**DTO-BioFlow**
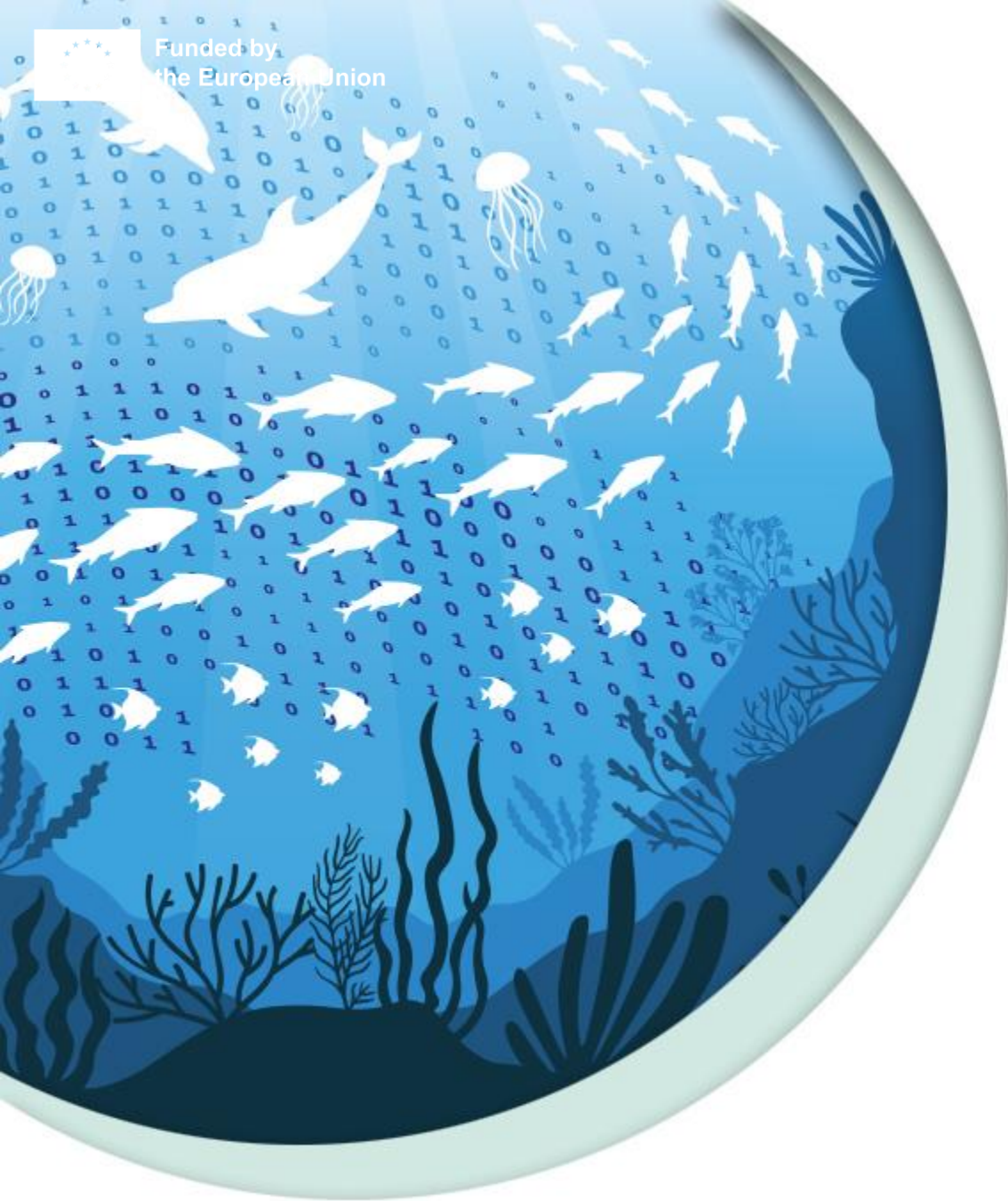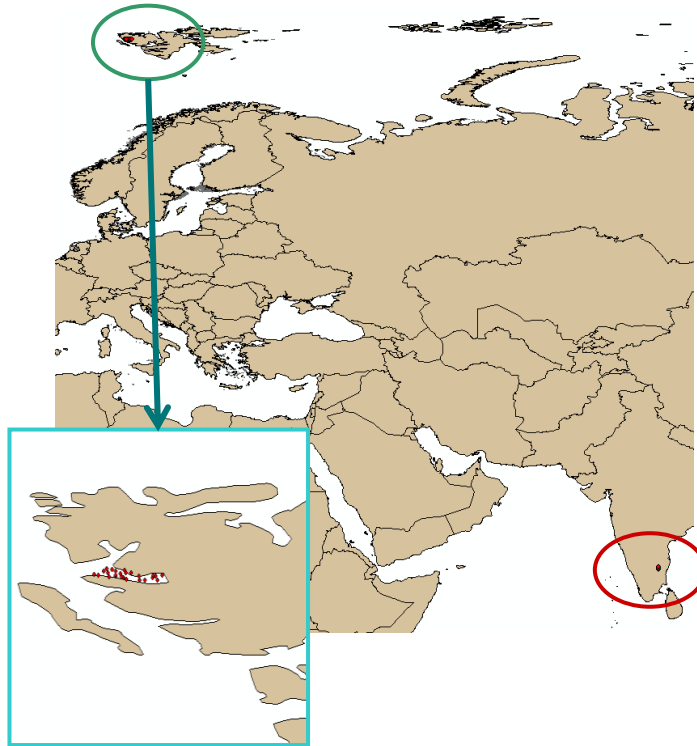
Integration of biodiversity monitoring data into the Digital Twin Ocean

