

1st **SeaDataCloud** Training Workshop  
Oostende, (20-27 June 2018)

# Biological data

Daphnis De Pooter & Paula Oset (VLIZ)





# Outline

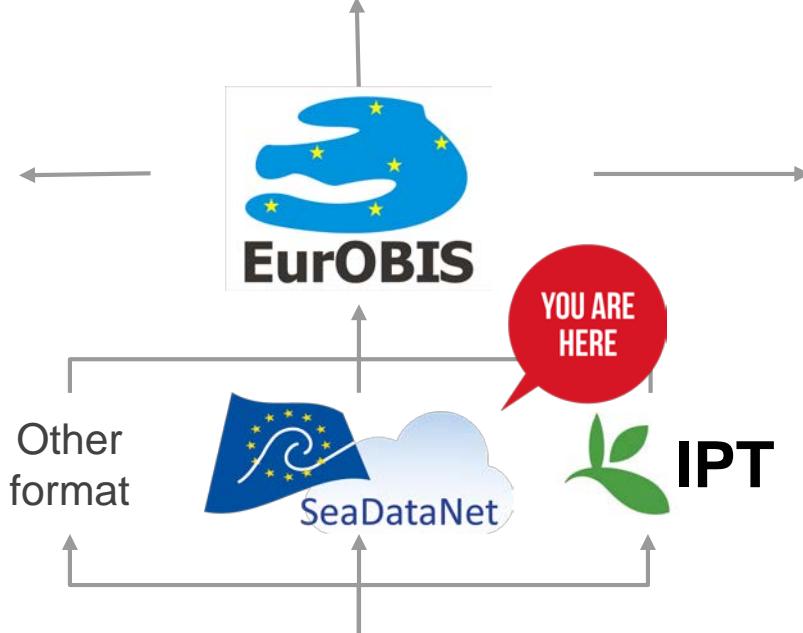
- EMODnet Biology data flow
- Introduction to biological data
- EMODnet Biology: OBIS Event Core format & Darwin Core
  - OBIS Event Core
  - Darwin Core terms
  - Standards: BODC vocabulary and taxonomy (WoRMS)
  - Exercise: Taxonomy and use of World Register of Marine Species
- SDC and biological data
  - Example 1: Data format: CDI + ODV biology variant
- Biological Data QC
  - Demo using R



