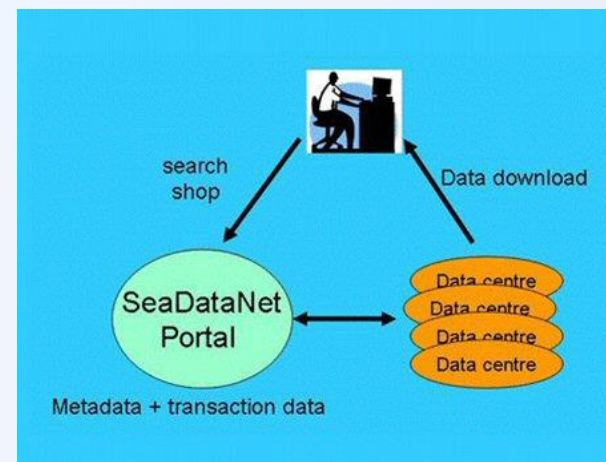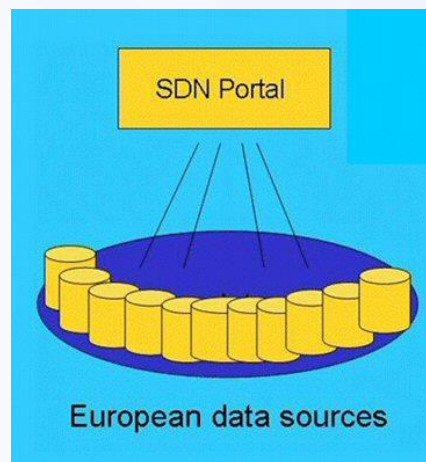Upgrading the Common Data Index (CDI)
Data Discovery and Access service

Dick M.A. Schaap - SeaDataCloud technical coordinator (MARIS, Netherlands)

# CDI Data Discovery and Access service

- one of the core services of the SeaDataNet portal

- It provides a highly detailed insight and unified access to the large volumes of marine and oceanographic data sets managed by the distributed data centres

- fine-grained index (ISO 19115 – ISO 19139) to individual data measurements (such as a CTD cast or moored instrument record)

- Giving access to associated data sets in unified formats (SDN ODV / NetCDF (CF))

- supported by Controlled Vocabularies, and Directories (EDMO, EDMERP, CSR, EDMED)
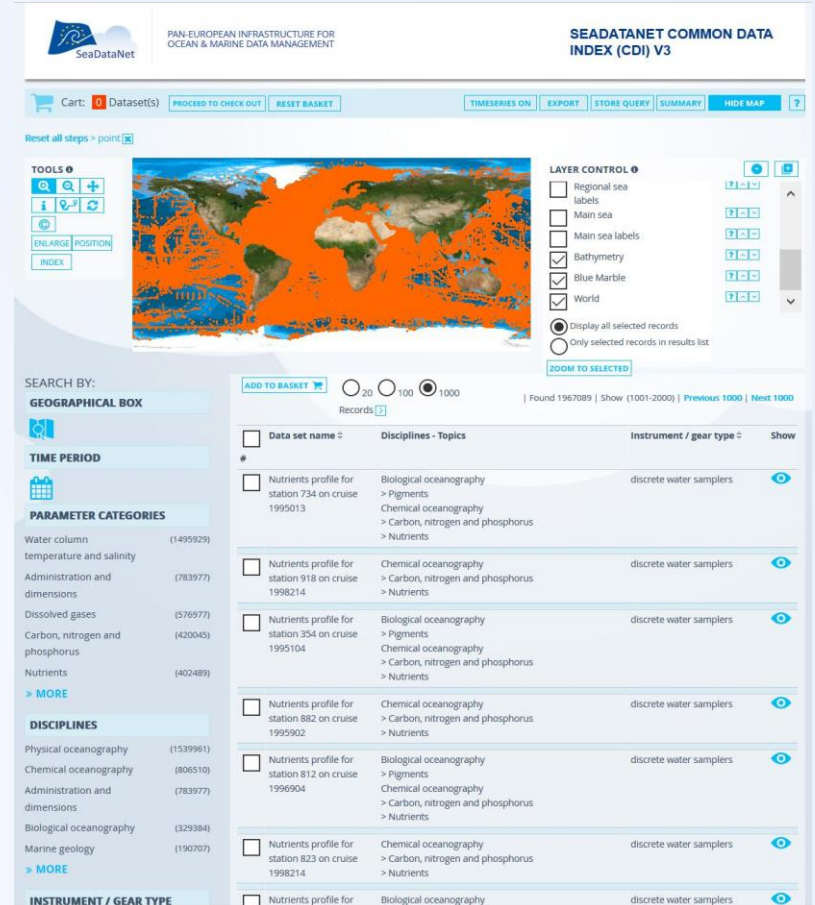
# 2 types of user interfaces