# DTO-BioFlow data training workshop:

# Quality Control in data management

# Why is quality control important?

**"Monitoring in Kongsfjorden area"**

**"Monitoring in Belgian part of the North Sea"**

Latitude & longitude switched

"+" & "-" signs switched

# Why is quality control important?

| | Species names before quality control | | | | | Species names after quality control | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| | # Species | # Rare species | $H'$ | $1-D$ | ES50 | # Species | # Rare species | $H'$ | $1-D$ | ES50 |
| Rocky shore data | | | | | | | | | | |
| ANE | 219 | 15 | 4.63602 | 0.98777 | 38.11 | 187 | 11 | 4.45772 | 0.98509 | 36.25 |
| Arctic | 646 | 69 | 6.00024 | 0.99666 | 46.33 | 378 | 44 | 5.38261 | 0.99403 | 43.67 |
| Mediteranean | 1,120 | 238 | 5.74091 | 0.99342 | 43.35 | 834 | 159 | 5.49015 | 0.99105 | 41.74 |
| North Sea | 251 | 29 | 4.50662 | 0.98424 | 35.89 | 163 | 25 | 3.95956 | 0.97469 | 30.14 |

*"From 6,172 unique taxon names […] to 4,525, mostly due to spelling variations and synonymy."*

*" … Such [taxonomic] quality control is highly needed, since a misspelled or obsolete name could be compared to the introduction of a rare species, with adverse effects on further (biodiversity) calculations…"*

***Source***: Vandepitte *et al.* (2010)

# What needs to be QCed

# What needs to be QCed

- Adherence to chosen standard

- Content format (dates, coordinates, etc.)

- Unique value fields (IDs)

- Duplicated records and redundant information

- Impossible values and outliers

- Mandatory data

- Completeness (units, geodetic datums, etc.)

- ...

# Quality control tools

There are many different tools you can use to help you perform quality checks on your data. These include:

- R package obistools
  - Checking required fields, coordinates, depth, outliers, dates, …
- R package and function Hmisc:: describe
  - Summary statistics
- LifeWatch & EMODnet Biocheck
- Lifewatch data services
- WoRMS taxon match tool
- GBIF data validator
- Excel Conditional Formatting tool - identify duplicated data
- …

# The LifeWatch & EMODnet Biocheck tool

≋ Shiny app:
https://rshiny.lifewatch.be/BioCheck/

≋ R package:
https://github.com/EMODnet/EMODnetBiocheck

# The LifeWatch & EMODnet Biocheck tool

IPT resource URL or
Darwin Core Archive file

# The LifeWatch & EMODnet Biocheck tool



≋ Explore dataset

# The LifeWatch & EMODnet Biocheck tool

≋ Identify issues

   ≋ List of checks under "About"

≋ Check and deal with all error and warning messages!