# Data flow



OCEAN BIOGEOGRAPHIC  
INFORMATION SYSTEM



Data providers



---

# Introduction to biological data

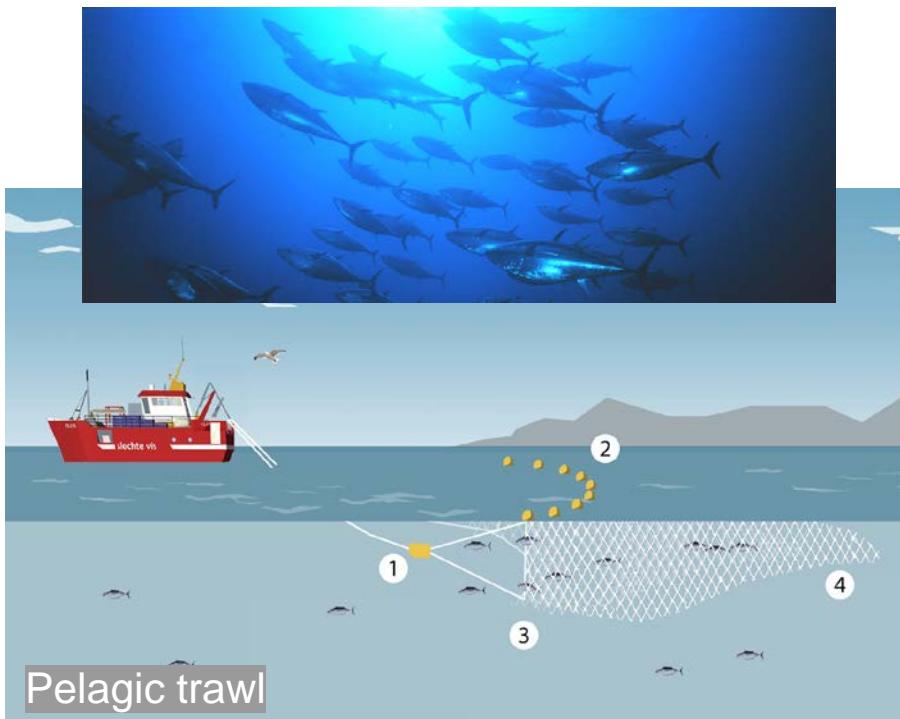


# Biodiversity

- Biodiversity, the variety of life found in a place on Earth. Common measure: the count of species in an area (species richness).
- Basic information:
  - What
  - Where
  - When
- Additional information:
  - How many/much
  - Under which environmental conditions



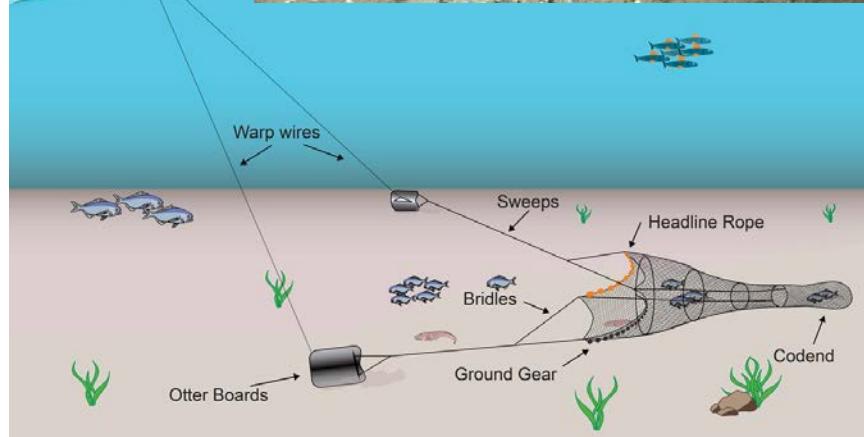
# Biological sampling:





# Biological sampling:

Bottom trawl



Van Veen grab



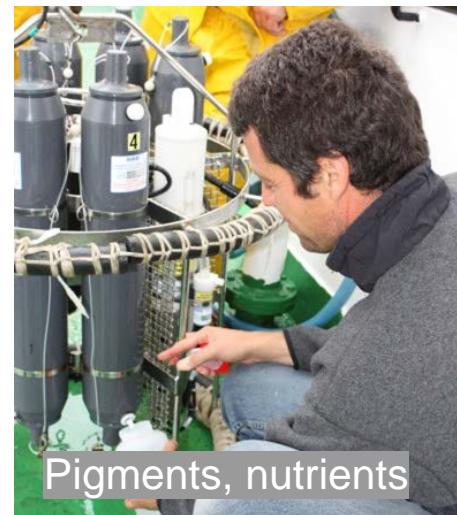


# Biological sampling:





# Biological sampling:





# Biological sampling:





# Biological sampling:





- How to capture these complexities and different types of information?



---

# Overview of biological data format



01

Data schema / structure

- 3 tables ([OBIS-Event core data format](#)):
  - Event core table.
  - Occurrence table.
  - Extended Measurements or Facts (eMoF).

02

Field nomenclature

- Darwin Core (DwC) standard [terms](#).
- Minimum of fields required per table.

03

Content - Controlled vocabulary  
and standards

- Date/time & Lat/lon.
- EventID and OccurrenceID.
- Taxonomic information: [LSID\\*](#)
- Other parameters: [BODC-NERC vocabulary](#).



