



SeaDataCloud

Upgrading the CDI Data Discovery and
Access service, adopting the cloud

Dick M.A. Schaap, MARIS, SDC Technical Coordinator

SDC Review, Brussels – Belgium, 6 December 2018
sdn-userdesk@seadatanet.org – www.seadatanet.org

CDI service for discovery and unified data access

SeaDataNet portal



Already 115 data centres connected and more underway

European data sources
data centres > 650 originators

**Search
and
Shop**



**Metadata
+ transaction data**



**Data
download**



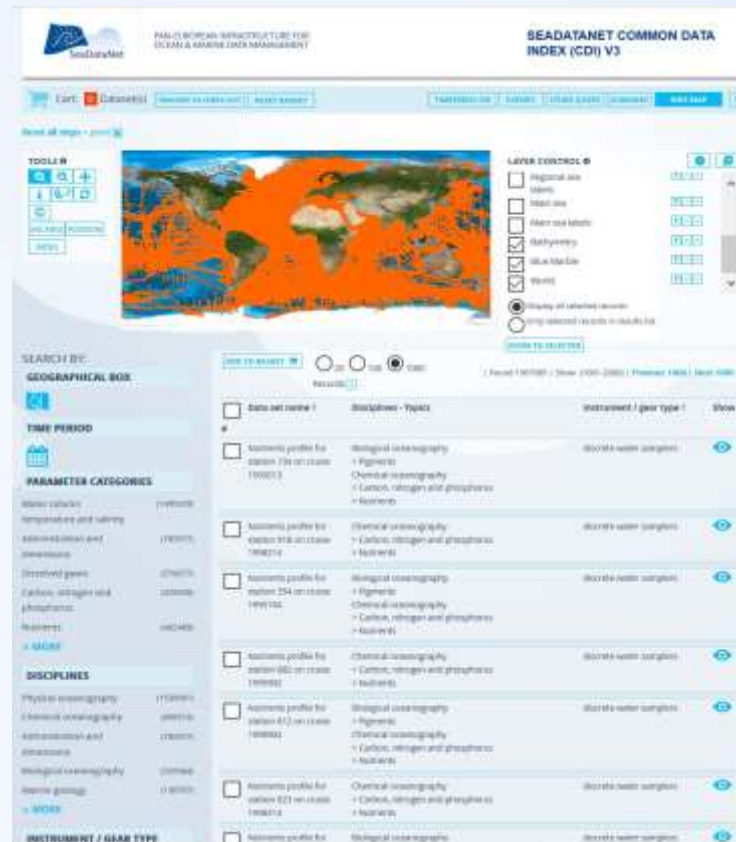
Current CDI user interfaces



The Extended Search interface features a world map at the top with a 'LAYER CONTROL' panel on the right. Below the map is a 'SEARCH' section with a 'Find search' input field and a 'SEARCH' button. The search criteria are organized into several sections:

- Disciplines - Topics:** A dropdown menu showing 'All' and 'Administration and dimensions'.
- Discovery parameters:** A dropdown menu showing 'All' and 'Acoustic backscatter in the water column'.
- Characteristics:** Fields for 'Character name', 'Instrument type', 'Platform type', 'Measuring unit type', 'Originator', 'CDI partner', 'Country', and 'Access restriction'.
- Water depth (m):** A range selector from '0' to '1000'.
- Temporal resolution:** A dropdown menu showing 'All'.
- Duration:** A range selector from '0' to '1000'.

