



**SeaDataNet**

*PAN-EUROPEAN INFRASTRUCTURE  
FOR OCEAN & MARINE DATA  
MANAGEMENT*

## **Status of the biological QC developments**

**Simon Claus, Daphnis De Pootere & Filip Waumans**

VLIZ- Flanders Marine Institute

## WP10.2.4: development of online Biology Data quality checks (QC)

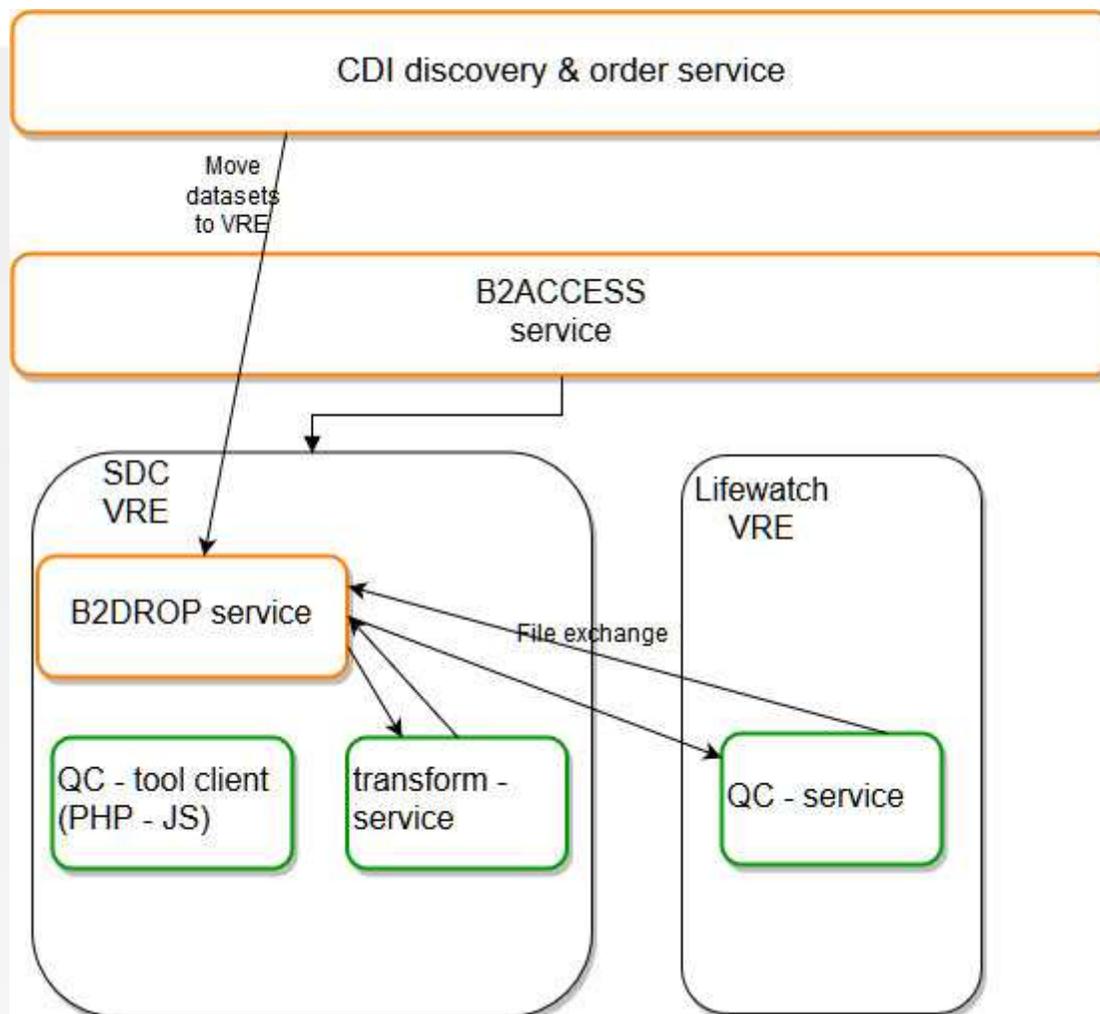
This task will develop online services and tools to analyse the quality and completeness of biology data. Records will be reviewed through a series of QC steps dealing with the

- 1) data format
- 2) completeness and validity of information
- 3) taxonomic quality
- 4) geographic quality
- 5) whether or not the record is a (possible) outlier.

