# Enhancing the EUDAT B2SAFE replication software
# Tasks 9.1/9.2
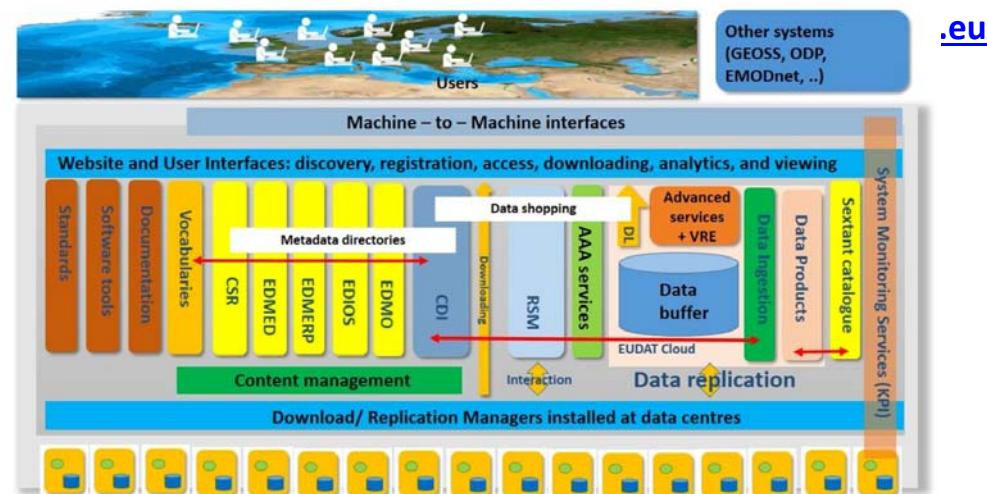
**SeaDataCloud**
**Technical Task Group Meeting**
**Riga, November 29th 2016**

Giuseppe Fiameni
CINECA

www.eudat.eu

# EUDAT & SeaDataCloud

- Strategic cooperation

- Main goal:
  - Strengthening SeaDataNet underlying IT infrastructure
  - Improving access to and discovery of SeaDataNet data & products

- Output:
  - A cloud environment for storing datasets and providing cloud computing capabilities
  - A Virtual Research Environment (VRE) with various analysis and visualisation tools
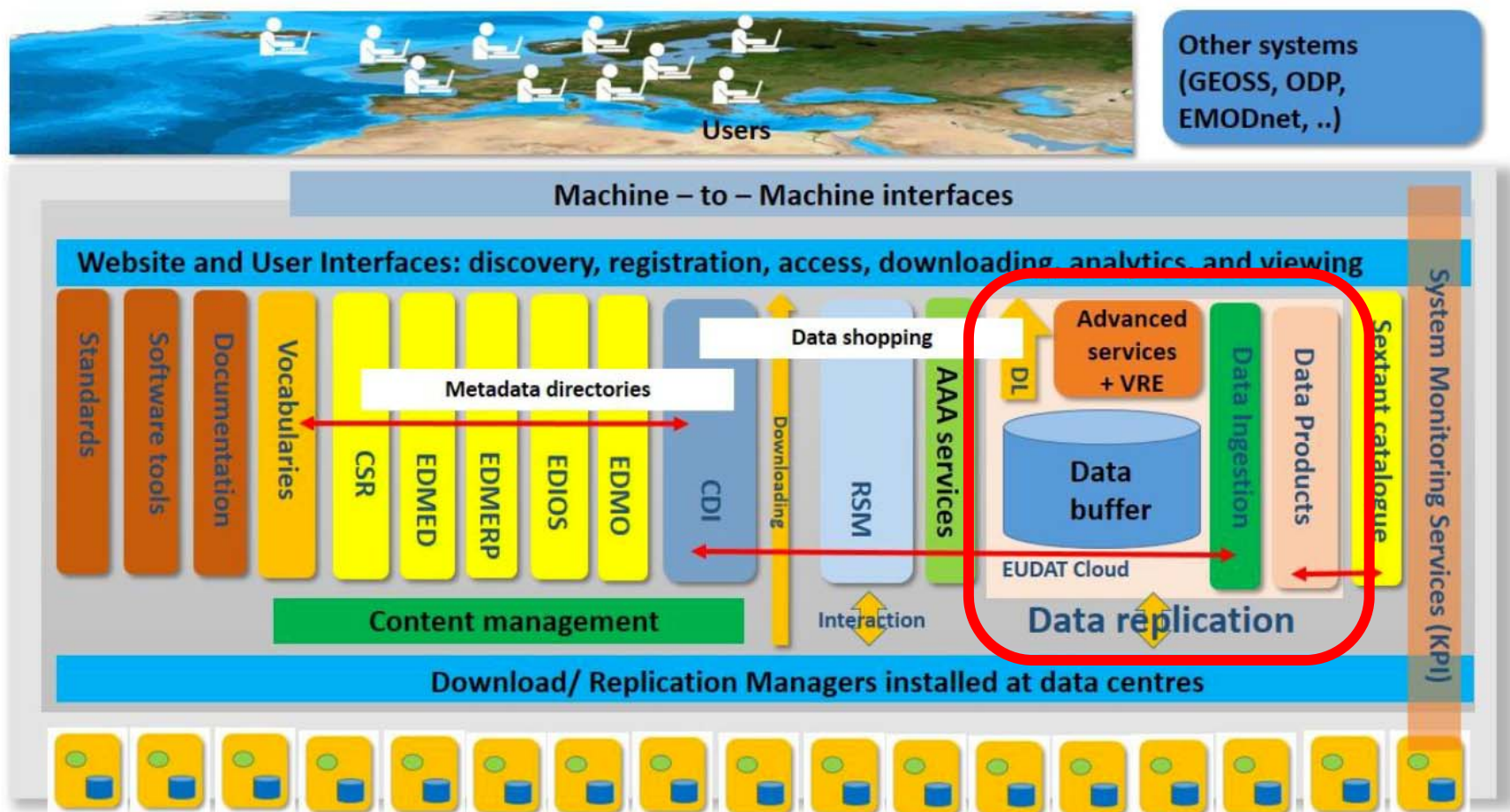