---

## OBIS Event core format

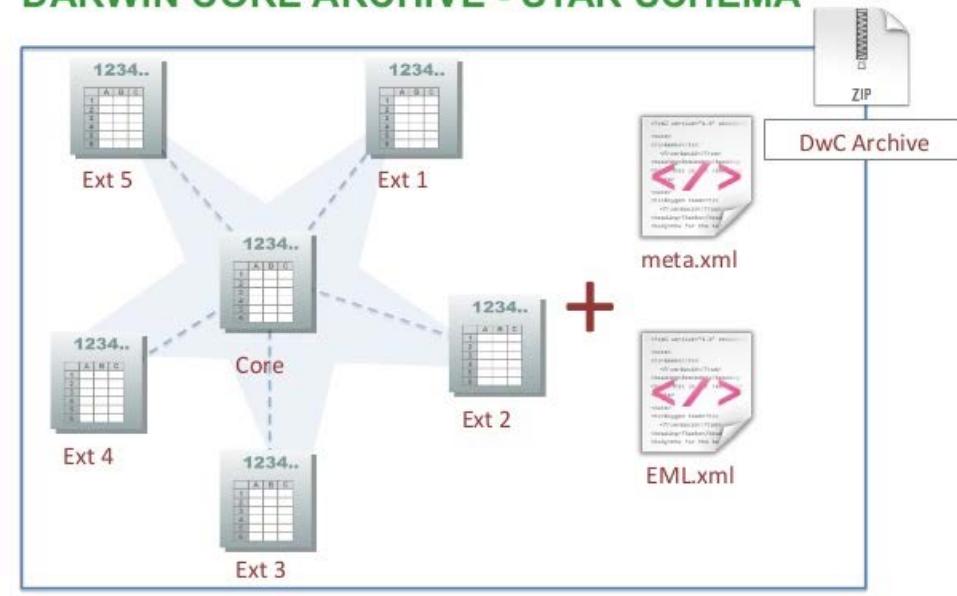


# Darwin Core schema

The conceptual data model of the Darwin Core Archive is a “**star schema**” (Robertson et al. 2014):

- **Core record**, such as an occurrence or an event, as the center of the star.
- **Extension records**, radiating out of the star, can optionally be associated with the core, linked by database keys such as an ID column.