- D10.9: Specification of Biology Data QC online and development plan M12=>M18**
- D10.10: Phase 1 of Biology Data QC online operational (M24)
- D10.11: Phase 2 of Biology Data QC online operational (M36)
- D10.12: Phase 3 of Biology Data QC online operational (M42)

## D10.10: Phase 1 of Biology Data QC online operational (M24)

- The development of the **transformation service** in a Docker container which will convert the BIODDEF ODV format in DwC-A format and vice versa
- The implementation of the **QC steps** mentioned in point 3 of the use case description
- Tuning with the overall architecture and VRE development.



[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)  
[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)



DATA SUPPORTING INFRASTRUCTURE  
FOR OCEAN & MARINE DATA  
MANAGEMENT

## 4. Transformation service

SPECIFIC DOCUMENTATION FOR **BIOLOGICAL DATA**

DESCRIPTION OF THE SEADATANET DATA TRANSPORT FORMAT FOR BIOLOGICAL DATA

[DOWNLOAD \(991.96 KB\)](#)

TEMPLATE AND EXAMPLES OF BIOLOGICAL DATA FILES (ZIP FILE)

[DOWNLOAD \(550.84 KB\)](#)



Darwin Core Archive

OBIS-ENV-DATA

```
RStudio
File Edit Code View Plots Session Build Debug Profile Tools Help
Go to file/function Addins
EMODnet_allservices.R bd bod prepdata.Rmd EMODnet_biology.R dome dome1 BaltizZpprepData.Rmd
Source on Save
36 close(con)
37
38 ##read datafile in dataframe variable skipping commentlines
39 data<-read.delim("data_from_collection_minimumObservationDepth.txt",
40                skip = commentlines, colClasses = "character", na.strings = "" )
41
42 #####
43                map header with DwC terms #####
44 #####
45 # map sdc --> dwc
46 map_sdc_terms <- c(LocationID= "Station name",
47                  eventDate= "yyyy.mm.ddThh.mm.ss.sss",
48                  Locality= "Station",
49                  decimalLongitude= "Longitude..degrees_east.",
50                  decimalLatitude= "Latitude..degrees_north.",
51                  minimumdepthmeters= "minimumdepth",
52                  maximumdepthmeters= "maximumdepth",
53                  samplingprotocol= "samplingprotocol.#{.}")
54
55 dwc_terms <- c("eventDate",
56              "decimalLongitude",
57              "decimalLatitude",
58              "scientificName",
59              "scientificNameID",
60              "occurrenceStatus",
61              "basisOfRecord",
62              "minimumDepthInMeters",
63              "maximumDepthInMeters")
64
```





- Is the dataset integrity ok?
  - Do all eventID's in the occurrence extension refer to an Event Record?
  - Do all eventID's in the eMoF extension refer to an Event Record?
  - Do all occurrenceID's in the eMoF extension refer to an occurrence Record?
  - In case of a biometrical parameter, is the eventID in the eMoF link the same one as in the occurrence extension?
  - Are there 'duplicate occurrences' meaning is the same taxon listed twice at the same EventID without any difference in any of the biometric parameters?
  - Are there 'duplicate measurements' meaning does the same measurement occur twice for the same occurrenceID or the same EventID?
- Are all mandatory fields present in the dataset?
- Are all mandatory field filled out?
- Does eventDate follow the required format?
- Does the scientificNameID follow the required format?
- Are there coordinates located on land? (buffer of 3km is taken under consideration)
- Are there depths at the location deeper than the depths stored by GEBCO? (a margin of 150m is taken under consideration)
- Do all measurementtypes have a MeasurementTypeID?
- Does the measurementTypeID / measurementValueID refer to an existing term in the BODC vocabulary?
- Do all measurementValues that refer to facts have a measurementValueID?
- Are there records where measurementValue is NULL?
- Are there records that refer to biological measurement where measurementValue = 0 (and where occurrenceStatus is not absent)?

[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)  
[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)



EUROPEAN INFRASTRUCTURE  
FOR OCEAN & MARINE DATA  
MANAGEMENT

Additional checks which are planned to be implemented are:

- Is there a sampling instrument present?
- Are there other sampling descriptors present?
- Provide an overview of the number of taxa per kingdom and class.
- Check for non-marine and non-brackish taxa.
- Has the unit provided the same base unit as specified by the measurementTypeID?
- add the name of the dataset + URL at the top + number of records per file
- List the non-matched taxa
- Add geographical cover and temporal cover from the eml to the overview page.

