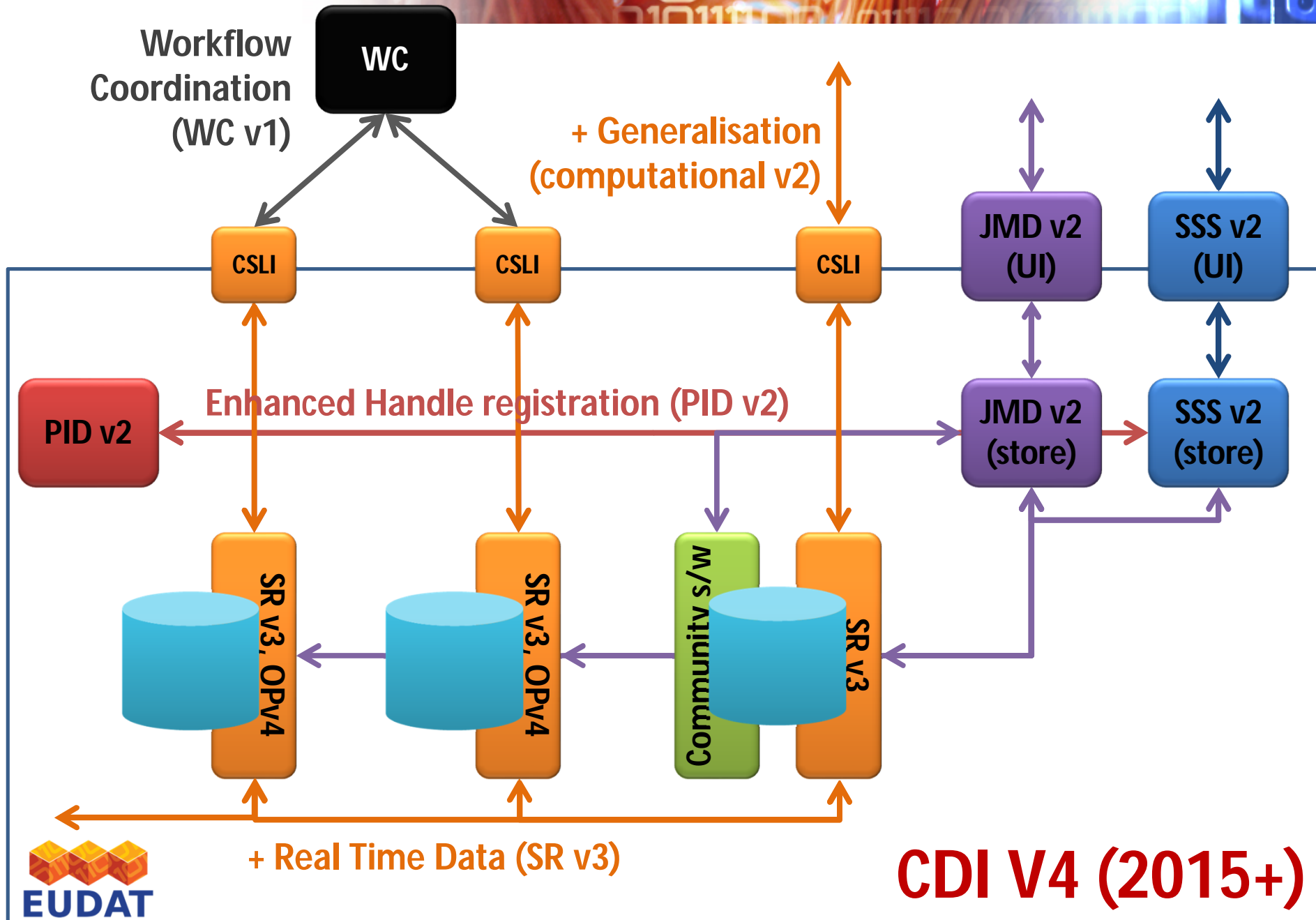
# CDI service development timeline

| Year | EUDAT Services | CDI |
|------|----------------|-----|
| 2012 | Safe Replication v1, Data Staging v1, PID v1, Operational Services v1; AAI (design) | V 1 |
| 2013 | AAI v1, Joint Metadata v1, Simple Store Service v1, Safe Replication v2, Data Staging v2, Operational Services v2; Common Service Layer Interface (design), workflow and computation services (design) | V 2 |
| 2014 | AAI v2, Joint Metadata v2, Simple Store Service v2, Data Staging v3, Operational Services v3, CSLI v1, computational and workflow services v1 | V 3 |
| 2015 + | PID v2, Safe Replication v3, computational and workflow services v2 | V 4 |

EUDAT

Data Staging (DS v1)

CSLI  CSLI  CSLI

Handle PID registration (PID v1)