Extended Search

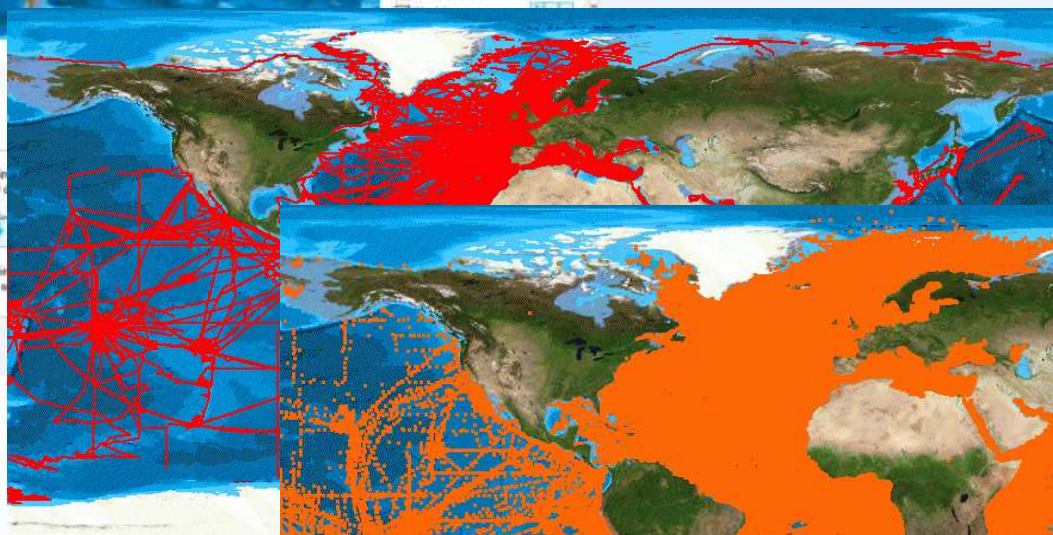
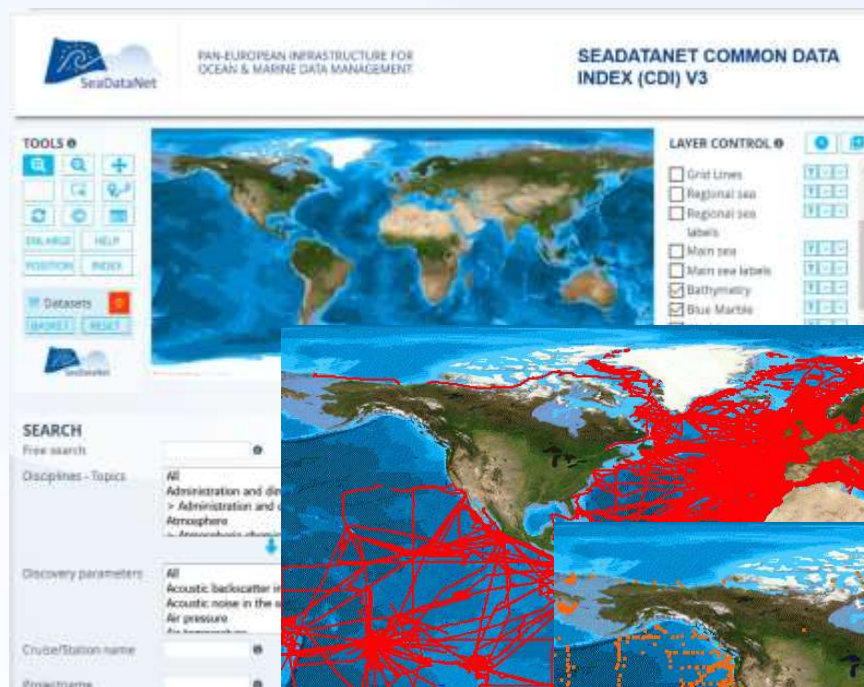


The Quick (facet) Search interface features a world map at the top with a 'LAYER CONTROL' panel on the right. Below the map is a 'SEARCH BY:' section with a 'Find search' input field and a 'SEARCH' button. The search criteria are organized into several sections:

- Geographical Box:** A dropdown menu showing 'All' and 'World'.
- Time Period:** A dropdown menu showing 'All' and '1980-1999'.
- Parameter Categories:** A list of categories including 'Acoustic backscatter in the water column', 'Chemical oceanography', 'Biological oceanography', 'Physical oceanography', 'Atmospheric and oceanic', 'Marine geology', and 'Marine biology'.
- Disciplines:** A list of disciplines including 'Acoustic backscatter in the water column', 'Chemical oceanography', 'Biological oceanography', 'Physical oceanography', 'Atmospheric and oceanic', 'Marine geology', and 'Marine biology'.
- Instrument / Gear Type:** A dropdown menu showing 'All' and 'Acoustic backscatter in the water column'.