.eu

**Implementation phases**

1) Replicate data from a subset of SeaDataNet centers to the EUDAT centers (technical cache)

2) Define relevant services for SeaDataNet (e.g. data subscription, etc.)

3) Develop a common cloud and computing environment integrated with the SeaDataNet portal and other supported services
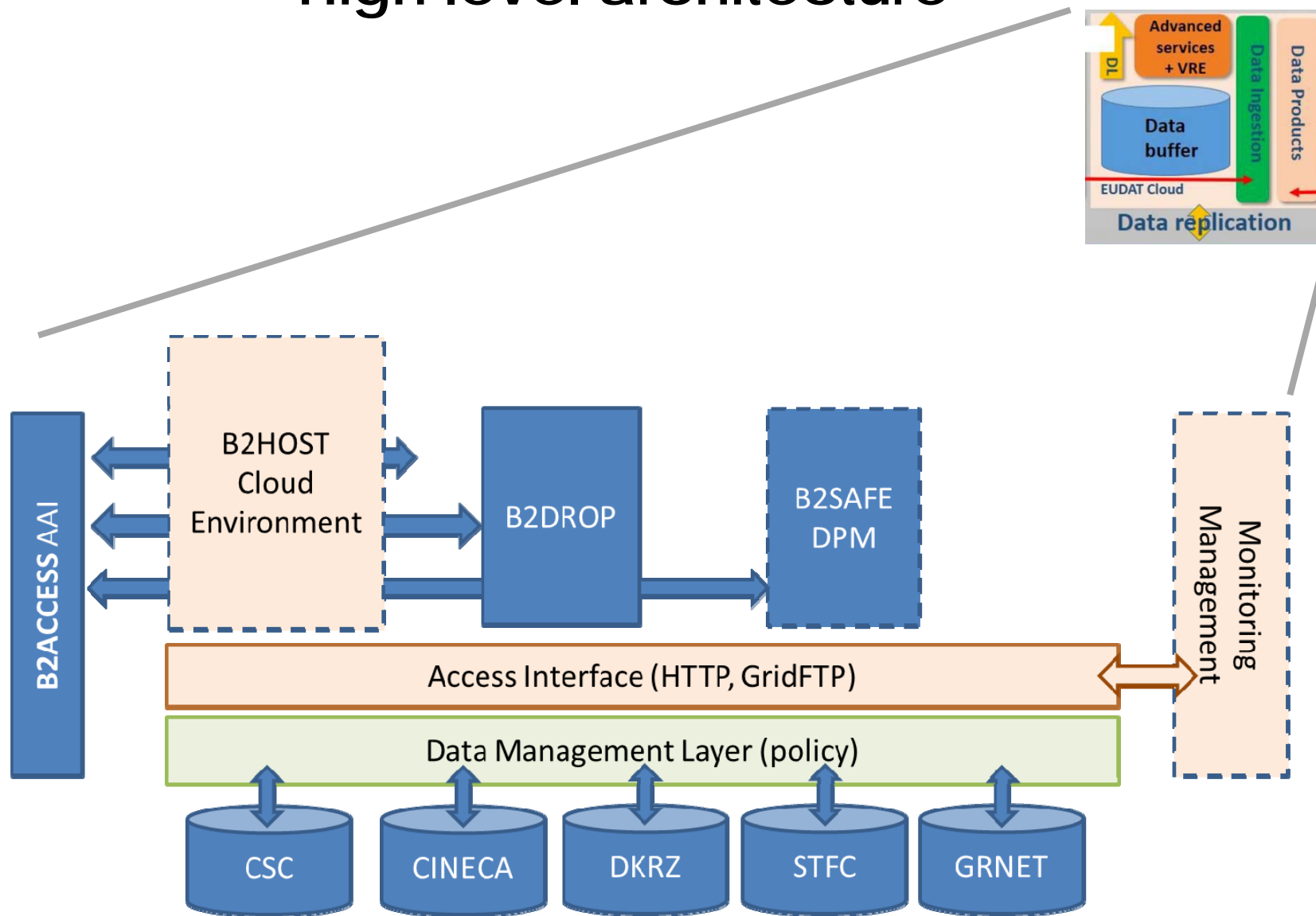
# WP9.1.1: Specification of the SeaDataNet European cloud environment

