



DTO-BioFlow

Integration of biodiversity monitoring data into the Digital Twin Ocean

FAIR data, why and how

including general overview of data standards used in EMODnet Biology



Research data management











"Data is a precious thing and will last longer than the systems themselves."





Tim Berners-Lee, inventor of the World Wide Web



Why it is important

Order, chaos or organized chaos?





"We are all data providers and data users"

under the current system. Students in PhD programmes spend up to 80% of their time on 'data munging', fixing formatting and minor mistakes to make data suitable for analysis – wasting time and talent. With 400 such students, that would amount to a monetary waste equivalent to the salaries of 200 full-time employees, at minimum. So, hiring 20 professional data stewards to cut time lost to data wrangling would boost effective research capacity. Many top



We asked our respondents how much time they spend on each of the above tasks, and for each item, enter a number representing the percentage of time spent on each task relative to the other tasks on this list. The percentage values had to add up to 100.

More than 70% of researchers have tried and failed to reproduce another scientist's experiments, and more than half have failed to reproduce their own experiments. Those are some of the telling figures that emerged from *Nature*'s survey of 1,576 researchers who took a brief online questionnaire on reproducibility in research.



How is your data analysis going?

Can't understand the data

your future self, by Julien Colomb, CC-BY-NC, derived from .NORM Normal File For.

Baker, M. (2016). 1,500 scientists lift the lid on reproducibility. Nature, 533(7604).<u>https://doi.org/10.1038/533452a</u> Mons, B. (2020). Invest 5% of research funds in ensuring data are reusable. Nature, 578(7796), 491-491.<u>https://doi.org/10.1038/d41586-020-00505-7</u>



Personal benefits

Work more efficient & organised



More references & credits to your work

Career recognition

Collaborations

Moral obligations

Efficient use of public resources



Facilitates data finding & re-use

 \rightarrow New research & new insights



Better data leads to better research

- \rightarrow Improved decisions-making
- → Increased transparency & trust in science



Why it is important

RDM



VLIZ



Research data management













- Rich metadata & available online
- Persistent identifier
- Retrievable
- Accessible ≠ OPEN
- Authentication & authorisation steps
- Metadata should always be accessible

- Machine readable components
 - Open formats
 - Recognized standards
 - Linked data
 - Integration ready

- Data 'provenance'
- Data usage licence



DATA & METADATA

Who created the data? What the data files contain? When the data were generated? Where the data were generated? Why the data were created? How the data were generated?







It is a spectrum

≠ Open data

Open data is data that anyone can access, use & share



The circle of life Research data management







The circle of life data







RDM in practice!







Data Management Plan

What?

- How data will be handled **during & after** a research project
- Formal & "living" document

Why?



Save time

Avoid problems

.

Anticipate costs



FAIR by design





Data Management Plan



Content of DMP







Collecting data

- Metadata
- Controlled vocabularies
- Standards





Metadata and Documentation

Collecting



On three levels

- **Project (directory structure)**
- Files (naming conventions, READme)
- Data standards & vocabularies





Data standards

Global & multidisciplinary standards:

"Set of guidelines or rules that specify how data should be structured, formatted, and represented to ensure consistency, interoperability, and efficient data exchange"



ISO

International Organization for Standardisation

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:



Collecting



Data standards

Collecting

Global & multidisciplinary standards:

DwC = Darwin Core

EML

= Ecological Metadata Language

eventID	parentEventID	eventDate	decimalLongitude	decimalLatitude
site_1 zone_1 zone_2 zone_3 quadrat_1 transect_1 transect_2	site_1 site_1 site_1 zone_1 zone_2 zone_3	2019-01-02 2019-01-03 2019-01-04	54.7943	16.9425
id		occurrenceID	scientific	Name
$quadrat_1$		occ_1	Ulva rigi	da
$quadrat_1$		occ_2	Ulva lact	uca
$transect_1$ $transect_1$ $transect_2$		occ_3 occ_4 occ_5	Plantae Plantae Gracilari	a
${ m transect}_2$		occ_6	Laurenci	a

Basic Metadata Geographic Coverage Taxonomic Coverage Temporal Coverage Keywords Associated Parties Project Data Sampling Methods Citations Collection Data External links

Additional Metadata



Controlled vocabularies

- List of terms where each term means just one thing
- Ensure standardisation

Example Marine Sciences

	Identif	ifier↑ P	referred label ↑	Alternative label ↑	Definition ↑
National Oceanograp British Oceanograp Centre BODC	bhy Centre hic Data SDBIC	D OL09 sı tr	Iry weight biomass of biological entity pecified elsewhere per unit volume of ne water body	WaterDryWtBiom_BE007117	The mass measured after drying at elevated temperatures until a stable mass is reached, of an identified biological object described elsewhere in the metadata occurring in a given volume of any body of salt or fresh water.
	SDBIC	A OL07 bi u	sh-free dry weight biomass of iological entity specified elsewhere per nit volume of the water body	WaterAshFreeBiom_BE007117	The mass lost on ignition of an identified biological object described elsewhere in the metadata occurring in a given volume of any body of salt or fresh water.
• Biomass	SDBIC	V OL04 si tř	Vet weight biomass of biological entity pecified elsewhere per unit volume of ne water body	WaterWetWtBiom	The mass as caught of an identified biological object described elsewhere in the metadata occurring in a given volume of any body of salt or fresh water.
	SDBIC	B OL12 si tř	iomass as carbon of biological entity pecified elsewhere per unit volume of ne water body by computation	WaterCarbonBiomassConv	The carbon biomass, calculated from the cell counts using literature conversion factors, of an unspecified biological entity in a given volume of any body of salt or fresh water.