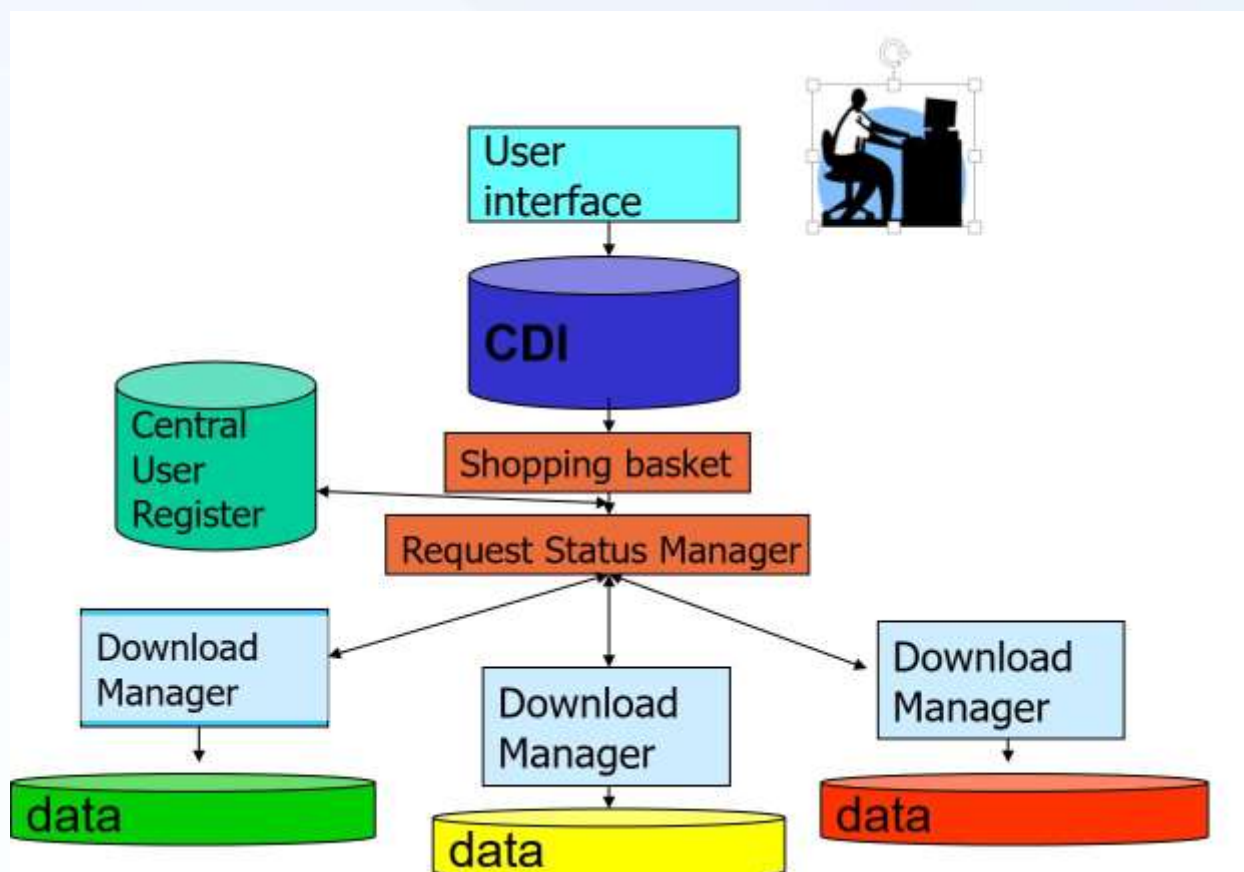
Quick (facet) Search

CDI service with global coverage



> 2.15 Million CDI entries for physics, chemistry, biology, geology and geophysics