- Cloud environment for storing data sets and providing cloud computing capabilities
- Replication of data sets between EUDAT centre and a subset of the total number of SeaDataNet data centres.
- An unified data management layer
- Optimize access to data through intelligent replication
- Central registry to track the deposited SeaDataNet data and its provenance.
- Knowing the relatively low volumes generally associated with marine in situ observations, the optimised solution might manage full replication of the datasets in every EUDAT centre.
- Monitor/logging API to support the central monitoring of data replicated and accessed to the EUDAT centres (task 8.5)

# EUDAT & SeaDataCloud
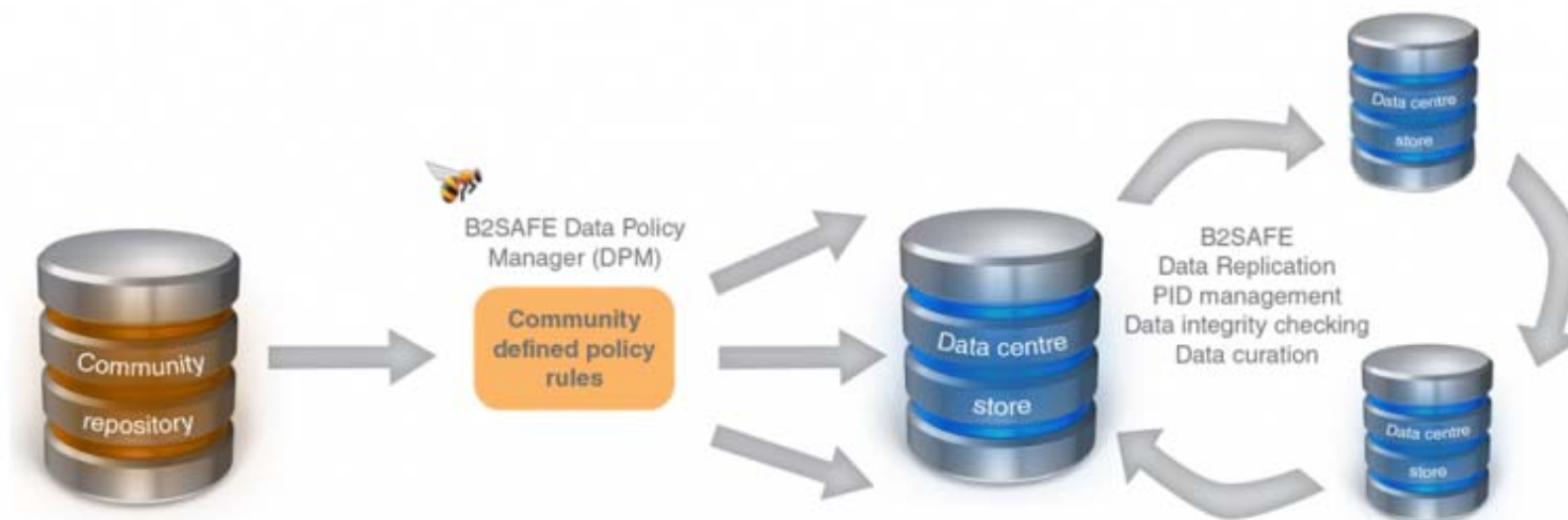
# High level architecture

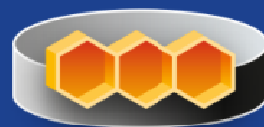# WP9.1.1: Specification of the SeaDataNet European cloud environment