## DARWIN CORE ARCHIVE - STAR SCHEMA



Slide source: GBIF G823 Nodes training & iDigBio, Florida 2015



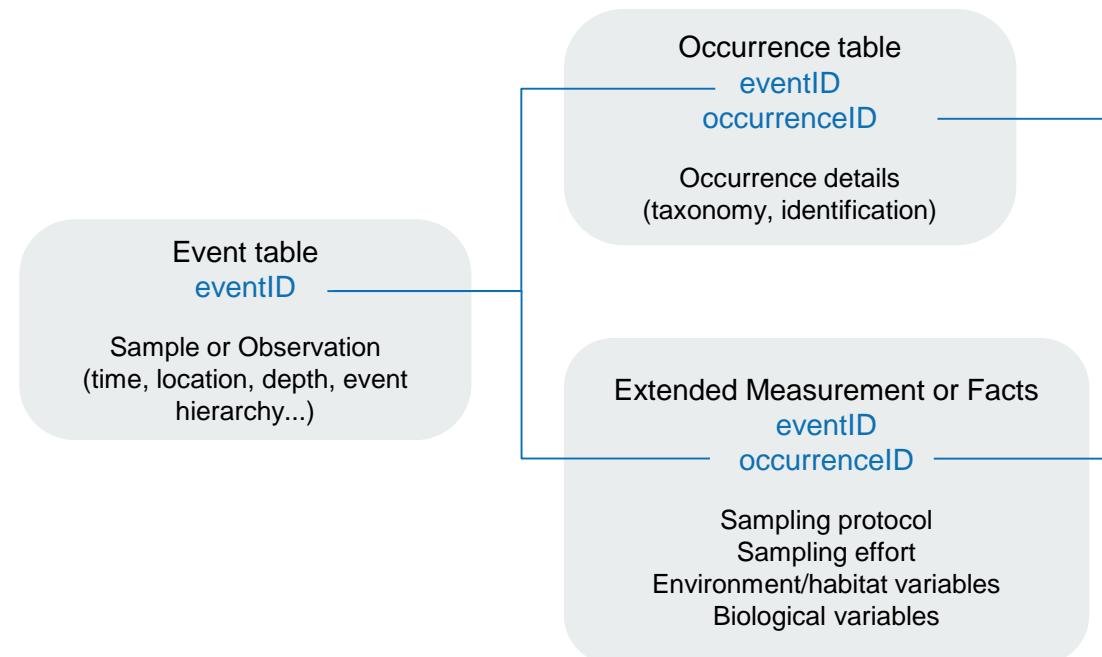


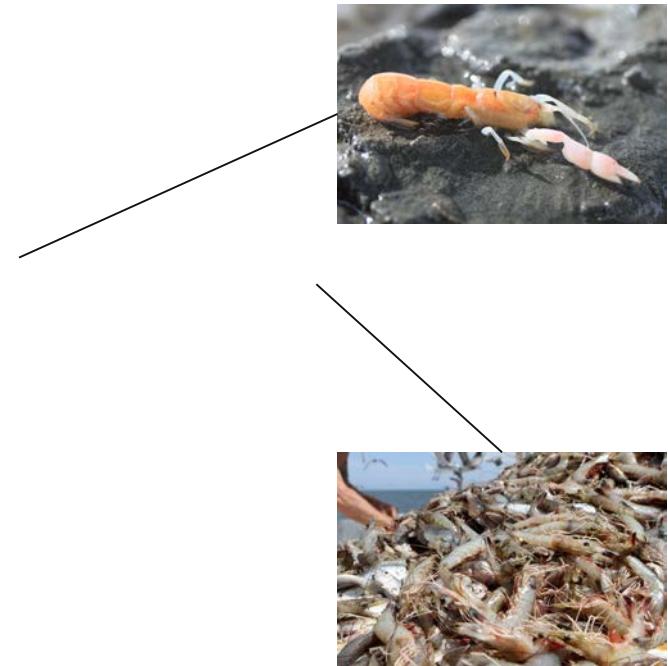
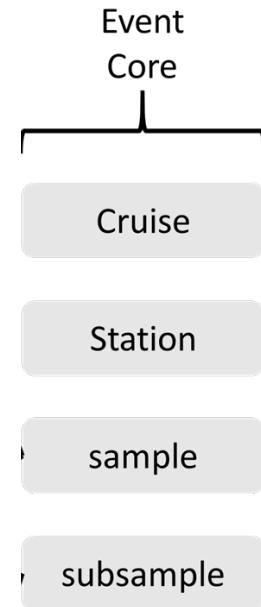
01

## Data schema / structure

Data is structured in 3 tables related to each other via the eventID and the occurrenceID. This structure allows to store not only **occurrences** but also **sampling information** and **additional biological and/or abiotic measurements**.

- 3 tables ([OBIS-Event core data format](#)):
  - Event core table.
  - Occurrence table.
  - Extended Measurements or Facts (eMoF).







---

## Field names: Darwin Core terms



## 02

### Field nomenclature

- Darwin Core (DwC) standard terms.
- Minimum of fields required per table.

The DwC terms that are most relevant to EMODnet Biology format are the following (those in **bold** are mandatory):

#### Event table

*datasetName, eventID, parentEventID, eventDate, institutionCode, habitat, type, minimumDepthInMeters, maximumDepthInMeters, decimalLatitude, decimalLongitude, coordinateUncertaintyInMeters, footprintWKT, modified*

#### Occurrence table

*eventID, occurrenceID, scientificName, scientificNameAuthorship, scientificNameID, kingdom, taxonRank, identificationQualifier, occurrenceStatus, basisOfRecord, modified*

#### Extended Measurement or Fact table

*measurementID, eventID, occurrenceID, measurementType, measurementTypeID, measurementValue, measurementValueID, measurementUnit, measurementUnitID, measurementAccuracy, measurementRemarks*



---

# **Content: controlled vocabulary and standards**



## 03

### Content - Controlled vocabulary and standards

- Date/time & Lat/lon.
- EventID and OccurrenceID.
- Taxonomic information: [LSID](#)
- Other parameters: [BODC-NERC vocabulary](#).

Besides the field names, the content or the data itself has to follow certain standards. For example, the date-related fields have to be ISO 8601 compliant, the latitude and longitude have to be in decimal degrees and referenced to the WGS84 projection (EPSG:4326 datum).

An overview of the required format for the content of the different fields is available [here](#).