Current CDI service architecture



Issues with current CDI service

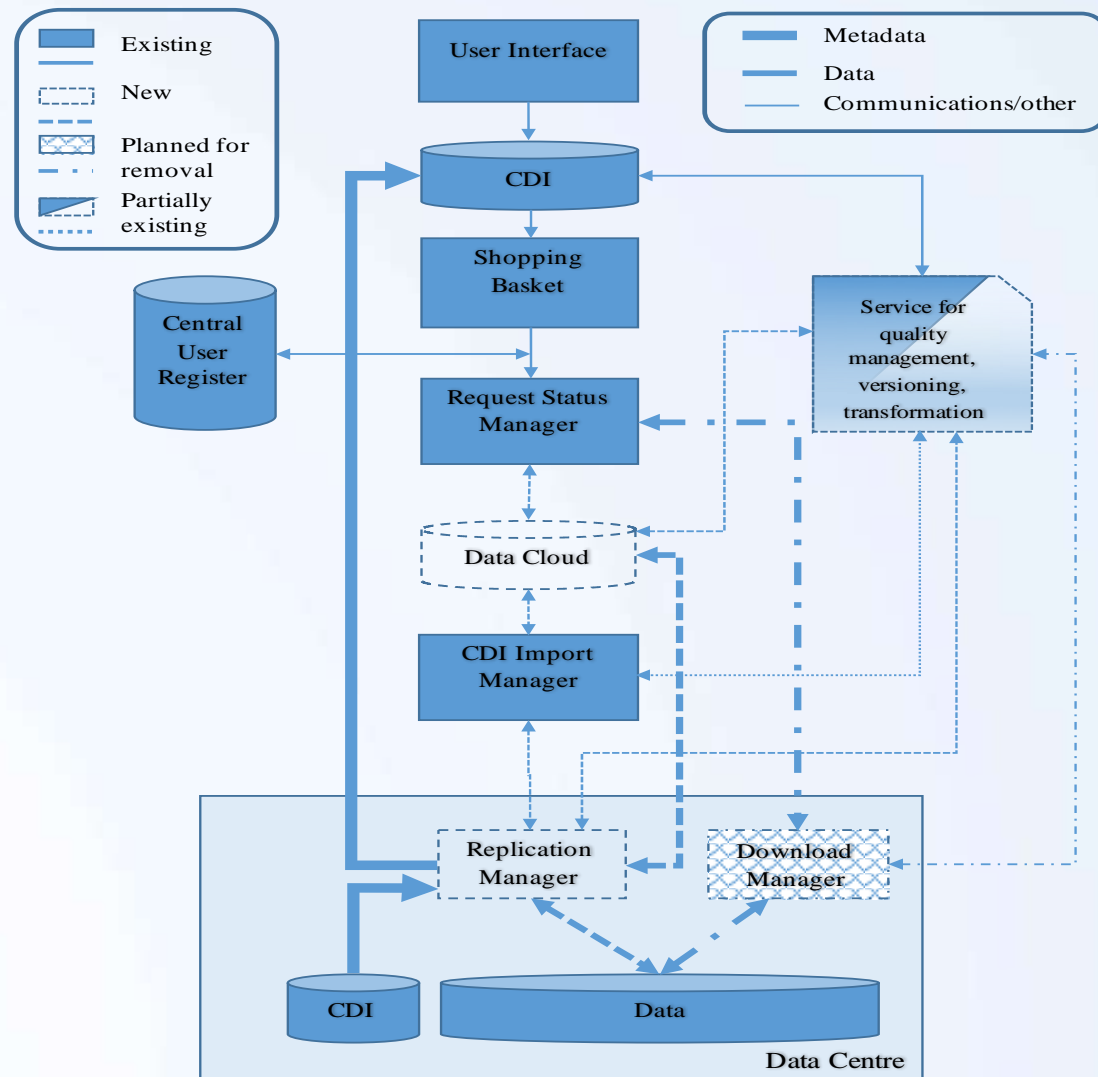
- **performance for users:** CDI data access service interacts with the distributed data collections and databases at the connected data centres.
 - user can submit a shopping basket with requests for data from multiple data centres.
 - user must await the automatic data preparation by each of these 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



Principles for upgrading the CDI service using the cloud

- To configure and maintain a **CLOUD** environment with High Performance Computing (HPC) facilities to host **copies of unrestricted data resources**
- Exchange by dynamic **replication** from the individual data centres, following their updating of the CDI catalogue service
- In the cloud buffer:
 - checking overall quality of metadata and data, as extra check on top of local QA-QC by data centres
 - checking integrity of data files and metadata relations.
 - results of checks to be reported back to data centres for amendments of their submissions and/or local configurations for mapping data and metadata.
- Include transformation services for converting data sets to SeaDataNet ODV and NetCDF formats and relevant INSPIRE data models.
- Introduce versioning of metadata and data as part of provenance

New CDI service architecture



New CDI service components

- **Local software tools** at data centres to prepare ingestions
- **Replication Manager (RM)** at data centres for exchanging to Import Manager and EUDAT cloud
- **Import Manager dashboard** to steer import and validation process
- **EUDAT cloud** with adapted EUDAT services
- Upgraded **CDI User Interface**, ordering and downloading facility