- Cloud environment for storing data sets and providing cloud computing capabilities → B2HOST
- Replication of data sets between EUDAT centre and a subset of the total number of SeaDataNet data centres. → B2SAFE
- An unified data management layer → B2SAFE + API
- Optimize access to data through intelligent replication → Extension of the B2SAFE
- Central registry to track the deposited SeaDataNet data and its provenance. → B2FIND
- Knowing the relatively low volumes generally associated with marine in situ observations, the optimised solution might manage full replication of the datasets in every EUDAT centre. → B2SAFE
- Monitor/logging API to support the central monitoring of data replicated and accessed to the EUDAT centres (task 8.5) → NEW
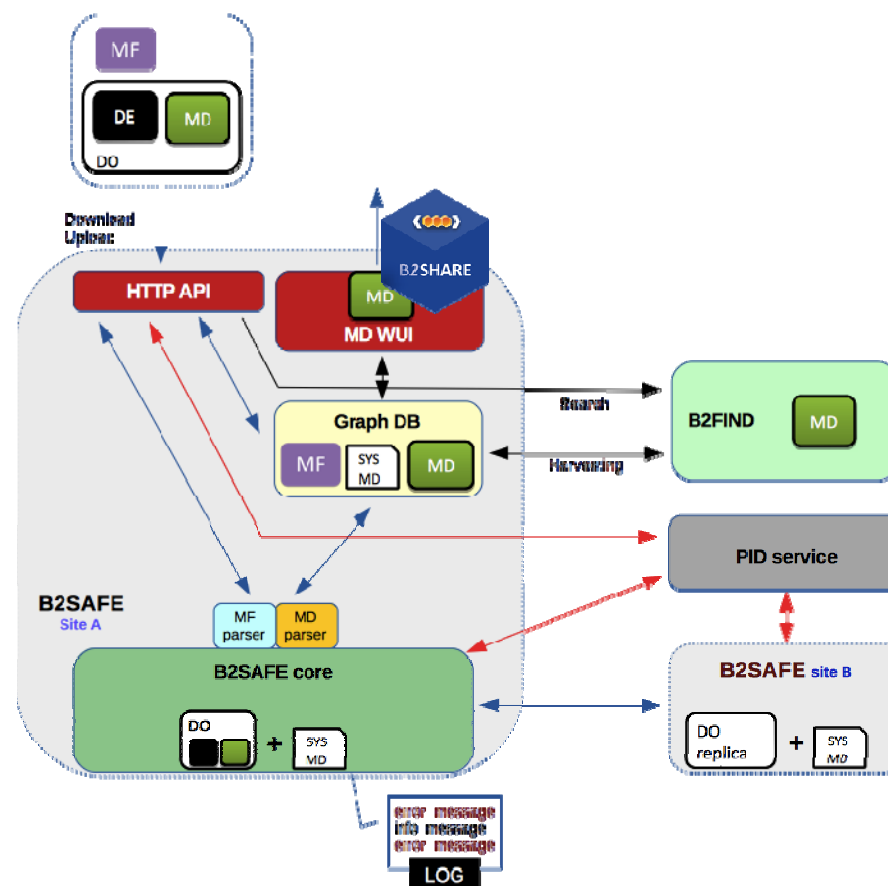
**B2SAFE**
Replicate Research Data Safely

- A robust, safe and highly available data management and replication service allowing community and departmental repositories to replicate and preserve their research data across EUDAT data nodes.
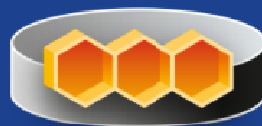
- Support for new PID record structure
- PoC on GraphDB as local MD store
- PoC on Role based authorization
- Implement support of the EUDAT data model in B2SAFE and B2STAGE-HTTP
- Extend MD store with WUI (e.g. B2SHARE)
- Extend metadata support in B2SAFE and BSTAGE-HTTP
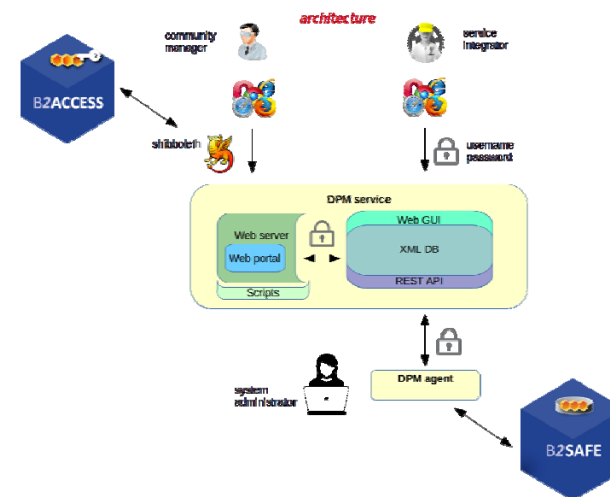
## Metadata Store