Extended Search                    Quick (facet) Search

SeaDataCloud

GEOSS portal

IODE ODP portal

*Aggregated collection*

Data discovery and access

SeaDataNet

Black Sea portal

Caspian portal

Geo-Seas portal

*Regional subsets*

> 110 data centres

NODCs; HOs; GEOs; BIOs; ICES; PANGAEA

*Thematic portals*

EMODnet
European Marine Observation and Data Network

Bathymetry

Physics

Chemistry

Geology

Biology

≈ 650 European data originators

*Engine behind various regional and thematic portals*

sdn-userdesk@seadatanet.org – w

# Web services

# 110 connected data centres

# CDI data population



CDIs per year

# CDI data population



Number of CDIs per group of parameters

# CDI data population



➢**2 million** CDI entries from **34** countries, **110** data centres and > **650** originators
➢physics, chemistry, geology, geophysics, bathymetry and biology;
➢from **1805 to 2017**; **86%** unrestricted or under SDN License

# discovery and unified data access

**SeaDataNet portal**

**Search and Shop**

**Data Cloud**

Already 110 data centres connected and more underway

**Data centres**

**European data sources**

data centres ⬅ ≈ 650 originators

**Metadata**

**+ transaction data**

# Issues with current CDI service

- **performance for users**: CDI data access service interacts with the distributed data collections and data bases at the connected data centres.
  - user can submit a shopping basket with requests for data from multiple data centres.
  - user must await the automatic data preparation by each of these data centres
  - user must download resulting data sets through the RSM as packages directly from each data centre, which implicates multiple download transactions
- **performance for users**: data centres are not always online, operational and have different machine capacities which might give extra delays
- **quality issues**: concerning formats of data files (ODV + NetCDF) and their consistency with CDI metadata.
- **installation and configuration** of the Download Manager software can be challenging due to different configurations, firewalls etc., which in practice results in having different versions installed

# Upgrading CDI service using the cloud

- Configure and maintain a cloud environment as a 'cache' to host copies of all data resources (from the distributed data centres)

- Exchange by dynamic replication from the individual data centres, following their updating of the CDI catalogue service

- In the cloud buffer new functions:

  – checking possible duplicates

  – Checking overall quality of formats

  – Checking integrity of data files and metadata relations.

  – Results of checks reported back to data centres for amendments of their submissions

- Develop a Virtual Research Environment (VRE) to facilitate collaborative and individual research by users

- Provide customised services (MySeaDataCloud) to let users have search profile, receive alerts on new available data, ingest and manage their own datasets