[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)  
[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)



DATA SUPPORT INFRASTRUCTURE  
FOR OCEAN & MARINE DATA  
MANAGEMENT

## QC service

- Tool developed in R, all scripts are available from the EMODnet github repository at <https://github.com/EMODnet/LifeWatch-EMODnet-Biology-QC-tool>.
- Is using the OBIS tools package <https://github.com/iobis/obistools>
- Available from the LifeWatch services at <http://rshiny.lifewatch.be/BioCheck/>

LifeWatch & EMODnet Biology QC tool

[Data overview](#) [Issues found](#) [Issues on map](#) [Invalid Event Records](#) [Invalid Occurrence Records](#) [Invalid MoF Records](#) [About](#)

Welcome to the LifeWatch and EMODnet data quality control tool. This tool allows you to assess to what extent a dataset published on an IPT or a DwC-A file meets the EMODnet Biology Data Quality Criteria. Provide the URL to the homepage of the IPT resource in the field IPTLink and press load. This may take a few minutes for large (+100k occurrence records) datasets. Then navigate through the different tabs above to get detailed QC information, or press the download button to resolve the issues later. The tab about specifies which QC checks are carried out, and which QC check are planned to be included. Please use the R functions at <https://github.com/EMODnet/LifeWatch-EMODnet-Biology-QC-tool> or contact [bio@emodnet.eu](mailto:bio@emodnet.eu) if you need a qc report for datasets with more than 500k occurrence records.

Link to IPT resource

Overview of dataset: Benthic Network

Overview of event and occurrence records

basisOfRecord	occurrenceStatus	countOccurrence
HumanObservation	present	113482

Overview of the measurement or fact records

type	measurementType	minValue	maxValue	measurementUnit	count	preflabel	definition
OccurrenceMoF	Sampling instrument	NA	NA	NA	52245	Sampling instrument name	The name of the gear or instrument used to collect the sample or make the in-situ measurement or observation.

[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)  
[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)



DATA SUPPORT INFRASTRUCTURE  
 FOR OCEAN & MARINE DATA  
 MANAGEMENT

## QC service

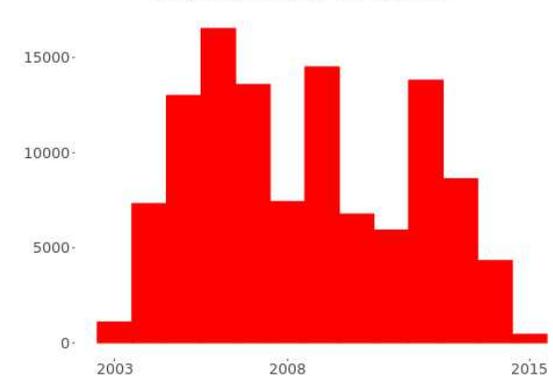
row	maximumDepthInMeters_error	minimumDepthInMeters_error
1	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
2	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
3	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
4	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
5	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
6	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
7	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
8	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
9	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters
10	Empty value for recommended field maximumDepthInMeters	Empty value for recommended field minimumDepthInMeters