- Improved UI to support chaining and filtering of policies
- Improved policy agent to list, update and test policies
- Improved policy management via XML DB
- Policy validation against DPMT constrains
- Support for complex workflows across multiple nodes
- Pilot service available
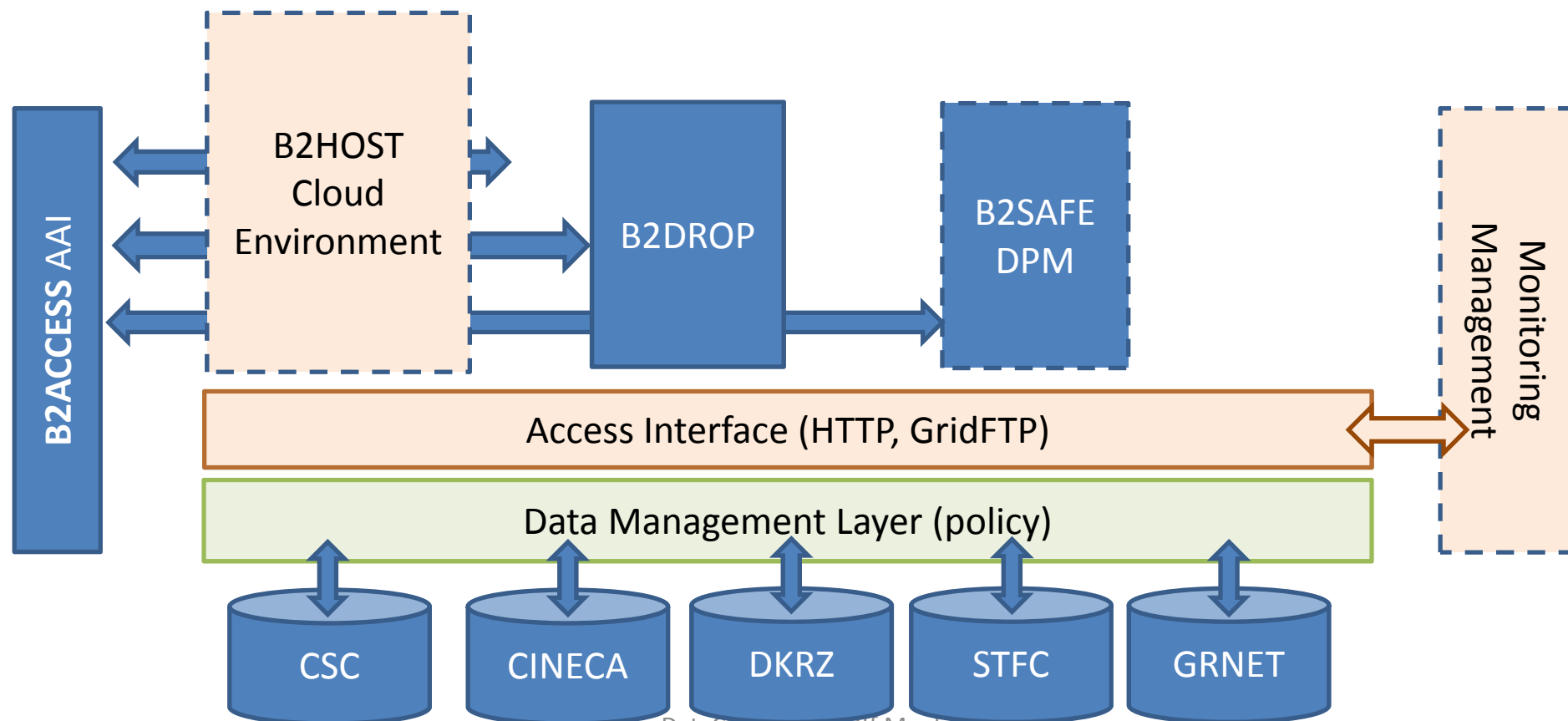- Handover to Operations being planned

## Data Policy Manager

# High level architecture

# WP9.1.2: Enhancing the EUDAT B2SAFE replication software

- The replication scheme will store copies of local SeaDataNet files (for unrestricted and SeaDataNet licensed data) in ODV format in the cloud. It will be triggered by updates in the CDI (Common Data index) directory. This task will also look into versioning of replicated data sets using current developments within EUDAT PID and metadata based versioning policies. The software development work entails enhancing the EUDAT B2SAFE service for registering the replicated SeaDataNet data and the logic for optimising the data location and building the replication and access monitoring API. The versioning of datasets will allow managing in the same repository different processing versions of the same observations as might happen for real time observation results and climatological research quality assessed results.

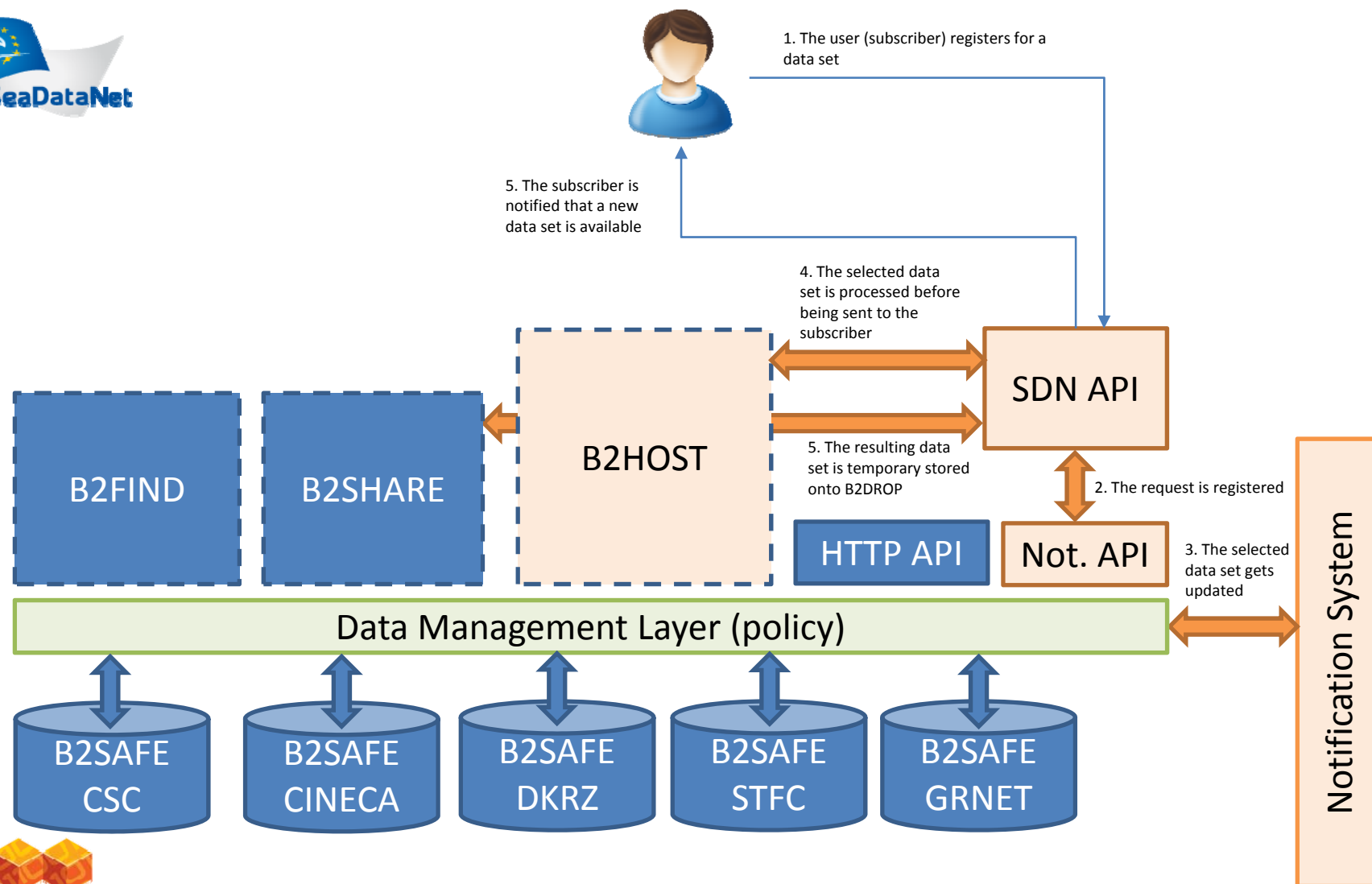# WP9.1.2: Enhancing the EUDAT B2SAFE replication software