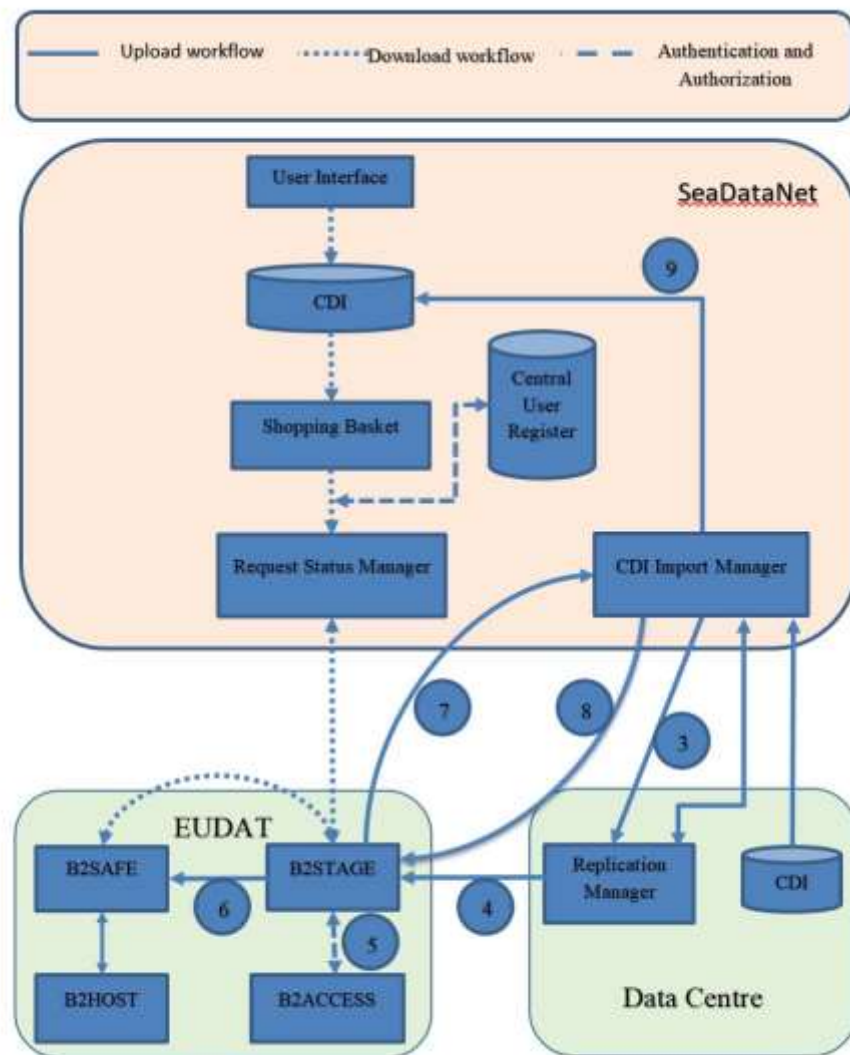
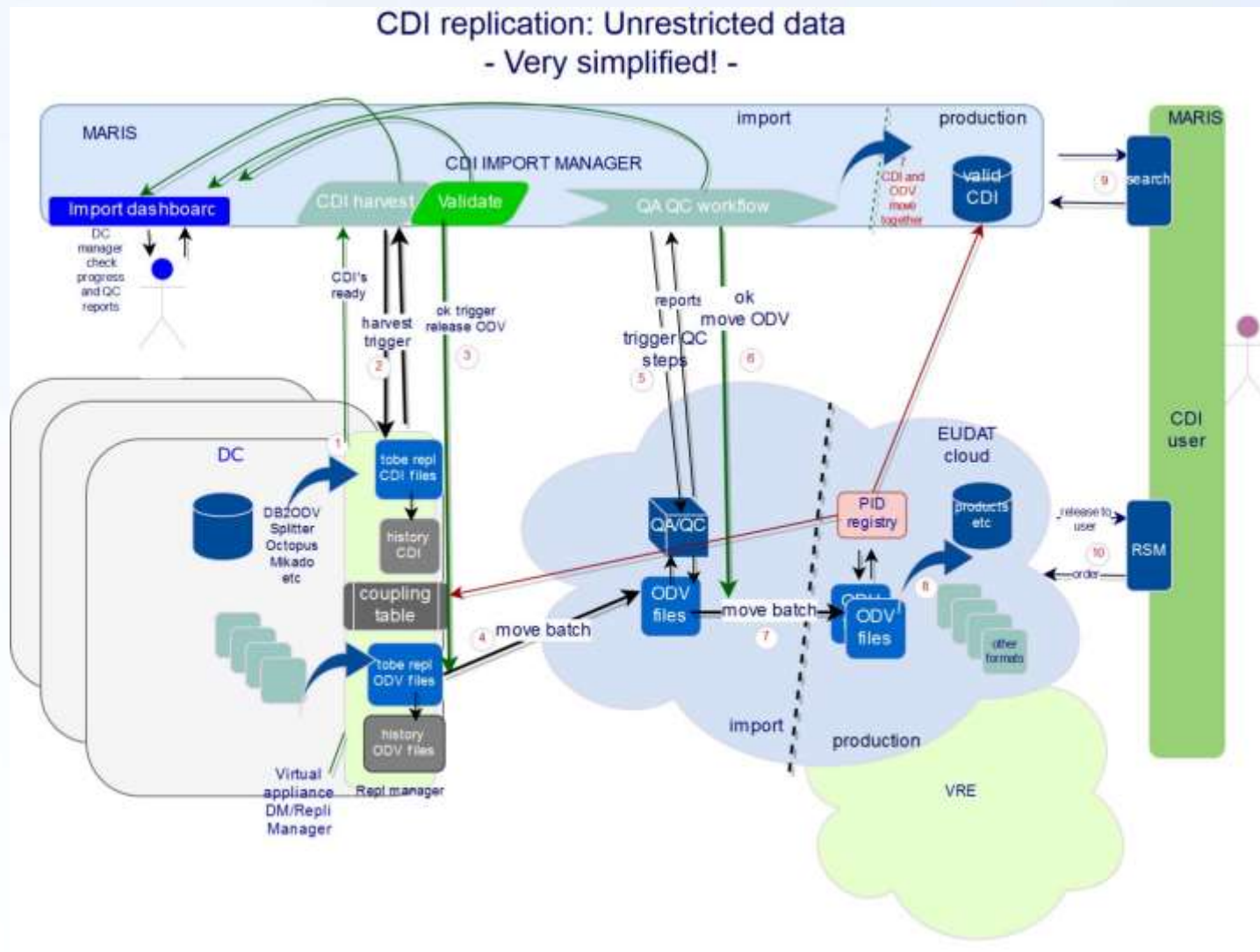