- EventID and OccurrenceID, coordinates and date (Common terms)
- Taxonomic information: LSID
- Other parameters (eMoF table): BODC-NERC controlled vocabulary.



---

---

## **BODC controlled vocabulary in the eMoF table**



## Controlled vocabulary (eMoF table)

- The eMoF extension is used to store:
  - information related to sampling method and sampling effort (via eventID).
  - measurements linked to a biological occurrence (via occurrenceID).
  - environmental measurements (via eventID).
- The MoF terms: measurementType, measurementValue and measurementUnit are completely unconstrained and can be populated with free text annotation.
  - Free text: to capture complex and as yet unclassified information
  - But heterogeneity (e.g. of spelling or wording) becomes a major challenge for effective data integration and analysis.



## Controlled vocabulary (eMoF table)

- Three fields to standardise the measurement types, values and units: **measurementTypeID**, **measurementValueID** and **measurementUnitID**.
- These terms are populated using **controlled vocabularies** from the NERC Vocabulary Server, developed by the British Oceanographic Data Centre (BODC)
  - [https://www.bodc.ac.uk/resources/vocabularies/vocabulary\\_search/](https://www.bodc.ac.uk/resources/vocabularies/vocabulary_search/)

MeasurementType	MeasurementTypeID
(free text)	(controlled vocabulary)
Body length	<a href="http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX">http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX</a>
Length	<a href="http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX">http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX</a>
Length (mm)	<a href="http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX">http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX</a>
length_in_mm	<a href="http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX">http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX</a>
Length of specimen	<a href="http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX">http://vocab.nerc.ac.uk/collection/P01/current/OBSINDLX</a>

Example of parameter standardization using controlled vocabulary.



## Controlled vocabulary (eMoF table)

id	Occurrence ID	Measurement Type	MeasurementTypeID	Measurement Value	Measurement Unit	Measurement UnitID
BF1M1	BIOFUN1_BF1M1_1	abundance	<a href="http://vocab.nerc.ac.uk/collection/P01/current/OCOUNT01">http://vocab.nerc.ac.uk/collection/P01/current/OCOUNT01</a>	26	Individuals	<a href="http://vocab.nerc.ac.uk/collection/P06/current/UUUU">http://vocab.nerc.ac.uk/collection/P06/current/UUUU</a>
BF1M1	BIOFUN1_BF1M1_1	density	<a href="http://vocab.nerc.ac.uk/collection/P01/current/SDBIOL02">http://vocab.nerc.ac.uk/collection/P01/current/SDBIOL02</a>	0.000329114	N/km <sup>2</sup>	<a href="http://vocab.nerc.ac.uk/collection/P06/current/NPKM">http://vocab.nerc.ac.uk/collection/P06/current/NPKM</a>
BF1M1	BIOFUN1_BF1M1_1	Wet Weight Biomass	<a href="http://vocab.nerc.ac.uk/collection/P01/current/SDBIOL05">http://vocab.nerc.ac.uk/collection/P01/current/SDBIOL05</a>	0.091139241	kg/km <sup>2</sup>	

Example of measurements linked to a biological occurrence (via occurrenceID)

id	measurementType	measurementTypeID	measurementValue	measurementValueID	measurementAccuracy	measurementUnit
BIOFUN_BF1M15	Trawling speed	<a href="http://vocab.nerc.ac.uk/collection/P01/current/APSAZZ01">http://vocab.nerc.ac.uk/collection/P01/current/APSAZZ01</a>	2.7		0.1	knots
BIOFUN_BF1A01	Gear	<a href="http://vocab.nerc.ac.uk/collection/Q01/current/Q0100002">http://vocab.nerc.ac.uk/collection/Q01/current/Q0100002</a>	Agassiz dredge	<a href="http://vocab.nerc.ac.uk/collection/L22/current/TOL0991/">http://vocab.nerc.ac.uk/collection/L22/current/TOL0991/</a>		

Example of information related to sampling method and sampling effort (via eventID)



## Useful links

<http://iobis.org/manual/emof-bodc/>

<http://iobis.org/vocab/>

[https://www.bodc.ac.uk/resources/vocabularies/vocabulary\\_search/](https://www.bodc.ac.uk/resources/vocabularies/vocabulary_search/)