PID v1

Community s/w

SR v1, OP v1

SR v1, OP v1

SR v1

Operational services (OP v1)

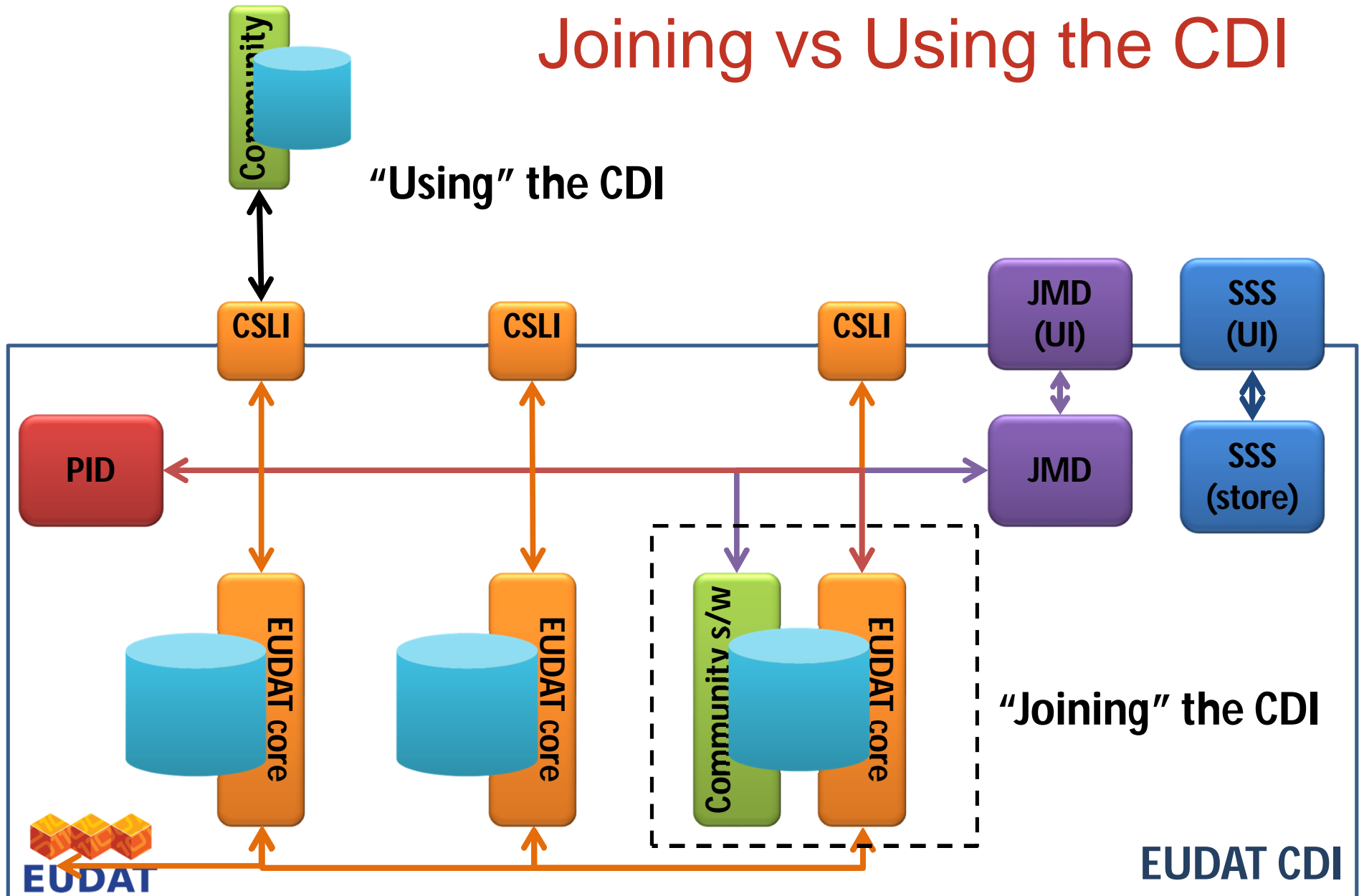Safe Replication (SR v1)

EUDAT

CDI V1 (2012)

# Current Node software suite

- **Safe Replication**
  - iRODS v3.x + EUDAT replication microservices + GSI (X.509) security
- **PID**
  - EPIC/Handle system (external service) + EUDAT EPIC client microservices for iRODS
- **Data Staging**
  - GridFTP iRODS DSI + EUDAT Data Staging script
- **Joint Metadata**
  - CKAN + Apache Solr + OAI-PMH
- **Simple Store**
  - Invenio + EUDAT faceted user interface layer

EUDAT

# Joining vs Using the CDI

# Many thanks for your attention!

## QUESTIONS?