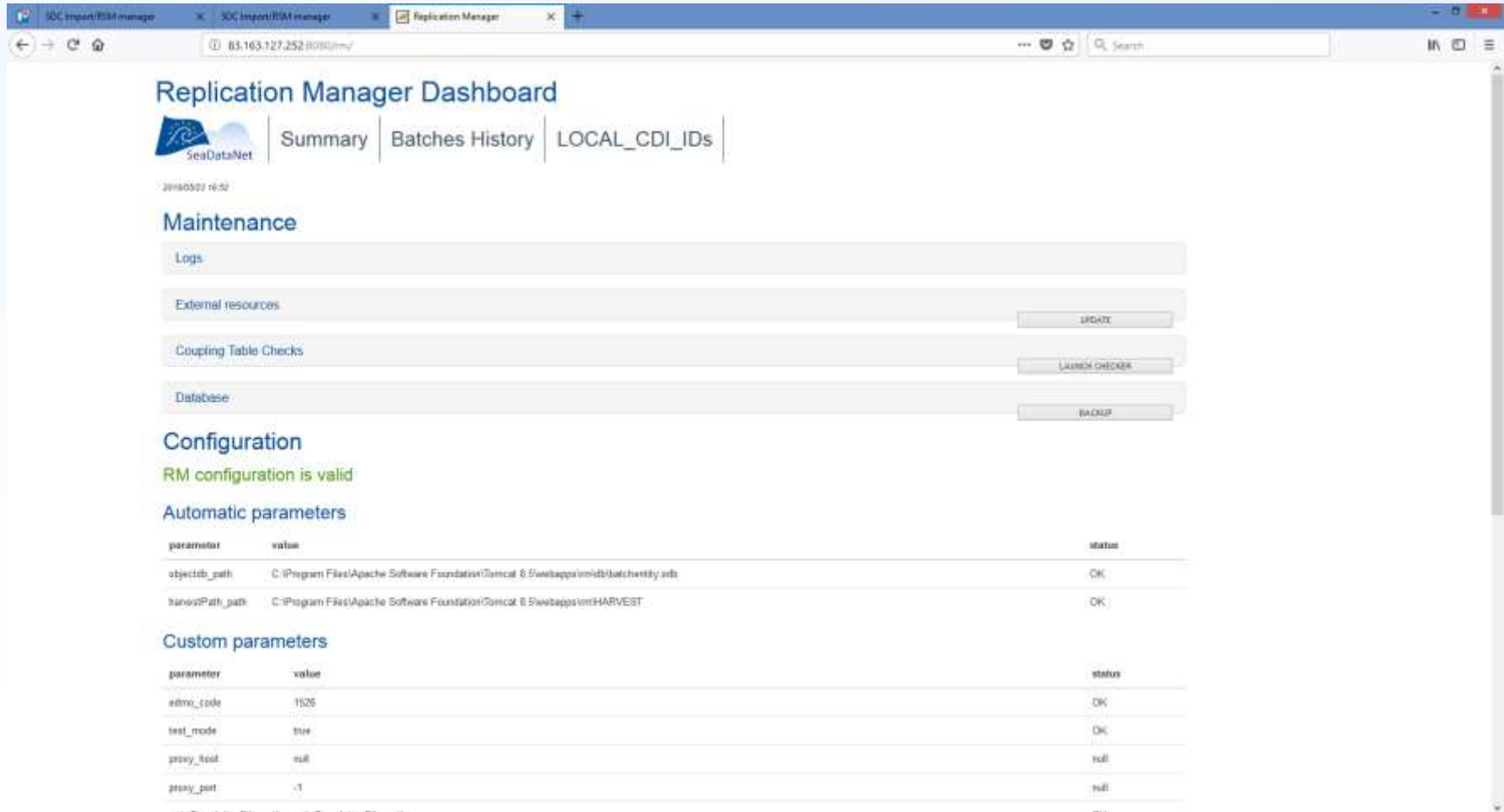
Figure ...: High level architecture

Unrestricted data flow



Unrestricted data ingestion process

1. Authentication
2. Publishing metadata
3. Checking metadata
4. Uploading a databatch
5. Unzip and run data checks
6. Move validated metadata and data in production



The screenshot shows the Replication Manager Dashboard in a web browser. The browser tabs include 'SOC Import/RM manager' and 'Replication Manager'. The address bar shows '83.163.127.252/ROB/...'. The dashboard has a navigation bar with 'Summary', 'Batches History', and 'LOCAL_CDI_IDS'. Below the navigation bar, the date '2019/03/27 16:52' is displayed. The 'Maintenance' section contains four items: 'Logs', 'External resources' (with an 'UPDATE' button), 'Coupling Table Checks' (with a 'LAUNCH CHECKER' button), and 'Database' (with a 'BACKUP' button). The 'Configuration' section shows 'RM configuration is valid'. The 'Automatic parameters' section contains a table with two rows: 'objectdb_path' and 'harvest_path', both with 'OK' status. The 'Custom parameters' section contains a table with five rows: 'edms_code', 'test_mode', 'proxy_host', 'proxy_port', and 'proxy_port', with statuses 'OK', 'OK', 'null', and 'null' respectively.

Replication Manager Dashboard

Summary | Batches History | LOCAL_CDI_IDS

2019/03/27 16:52

Maintenance

Logs

External resources UPDATE

Coupling Table Checks LAUNCH CHECKER

Database BACKUP

Configuration

RM configuration is valid

Automatic parameters

parameter	value	status
objectdb_path	C:\Program Files\Apache Software Foundation\Tomcat 8.5\webapps\rm\mid\batchentity.sdb	OK
harvest_path	C:\Program Files\Apache Software Foundation\Tomcat 8.5\webapps\rm\HARVEST	OK

Custom parameters

parameter	value	status
edms_code	1526	OK
test_mode	true	OK
proxy_host	null	null
proxy_port	-1	null

Import Manager dashboard

Import Batch

Replication Manager

importmanager.seadatanet.org/import_manager_v5/content.asp?screen=4

Login code (1526) Rijkswaterstaat Water, Traffic and Environment

Import Batch

Free search

SEARCH CLEAR

Found 19 Show (1-10)