[http://seadatanet.maris2.nl/v\\_bodc\\_vocab\\_v2/search.asp?lib=P01](http://seadatanet.maris2.nl/v_bodc_vocab_v2/search.asp?lib=P01)

(abundance%entity)

[Example dataset fully processed](#)



---

**Taxonomic information:**  
**World Register of Marine Species (WoRMS)**

---



## Why a taxonomic standard?

**Taxonomy** = the academic discipline of defining groups of biological organisms on the basis of shared characteristics and giving names to those groups.

### Spelling errors

*Actinobacillus actimomycetemcomitans*  
*Actinobacillus actimycetemcomitans*  
*Actinobacillus actinmymcetemcomitans*  
*Actinobacillus actinomicetemcomitans*

...

### Synonymy

*Halichondria panicea* Pallas, 1766  
*Hymeniacidon parfitti* Parfitt, 1868  
*Halichondria paciscens* Schmidt, 1875  
*Menanetia minchini* Topsent, 1896  
>60 synonyms -> [check](#)

### Spelling variations (correct)

*Agalinus paupercula borealis*  
*Agalinus paupercula* var. *borealis*  
*Agalinus paupercula* var. *borealis* Pennell  
*Agalinus paupercula* (Gray) Britt. var. *borealis*  
Pennell

...

### Homonymy

*Alebion*





# World Register of Marine Species

taxonomy

Expert editors

traits

distribution

The World Register of Marine Species (WoRMS) is the taxonomic backbone of EMODnet-Biology and OBIS.

→ All taxa in your dataset need a scientificNameID from WoRMS!

**WoRMS**  
World Register of Marine Species

Home About Subregisters Users Photogallery Documents LifeWatch Contribute

Quick search...

Taxa  Literature  Distribution  Specimen  Editors  Statistics  Tools  Manual  Log in

The World Register of Marine Species aims to provide the most authoritative list of names of all marine species globally, ever published



Latest taxon additions  
Updated: 2018-04-27T15:14

Anapaleoerchis hamjimai Fujino & Kifune, 1991	
Added: 2018-04-27	
Nanina kintana Morgan, 1885	
Added: 2018-04-27	
Dyakia kintana (Morgan, 1885)	
Added: 2018-04-27	
Dyakia lindstedti (Pfeiffer, 1856)	
Added: 2018-04-27	
Helix lindstedti Pfeiffer, 1856	
Added: 2018-04-27	

News

WoRMS press release: Celebrating a Decade of the World Register of Marine Species  
Added on: 2018-04-23 09:01:17 by Vandepitte, Leen



In 2018, to celebrate a decade of WoRMS' existence, it was decided to compile a list of our top marine species, both for 2017 and for the previous decade. ....

[Read more](#)

Tweets by @WRMarineSpecies

WoRMS Retweeted



## LSID: taxonomic information

All the occurrences are given a unique scientificNameID. This is done by matching the Scientific Names of your occurrence table with the [World Register of Marine Species](#), using the [taxon match tool](#). Information on how to use the taxon match tool [here](#).

After matching, the tool will return you a file with the AphiaIDs, LSIDs, valid names, authorities, classification and any other output you have selected. The WoRMS **LSID** is used for DwC field **scientificNameID**.

ScientificName	scientificNameID
Alepocephalus rostratus	urn:lsid:marinespecies.org:taxname:126684
Bathypterois mediterraneus	urn:lsid:marinespecies.org:taxname:299942
Coelorinchus mediterraneus	urn:lsid:marinespecies.org:taxname:280313
Galeus melastomus	urn:lsid:marinespecies.org:taxname:105812
Lepidion lepidion	urn:lsid:marinespecies.org:taxname:126495
Mora moro	urn:lsid:marinespecies.org:taxname:126497

**Important:** DO NOT provide the “accepted” scientificName and scientificNameID,

Please DO provide the original taxon name and the scientificNameID of the match



## Taxon match guided exercise

[Download exercise](#)