- The replication scheme will store copies of local SeaDataNet files (for unrestricted and SeaDataNet licensed data) in ODV format in the cloud. It will be triggered by updates in the CDI (Common Data index) directory. This task will also look into versioning of replicated data sets using current developments within EUDAT PID and metadata based versioning policies. The software development work entails enhancing the EUDAT B2SAFE service for registering the replicated SeaDataNet data and the logic for optimising the data location and building the replication and access monitoring API. The versioning of datasets will allow managing in the same repository different processing versions of the same observations as might happen for real time observation results and climatological research quality assessed results.

# High level architecture

1. The user (subscriber) registers for a data set

5. The subscriber is notified that a new data set is available

4. The selected data set is processed before being sent to the subscriber

**SDN API**

**B2HOST**

5. The resulting data set is temporary stored onto B2DROP

2. The request is registered

**B2FIND**

**B2SHARE**

**HTTP API**

**Not. API**

3. The selected data set gets updated

**Data Management Layer (policy)**

**Notification System**

**B2SAFE CSC**

**B2SAFE CINECA**

**B2SAFE DKRZ**

**B2SAFE STFC**

**B2SAFE GRNET**

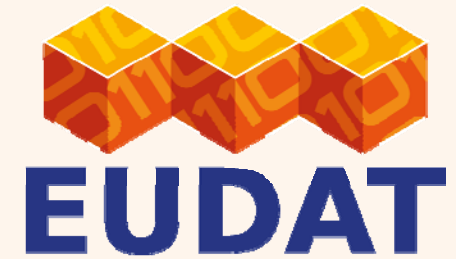# WP9.1.2: Enhancing the EUDAT B2SAFE replication software

- **licensed data** → data management policies may vary largely according to the nature of the license
- **triggering** → real-time triggers can be difficult to manage both sides especially at WAN level
- **versioning of replicated data sets** → multiple versions of the same data set, or data set replicas?
- **optimising the data location** → we need to understand at which granularity level the localization of data would be useful and possible
- **repository** → The Master Copy of the replicated data is usually managed by the community

# How to move forward

- **Refine requirements through piloting a real use case** (*M6 - Specification of the SeaDataNet European cloud environment*)

- Define the services set, e.g. which EUDAT services should be involved.

- Training on EUDAT services

- Define "data projects" for the enabling of the various SDN centres. This is a fundamental step to start with the data replication

  - *Who works with whom? Start small and then grow over the time*

# Interoperability with other services

- A basic VRE infrastructure will be developed by integrating B2DROP, B2SAFE and B2HOST
- B2DROP based workspaces with B2SAFE data interoperability will allow group collaborations
- A B2HOST based hosting environment that allows mounting B2SAFE stored data into the Virtual Machines or containers, allowing execution of processes close to the data
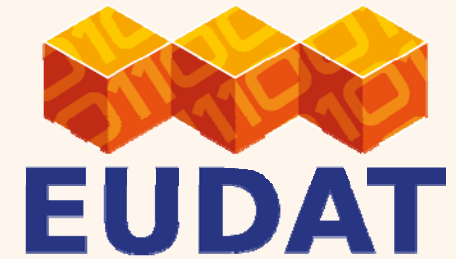
# Thanks for your attention.

SeaDataCloud - Kick-Off Meeting

# Extra Slides

SeaDataCloud - Kick-Off Meeting

# Functional requirements

- Through accessing a central file catalogue. i.e. the B2FIND, an authenticated user should be able to register for one or more "events" that may occur to a given data set and selects which actions should be executed on the data when either of the events occurs. Supported actions may include: *data sub-setting, data delivery, policy activation (through DPM), data curation, etc.*

- When one of the events occurs, all users registered to it receive a notification, i.e. an email message, containing a link to the new data set, stored onto B2DROP, resulting from the processing actions.