Key	Input_date	Last_update	Batch_test	Batch_status	Metadata_files_count	Metadata_files_ok	N_code	Edit
325	19/6/2018 17:16 PM	19/6/2018 17:16 PM	True	(1000) Batch created waiting for metadata download	15	0	-	
324	19/6/2018 15:18 PM	19/6/2018 15:37 PM	True	(1000) Batch in production; Ready for a new batch	15	13	-	
323	15/6/2018 13:33 PM	19/6/2018 15:12 PM	True	(2180) Removed from import database	15	15	-	
322	15/6/2018 12:22 PM	15/6/2018 12:22 PM	True	(1000) Batch in production; Ready for a new batch	15	15	-	
317	11/6/2018 13:21 PM	11/6/2018 13:57 PM	True	(1000) Batch in production; Ready for a new batch	100	80	-	
316	11/6/2018 11:55 AM	11/6/2018 12:38 PM	True	(1000) Batch in production; Ready for a new batch	2	1	-	
314	11/6/2018 11:37 AM	11/6/2018 11:53 AM	True	(1000) Batch in production; Ready for a new batch	2	1	-	
313	11/6/2018 10:20 AM	11/6/2018 11:36 AM	True	(1000) Batch in production; Ready for a new batch	26	23	-	
311	7/6/2018 10:32 AM	11/6/2018 10:19 AM	True	(1000) Batch in production; Ready for a new batch	3	3	-	
310	6/6/2018 12:34 PM	7/6/2018 10:30 AM	True	(2180) Removed from import database	14	14	-	

Restricted data flow

- CDI metadata same flow
- Data files extracted at same time as metadata and stored as local copies (to maintain versions)
- Shopping requests in RSM are evaluated by data centre
- If ok, data files will be temporarily released from local to cloud directory of user, upon request RSM



Benefits for users

- The performance will be increased, discovery and data requests improved, and downloading made more easy as each shopping request will provide one integrated download package instead of multiple packages from multiple data centres.
- Overall quality and coherence (data – metadata) will improve
- Tracking and tracing of data transactions will continue to be administered by an upgraded and much faster RSM service to oversee shopping requests and deliveries. The user RSM will be integrated as MySeaDataCloud service in the CDI user interface.
- Versioning of metadata and data will facilitate repeated analysis of e.g. environmental assessments in MSFD context after many years, and for scientific papers.