# Upgrading CDI service using the cloud



**SeaDataCloud**
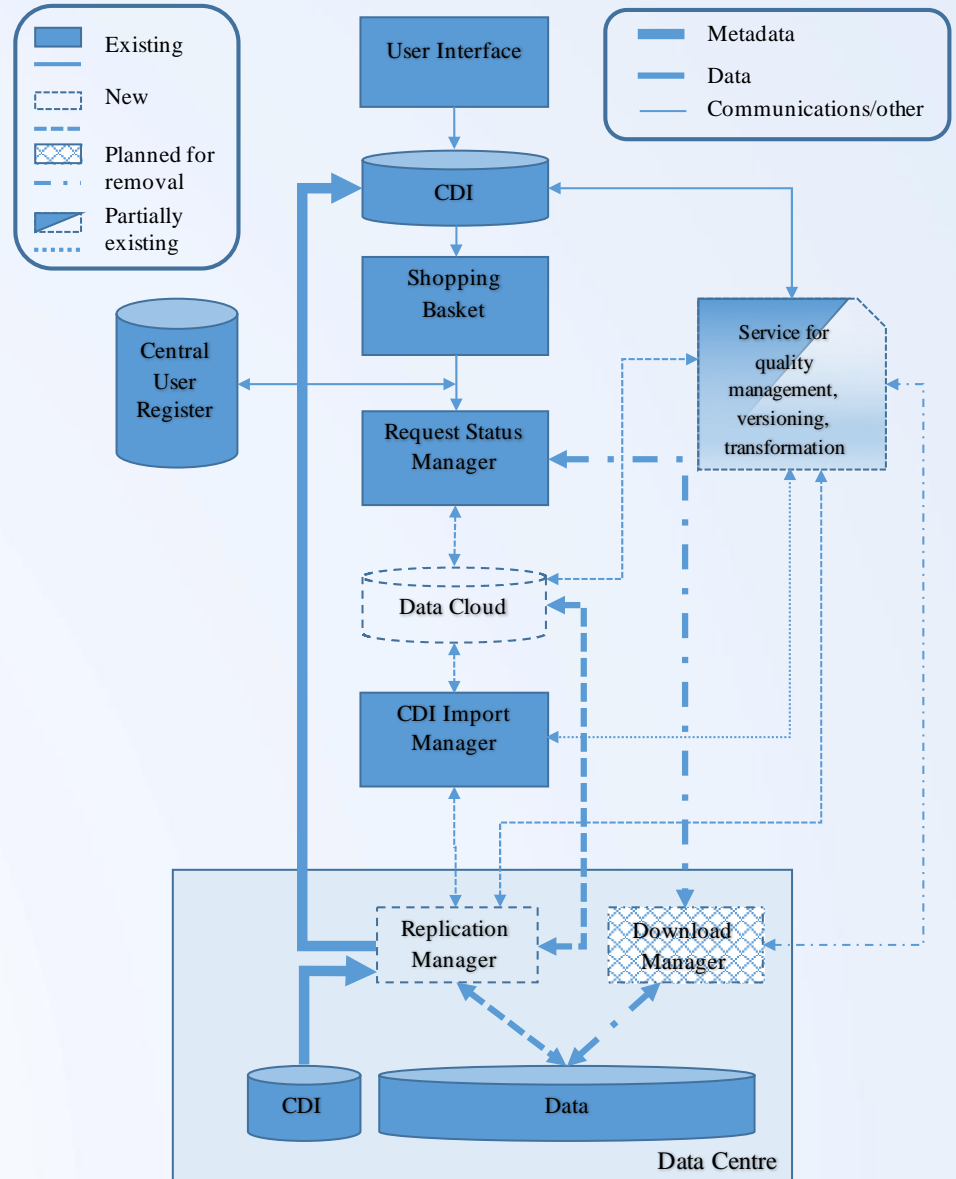