- The user should be to cancel submitted registrations.
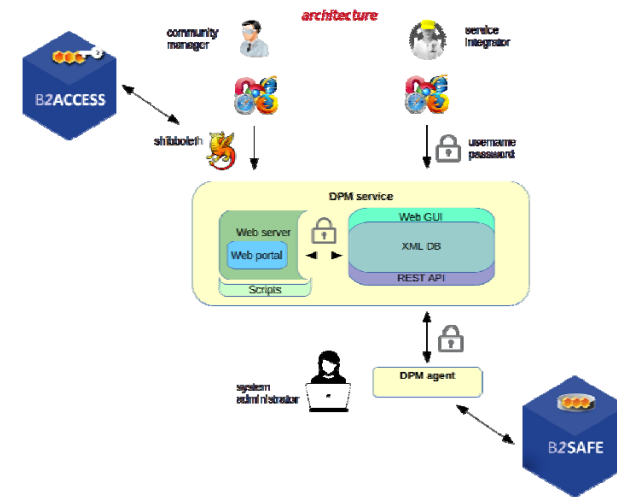
# Non-functional requirements

- Data sets are all open access (?)
- The list of possible executable actions must be validated by the administrator of the system before being offered to the user.
- The central file catalogue, or API, should provide an intuitive and easy-to-use interface to permit the selection of data sets and actions.
- In case of data extraction, the new generated data must contain the information about the provenance of original data sets.

# Some assumptions

- Notifications are sent on daily basis (push model)
- A maximum number of subscriptions should be allowed
- Subscriptions could be simple triplets of data set PID, query and user/email. So user/email wants to know when result of query on data set changes. Queries can be anything
- Subscription are created by external components leveraging EUDAT HTTP API.
- The system would fire up once a day, fetch active subscription from the service, process the queries, upload changed result sets to B2DDROP and report results back to subscription service.
- Based on the detected changes, the subscription service sends emails to users. The emails contain a link to updated result file, plus possibly a short description of the change.
- Importantly, the subscription service should support unsubscribe link. That eliminates the need for user web portal, at least in the first phase.

# Data Policy Manager

- Improved UI to support chaining and filtering of policies
- Improved policy agent to list, update and test policies
- Improved policy management via XML DB
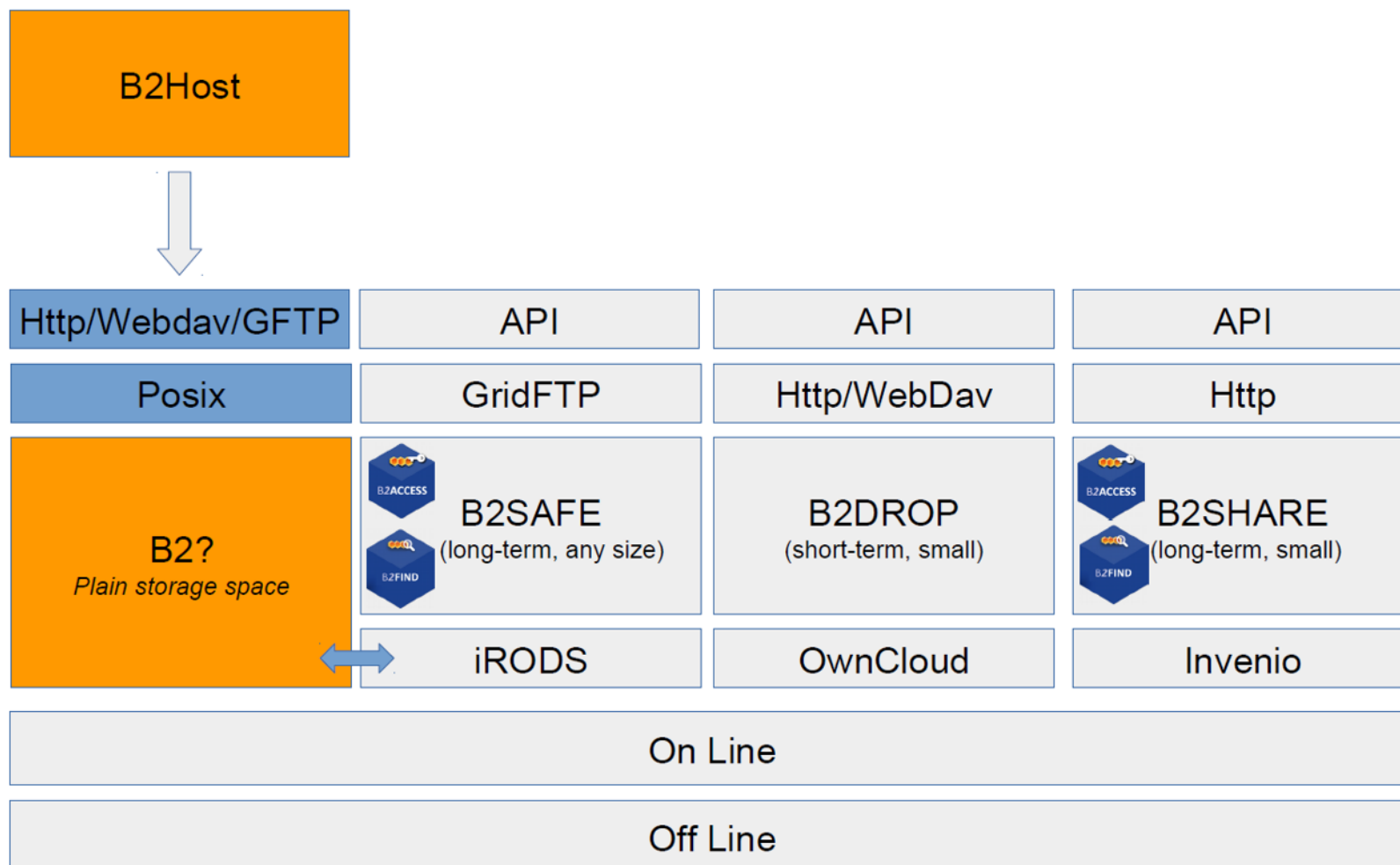- Support for complex workflows across multiple nodes, such as data distribution