Taxonomic standard

Collecting

WoRMS provides the most authoritative list of names of all marine

species globally, ever published





Geographic standard

Standard list of marine georeferenced place names & areas





Lat: 56.94 Lon: -76.22 Click on the map to get featu









Data curation



Curation steps

Data exploration



Data enrichment

Data validation

Reproducible procedures

Keep raw data intact

Document transformation

Version Control

Document Quality Control procedures

Use Open formats





Data curation

X

Processing & analysing

Name	Phone	Birth date	Country
John Smith	445-881-4478	August 12, 1989	Belgium
Fitch, Marie	(876)546-8165	June 15, 72	US
Deere, Alan	+1-189-456-4513	11/12/1965	USA

Name	Phone	Birth date	Country
John Smith	445-881-4478	1989-08-12	Belgium
Marie Fitch	876-546-8165	1972-06-15	USA
Alan Deere	189-456-4513	1965-11-12	USA





File naming conventions

Recommendations:

Be consistent

...

- Use standards (e.g. YYYYMMDD)
- Do not use special characters or spaces
- Avoid words like 'draft', 'final'... use version numbers instead (v01, v02)

Examples of files without a naming convention: Meeting notes jan 10.doc Third_test.xls ProjectProposalFirstVersion.doc Project-data.xls

Examples of files with a naming convention:
20230110_OT_ODM_exercise1_v01.doc
20230110_OT_ODM_exercise1_v03.doc
20230109_OT_ODM_EvaluationResults.xls
20230109_OT_ODM_RDLC.jpg









Data archiving

Marine Data Archive - MDA

- = trusted data repository for marine,
- coastal and estuarine research
 - Closed repository for personal files & projects / collaboration
 - **Open** repository for data
 - publication

Preserving & archiving

Marine Data Archive

🔄 Intro Archive Manual Policy Register Contact FAQ



MDA... a secure, online system to **archive data files** in a **welldocumented manner**.

Log in

https://mda.vliz.be/







Searchable resources

Repositories

- Archiving and sharing
- Generic, discipline specific or institutional

Catalogue

• Description (rich metadata) of and link to data

Portal

- Archiving and sharing + interactive tools (visualisation, combining data, ...)
- Often thematic

Sharing

Searchable resources

VLIZ

IMIS

SeaDataNet

and more ...

Belgian Marine Data Centre

Integrated Marine Information System

- = catalogue with metadata information about:
- All datasets (open / not open)
- Related to marine and coastal research / topics
- Link to data

or contact person

https://www.vliz.be/imis

VLIZ

Reusing data

Provenance and documentation

Usage license and credit

© creative commons

	Data reu
IMIS	log in
Publications Institutes Persons Datasets Projects Maps	
report an error in this record]	<u></u>
∟ifeWatch observatory data: zooplankton observations by imaging (ZooScan) in the Belgian Part of th	e North Sea
Sitable as data publication	
Flanders Marine Institute (VLIZ), Belgium (2023): LifeWatch observatory data: zooplankton observations in the Belgian Part of the North Sea. https://doi.org/10.14284/584	Ownload Data
Previous versions (6) view	
Contact: data@vliz.be	
Access data	

Also accessible through: http://rshiny.lifewatch.be/zooscan-data/

Availability: (cc) BY This dataset is licensed under a Creative Commons Attribution 4.0 International License.

Notes: Images are available upon request via LifeWatch Belgium (info@lifewatch.be)

Description

In the framework of the Lifewatch marine observatory a number of fixed stations on the Belgian Part of the North Sea (BPNS) are visited on a monthly or seasonal basis using the RV Simon Stevin. A grid of nine stations covers the coastal zone and are sampled monthly. Eight additional stations, located further at sea, are sampled on a seasonal basis. This dataset contains zooplankton observations in the Belgian Part of the North Sea (BPNS) since 2012. Zooplankton is sampled by vertical WP2 net tows, samples scanned with ZooScanner and identification with plankton analyser software, followed by manual validation.

- Wilkinson, M., Dumontier, M., Aalbersberg, I. et al. The FAIR Guiding Principles for scientific data management and stewardship. Sci Data 3, 160018 (2016). <u>https://doi.org/10.1038/sdata.2016.18</u>

DTO-BioFlow

Integration of biodiversity monitoring data into the Digital Twin Ocean

THANKS!