New modern interfaces
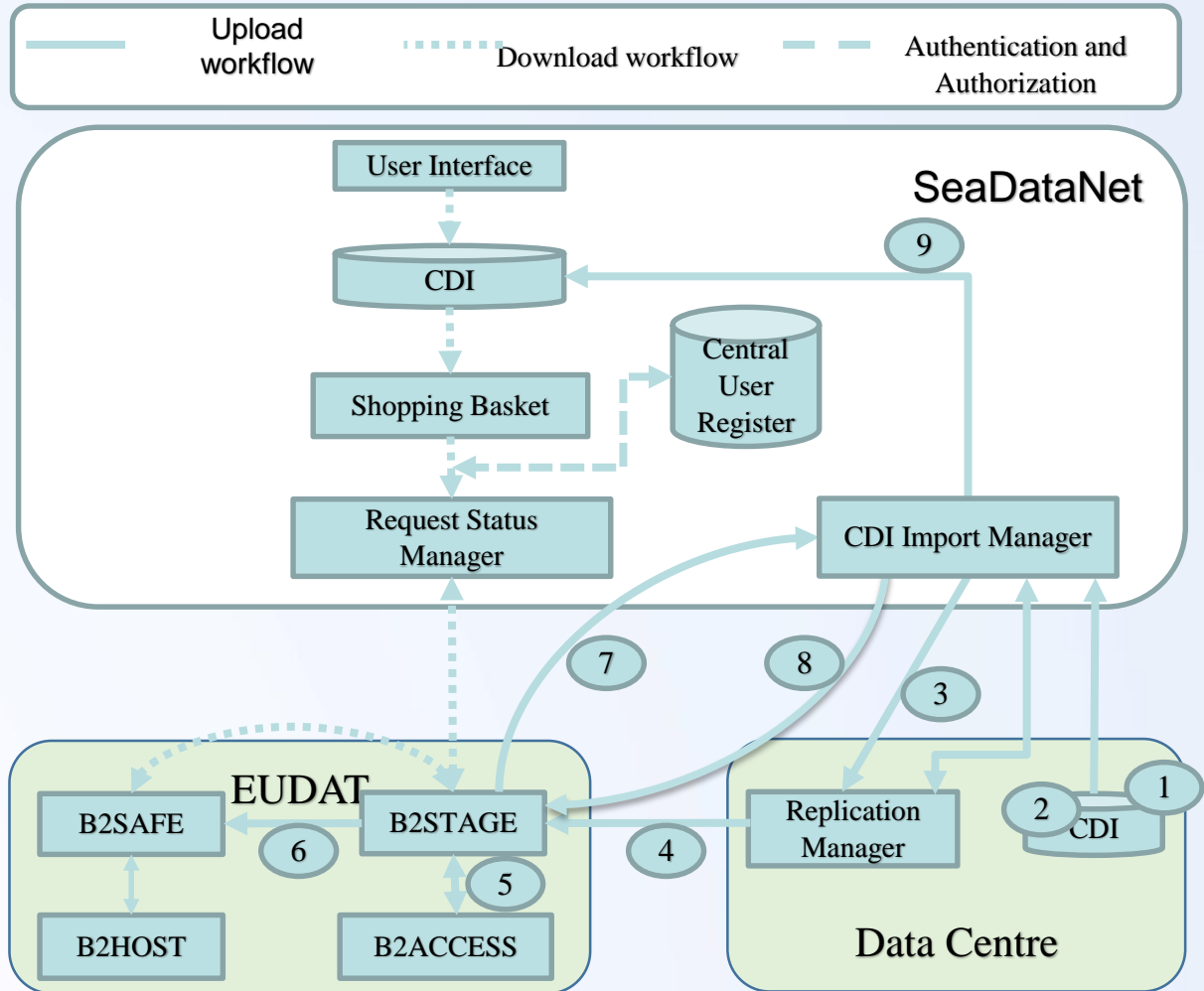
Transformation services