Previous 1 2 3 4 5 ... 10687 Next

Geographical cover of the dataset



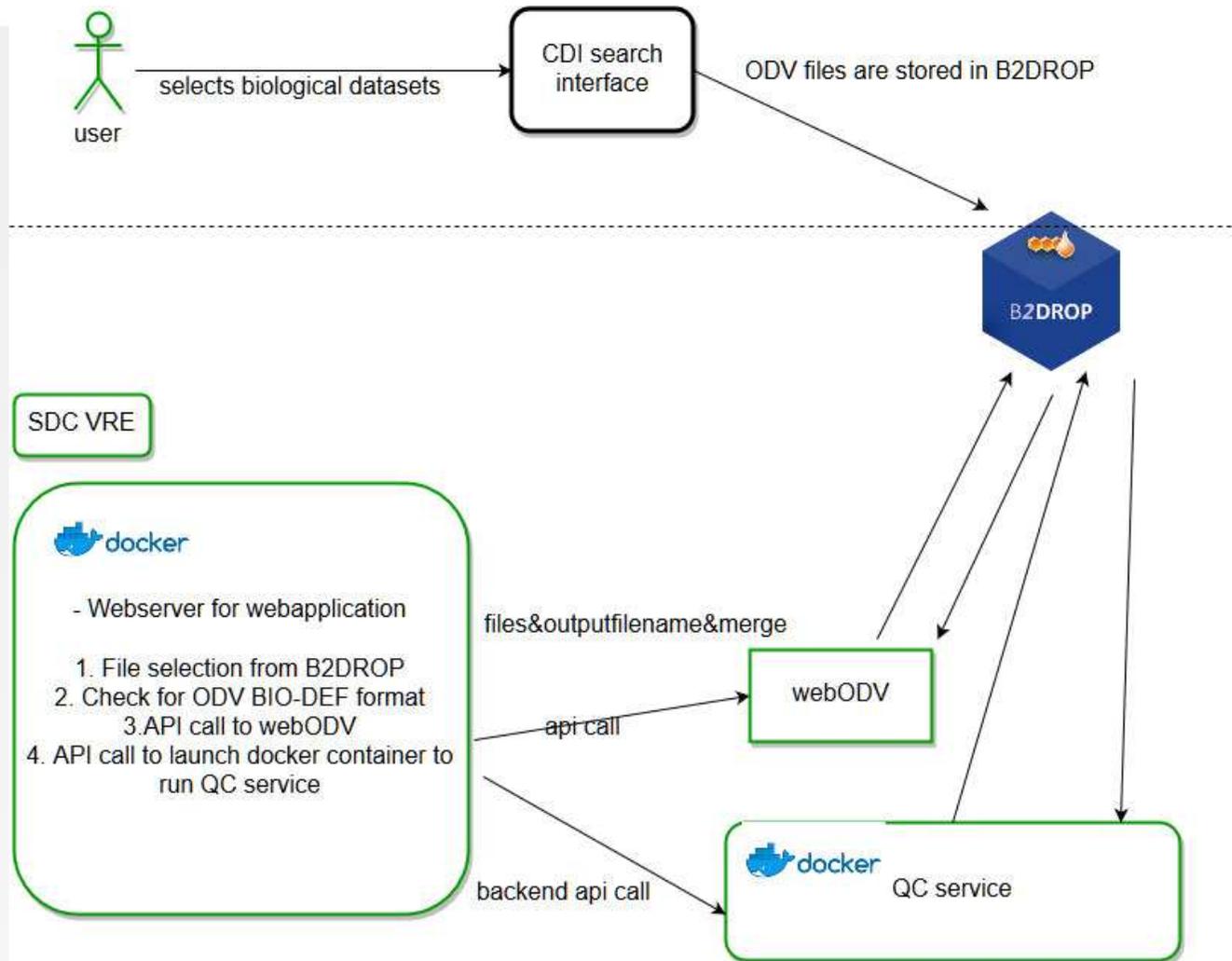
Temporal cover of the dataset





SeaDataNet

PAN-EUROPEAN INFRASTRUCTURE  
FOR OCEAN & MARINE DATA  
MANAGEMENT



A photograph of a sunset over a beach. The sun is low on the horizon, partially obscured by clouds, creating a warm, golden glow. The sky is filled with soft, textured clouds. In the foreground, the ocean waves gently wash onto a sandy beach. A dark, vertical post stands in the water on the right side of the frame, with a seagull perched on top. The overall mood is peaceful and serene.

Thank you for your attention!

[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)  
[http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine\\_only=true](http://marinespecies.org/rest/AphiaIDByName/abra%20alba?marine_only=true)



DATA SUPPORT INFRASTRUCTURE  
FOR OCEAN & MARINE DATA  
MANAGEMENT

## The different components are:

### 1. CDI discovery service

This service allows users to select BIO-DEF ODV files based via the CDI system.

### 2. B2ACCESS/ Marine ID service

These service will be used for authorization and authentication to the VRE. The user logs in via his/her Marine ID account, that maps to B2ACCESS by which the SDC VRE components will be secured.

### 3. File repository

The B2DROP environment will be used as a shared file repository. In the SDC VRE the user collects the biological datasets in the personal space. Results of transformations and QC are collected again in the same space.



## The different components are:

### 4. Transformation service

Biodiversity services such as the QC service work preferable with Darwin Core standards. This transformation service running within the SDC VRE platform (in a Docker container like the other transformation and processing services on the EUDAT platform) will convert the BIODF ODV files into Darwin Core, ready to serve the files later to the QC service.

### 5. Client application

Within the SDC VRE, the GUI of the SeaDataNet Biology quality assessment tool will be developed as a HTML, JavaScript, PHP web application running in a container on the EUDAT B2HOST platform. The front-end application will communicate with several services on the LifeWatch platform.

### 6. QC service

This service will execute the QC procedures and add the information to the ODV-file(s). This service will be developed as a REST-service calls by the user interface (component 5).