# The LifeWatch & EMODnet Biocheck tool

≋ Explore records with issues

# Demo

**Biocheck QC tool**



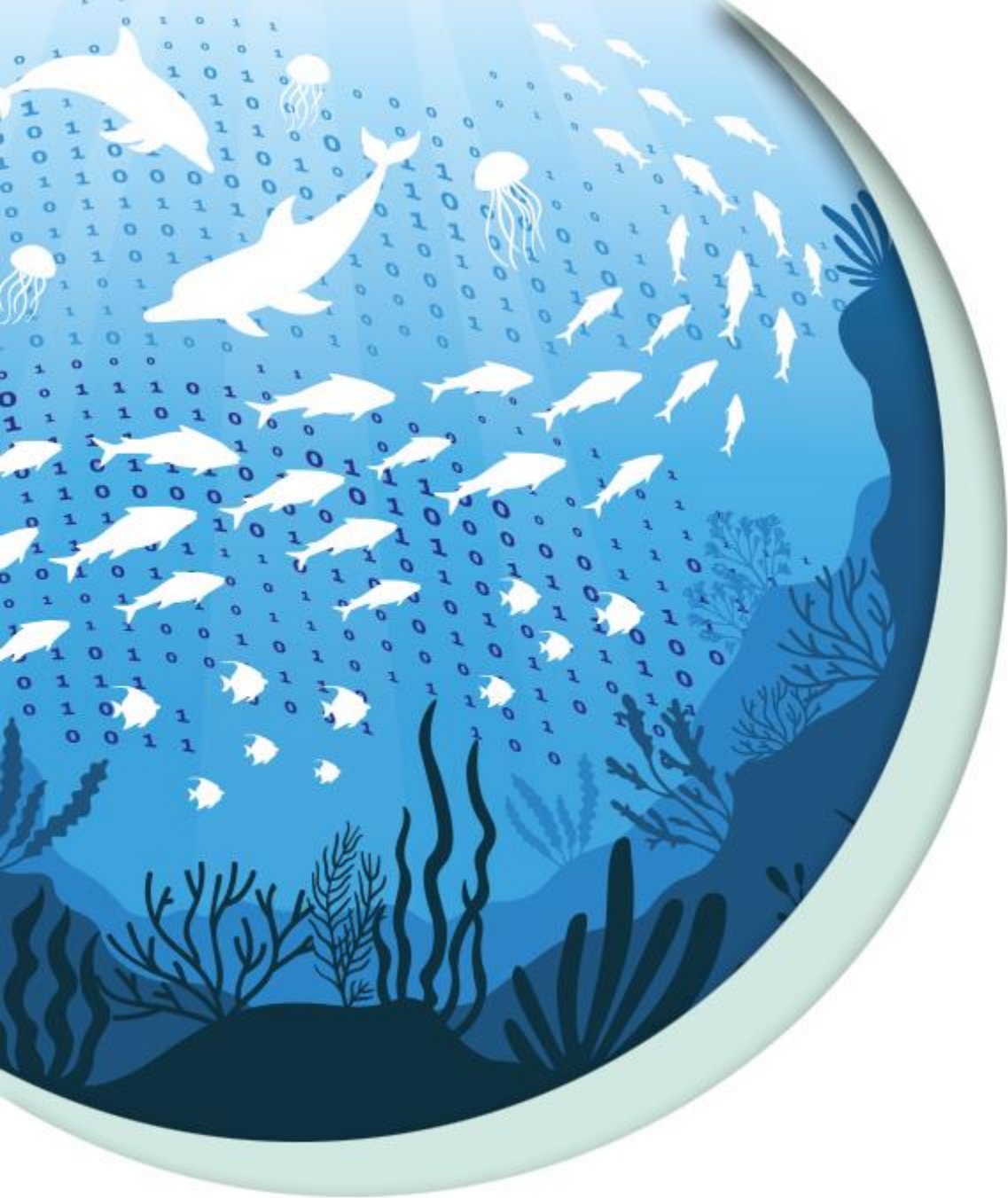https://rshiny.lifewatch.be/BioCheck/

# Relevant sources

≋ Vandepitte *et al.* (2010). Data integration for European marine biodiversity research: creating a database on benthos and plankton to study large-scale patterns and long-term changes. *Hydrobiologia 644: 1-13*

≋ The LifeWatch & EMODnet Biocheck tool

≋ The GBIF data validator

THANKS!