Buffer quality control

Data caching

Import manager

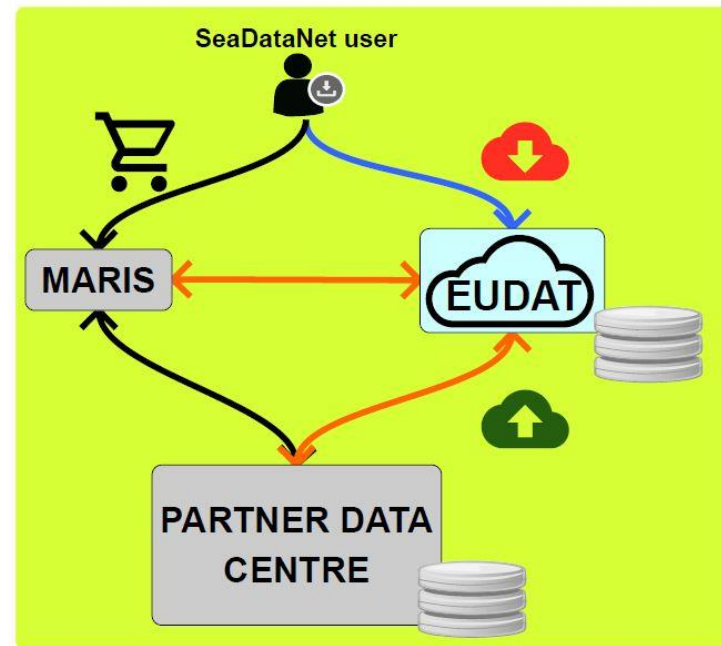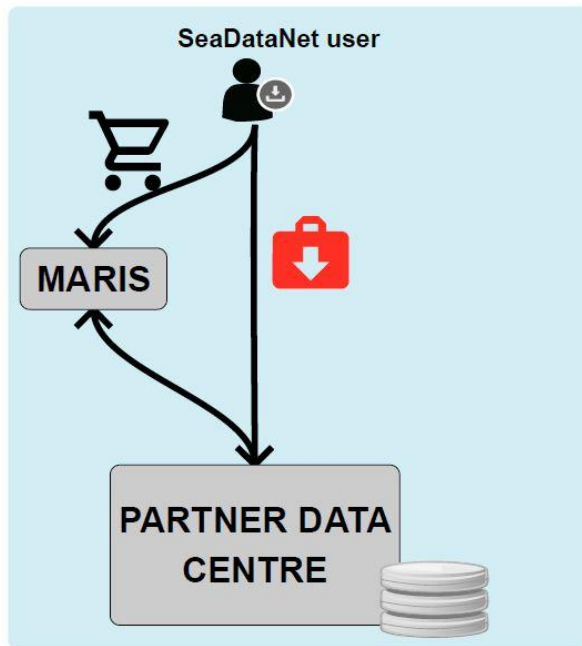Replication Managers at Data nodes

# New service components - workflow

# Replication Manager



3 parts instead of 2:    **EUDAT** is a new element in the workflow

MARIS + DATA CENTRE    →    MARIS + DATA CENTRE + **EUDAT**

# Architecture



old DM -> **new RM**:

- DM_Servlet (tomcat app)
    -> RM API ( tomcat REST service)
    -> RM CMS

- DM Batch -> RM Data

- DM Checker -> RM Checker

- DM ToolsBatch -> RM Tools
  (cleaner is deleted, update vocabs stays)

- -> RM CDI

These are not batches any more, but java libraries.

Modus 2 database is optionnal.
Data centres with modus 2 database can use the same server,
even same database, for modus2 and coupling,
but this is not mandatory.