Benefits for data centres


- Data centres will have a **Replication Manager** module and an **Import Dashboard** to trigger and control themselves the import of new and updated metadata and data sets (unrestricted) into the CDI service
- Data providers can oversee all relevant transactions for their data centre in the upgraded and much faster RSM system and generate relevant reports
- The system will also support handling restricted data sets
- Data centres will be outfitted with a Replication Manager (RM) replacing the Download Manager. The RM has less complexity and is easier to configure.
- Alternatively, Data centres can make use of the '**indirect solution**' which will be provided with improved functionality, handling both unrestricted and restricted data sets

New GUI on top of a new system



- Faster, using Elastic Search for search and indexing, and GeoServer for latest mapping technology
- Easier, including full text search next to controlled terms
- Modern design with large map and sliding windows
- **MySeaDataNet** integrated in GUI for customized services, such as SSO, RSM access, search profiles, prepared for VRE data pooling
- Developed and refined as prototype (<http://sdc.maris.nl>) interacting with users and data managers, following the SeaDataCloud Training Workshops in June 2018

New GUI impressions



PAN-EUROPEAN INFRASTRUCTURE FOR OCEAN & MARINE DATA MANAGEMENT

STANDING ORDERS | **ORDERS HISTORY** | **SAVE**

ORDER DETAILS

Stukje tekst over de betekenis van de kleuren console auctor, Phasellus sagittis purus nec augue malesuada dolor quis volutpat. Nulla pharetra lectus ipsum, sit

Waiting for processing | Approval pending

ID	Name	Unrestricted
17501	IMDIS2018 test	2
17500	ARGO floats	25
17499	IMDIS2018 test CTDs	100
17498	IMDIS2018 test	1
17493	Bathymetry order	2
17492	Test ARGO op NetCDF	
17491	Test ARGO floats	7

DOWNLOAD THE ORDER

Order ready to download

Order number: 17501
 Last update: 11/4/2018 11:20:59 PM
 Input date: 11/4/2018 11:20:59 PM
 Organisation: MARIENE INFORMATIE SERVICE MARIS B.V.
 Contact name: D.m.a. SCHAAP

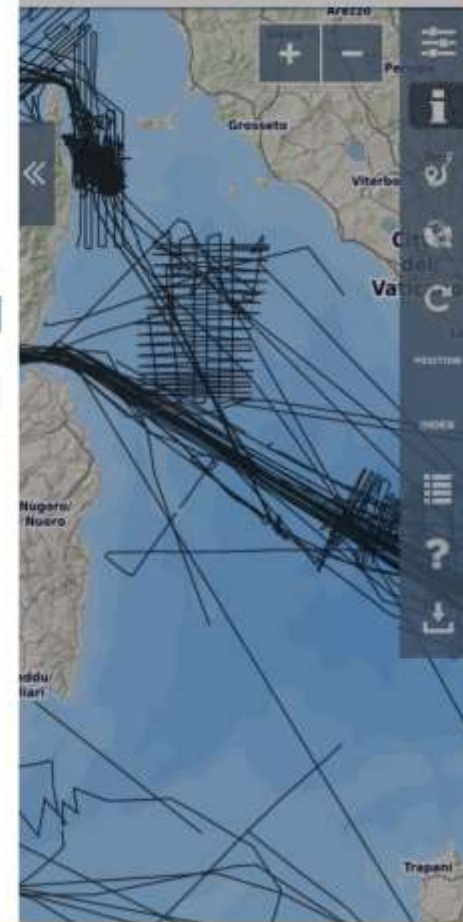
Unrestricted

[RECORDS \(1.05 MB\)](#)

[METADATA CSV FILE \(1.19 KB\)](#)

ello D.m.a. SCHAAP | My CDI

DATASET BASKET



Map showing data tracks in the Tyrrhenian Sea, with locations like Grosseto, Viterbo, and Trapani visible.



www.seadatanet.org