# Data Access

# EUDAT in SeaDataCloud

Five selected partners:
- CINECA, Casalecchio di Reno BO, Italy
- DKRZ Hamburg, Germany
- CSC, Helsinki Finland
- GRNET, Athens Greece
- STFC, Rutherford, UK

Mainly to cooperate to the developments of upstream and downstream services:
- WP2: Project Network coordination (all 6PM)
- WP3: Training of data providers (CSC and CINECA 3PM)
- WP5: Expansion of governance of metadata and data content (all 15PM)
- WP6: Core and advanced services (all 35PM)
- WP7: Tuning of requirements and overall integration (CSC 2PM)
- WP8: Governance of standards and development of common services (all 45PM)
- WP9: Development of upstream services (all 86PM)
- WP10: Developments of downstream services (all 65PM)

# EUDAT in SeaDataCloud

- WP2: Project Network coordination (all 6PM)
- WP3: Training of data providers (CSC and CINECA 3PM)
- WP5: Expansion of governance of metadata and data content (all 15PM)
  - WP5.4: Installing, configuring and taking into operation local components for the upgraded CDI data discovery and access service – led by MARIS, IFREMER and EUDAT centres
    - WP5.4.1: Preparing an implementation plan and instructions for the gradual upgrading of the CDI service (MARIS, IFREMER and EUDAT centres)
    - WP5.4.2: Installing and configuring the upgraded components for the population of the CDI Data Discovery and Access service at the >100 data centres (MARIS, IFREMER and EUDAT Centres)

- WP6: Core and advanced services (all 35PM)
- WP7: Tuning of requirements and overall integration (CSC 2PM)
- WP8: Governance of standards and development of common services (all 45PM)
  - WP8.5: Upgrading the operational Monitoring system - led by HCMR, GRNET, STFC and OGS

# EUDAT in SeaDataCloud

- WP9: Development of upstream services (all 86PM)
  - WP9.1: Upgrading the CDI service making use of the cloud – led by MARIS, EUDAT and IFREMER
    - WP9.1.1: Specification of the SeaDataNet European cloud environment (EUDAT)
    - WP9.1.2: Enhancing the EUDAT B2SAFE replication software (EUDAT, MARIS and IFREMER)
    - WP9.1.3: Integrating and adapting the EUDAT B2HOST service (EUDAT and MARIS)
    - WP9.1.4: Configuring the upgraded CDI service (MARIS and EUDAT)
    - WP9.1.5: Deploying, testing and taking into operation the upgraded CDI service (MARIS and EUDAT)
  - WP9.2: Developing and deploying additional services in the cloud for ensuring integrity and conformity of the CDI related cloud data resources - led by IFREMER, MARIS, AWI and EUDAT
  - WP9.4: Developing integrated online services for ingesting autonomous observatory data - led by 52North, IFREMER, CSIC, NERC-BODC, OGS and EUDAT
    - WP9.4.2: Development of the ingestion service (52North, IFREMER, OGS, CSIC and EUDAT)
  - WP9.8: Development of a solution for coordinated distributed DataCite DOI minting service – led by IFREMER and EUDAT

# EUDAT in SeaDataCloud

- WP10: Developments of downstream services (all 65PM)
  - WP10.1: Development of the Virtual Research Environment (VRE) – led by EUDAT, SYKE, MARIS and IFREMER
    - WP10.1.1: Specification of the VRE architecture in the cloud (EUDAT, SYKE, MARIS and IFREMER)
    - WP10.1.2: Development of the VRE (EUDAT)
    - WP10.1.3: Deployment, testing and taking into operation of the VRE (EUDAT, SYKE and IFREMER)
  - WP10.2: Developing advanced services for search, processing, analytics, quality, visualisation – led by IFREMER, MARIS, ETT, Deltares, AWI, ULG, 52North, VLIZ and EUDAT
    - WP10.2.1: enable sub-setting (search on the cloud database at data level: x,y,z,t,observed properties) (IFREMER and EUDAT)
  - WP10.3: Developing MySeaDataCloud – led by EUDAT, MARIS and CNR