# EUDAT Cloud

- Import area distributed across the five EUDAT partners

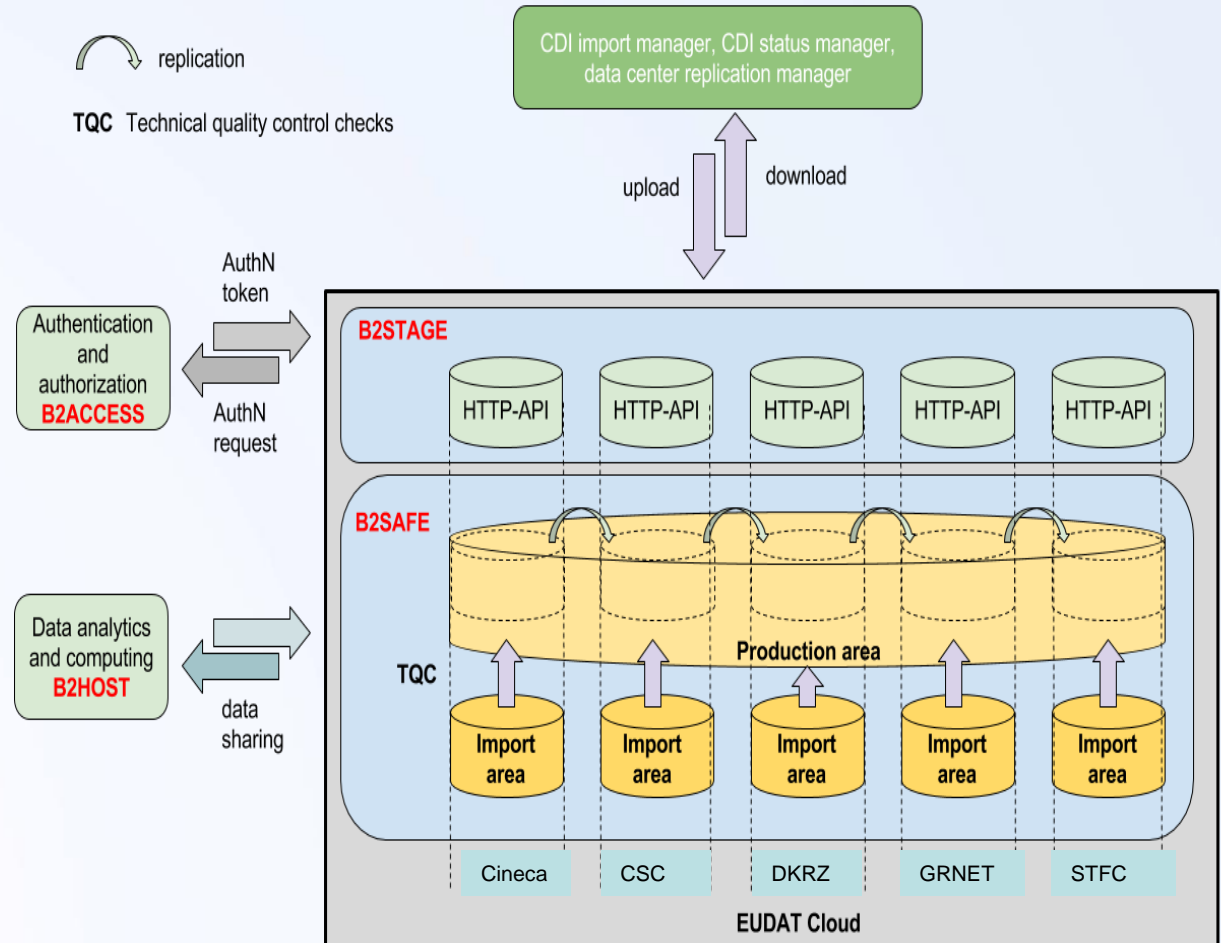- Production area replicated across the five EUDAT partners

# Benefits for CDI service and its users

- Cloud buffer in combination with the CDI service will
  - speed up the performance,
  - expand discovery and ease of use of the data access and downloading
  - provide users with one integrated download package instead of multiple packages from multiple data centres.
- Overall quality and coherence (data – metadata) will improve
- Data replication will be triggered per data centre by CDI updates. The replication module might have less complexity than the present Download Manager module
- A system of versioning will be introduced which is required in the context of the MSFD for facilitating repeated analysis of environmental assessments after many years, and for scientific papers.