Go to <http://marinespecies.org/aphia.php?p=match>



## Taxon match guided exercise

[Download solution](#)



# SDC and Biological Data

- SDC Data to EurOBIS
  - Automated transformation should be possible
    - Same terms
    - Similar structure
- Biological Data Exchange Format (BioDEF)





# Biological Data Exchange Format (BioDEF)

- CDI metadata
- ODV biology variant version 2.0
  - Semantic header
  - 9 Mandatory ODV fields
  - 9 Mandatory BioDEF fields
  - Conditional BioDEF fields
  - Optional BioDEF fields
  - Quality flags



# BioDEF Data Format

HEADER  
describing fields

```
//<subject> ...<object>...<units>....<instrument>...
//<subject> ...<object>...<units>....<instrument>....
//....
```

Fields 1 - 9

Fields 10 - 27

Fields ...

DATA  
TABLE

9 mandatory  
ODV fields

9 Mandatory  
BioDEF fields  
+ QC flags

# conditional and  
optional fields  
+ QC flags

Extendable!



# BioDEF Data Format

## 9 Mandatory ODV fields

- Cruise
- Station
- Type
- yyyy-mm-ddThh:mm:ss.sss
- Longitude [degrees\_east]
- Latitude [degrees\_north]
- LOCAL\_CDI\_ID
- EDMO\_code
- Bot. Depth [m]



# BioDEF Data Format

## 9 Mandatory BioDEF fields

- MinimumDepthOfObservation
- MaximumDepthOfObservation
- SampleID
- SamplingEffort
- ScientificName
- ScientificNameID
- Sex
- LifeStage
- ObservedIndividualCount

### 4 possible P01 codes

1. Area sampled of the bed
2. Volume sampled of the water body
3. Length of sampling track
4. Sample duration



# BioDEF Data Format

## Conditional BioDEF fields

- EventStartTime
  - EventEndTime
  - EventStartLongitude
  - EventEndLongitude
  - EventStartLatitude
  - EventEndLatitude
- 
- Time range for sampling
- 
- SubsampleID
  - SubSamplingCoefficient
  - Samplingprotocol
  - Occurrencestatus
- Sampling along track (e.g. trawling)
- 
- In case of subsampling



# BioDEF Data Format

## Optional BioDEF fields

- Abundance per unit area
  - Wet weight biomass per unit area
  - Coverage
  - ....
  - Length of biological entity
  - Size class of biological entity
  - ...
  - Sediment type category
  - JNCC habitat type (version 04.03)
  - ...
- 
- Biota quantification
- Biometrics
- Abiotic data



# BioDEF Data Format

## Template and examples



# QC Procedures Biological data

## To Verify

- Are all mandatory fields **present**?
- Are all values of mandatory fields **filled out**?
- Are all values in correct **format**?
- Are all values **possible**?



On land?

Depths possible?

...



QC flag → fitness  
for purpose

## PUBLIC SERVICE ANNOUNCEMENT:

OUR DIFFERENT WAYS OF WRITING DATES AS NUMBERS CAN LEAD TO ONLINE CONFUSION. THAT'S WHY IN 1988 ISO SET A GLOBAL STANDARD NUMERIC DATE FORMAT. THIS IS **THE** CORRECT WAY TO WRITE NUMERIC DATES:

2013-02-27

THE FOLLOWING FORMATS ARE THEREFORE DISCOURAGED:

02/27/2013 02/27/13 27/02/2013 27/02/13  
20130227 2013.02.27 27.02.13 27-02-13  
27.2.13 2013. II. 27. 2<sup>7</sup>/<sub>2</sub>-13 2013.158904109  
MMXIII-II-XXVII MMXIII <sup>LVII</sup><sub>CCCLXV</sub> 1330300800  
((3+3)×(111+1)-1)×3/3-1/3<sup>3</sup> 2013 2013  
10/11011/1101 02/27/20/13 01237 2273  
<sup>2</sup><sub>5</sub> <sup>3</sup><sub>67</sub> <sup>4</sup><sub>8</sub> Mississipi



# QC Procedures Biological data - demo

[Download R script](#)

R-package documentation:

<https://github.com/iobis/obistools>