# WP3
# developing & testing analytical Blue Cloud workbenches for generating highly qualified data collections (EOVs)

Dominique Obaton

IFREMER

## What ?

**TOOLS**

To build sustained pipelines that will enable integration and combination of datasets from various sources for further analysis as metadata homogenisation, duplicate management, quality control

**and DATA COLLECTIONS as results**

To illustrate what is possible through examples

To get first user feedback from "upgraded" datasets proposed. For a continuous improved loop

**Blue-Cloud2026**

## Why?

→ There is a need to improve different datasets: EMODnet, Copernicus, SeaDataNet, World Ocean Database (WOD2018), ELIXIR …

- ❖ Metadata
- ❖ Duplicates
- ❖ Quality control (QC)

To ensure whatever data collection an expert / scientist /data manager uses (can be a part of different datasets), QC is good & homogenous.

→ This will increase confidence of datasets usage

- ❖ whatever the source

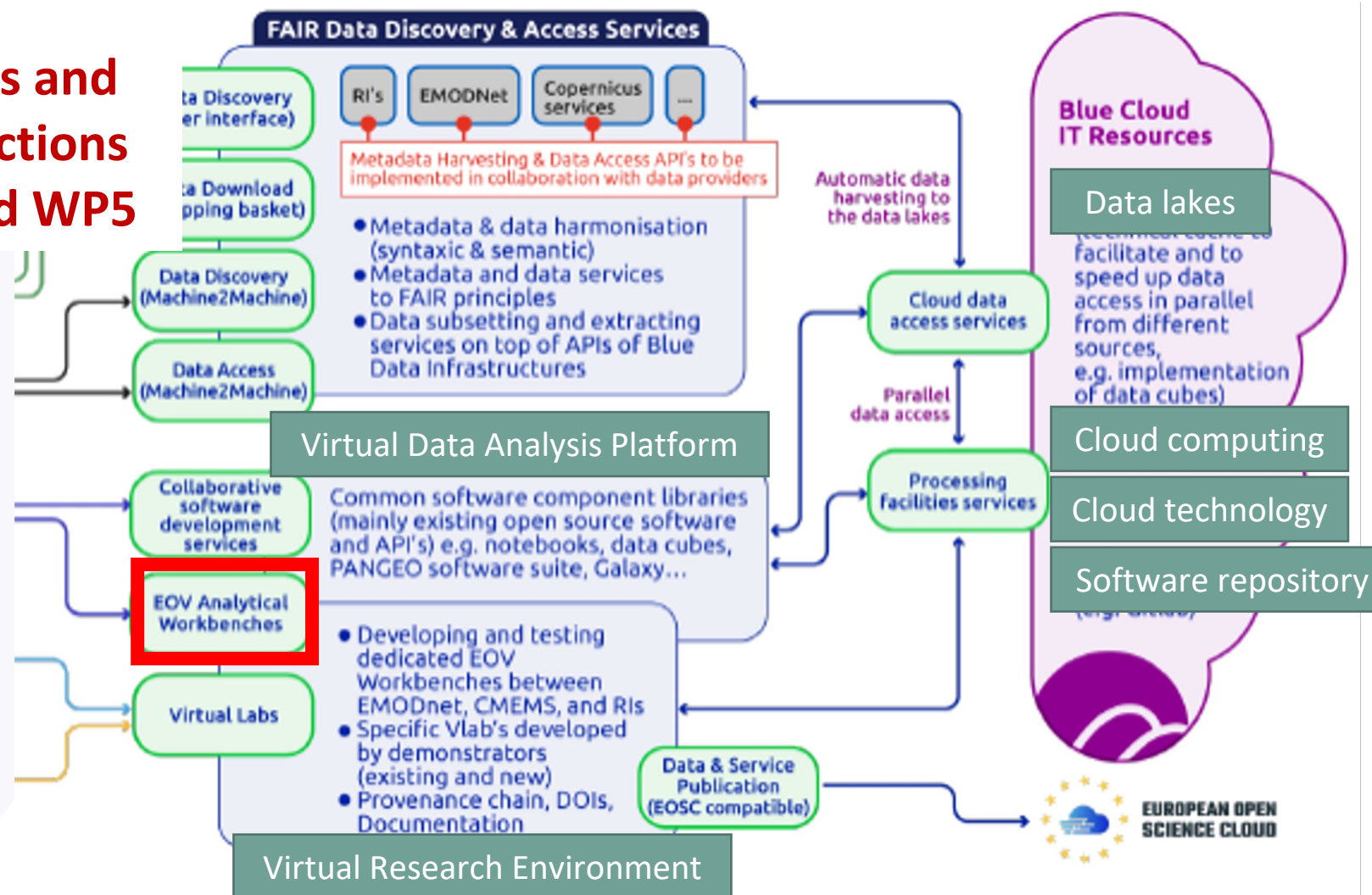→ Then this will increase the number of users

**How?**

*CHALLENGES*

**Big challenges and strong interactions with WP2 and WP5**

→ Possible thanks to a Blue Cloud IT resources with

- Data lakes
- Cloud computing
- Cloud technology
- Software repository

→ On a virtual environment

- Virtual data analysis platform
- Virtual research environment



N.B. similar development in the FAIR-EASE project for the whole Earth's ecosystem. **To work in strong interactions**

eosc | Blue-Cloud2026

**What we are now able to do**
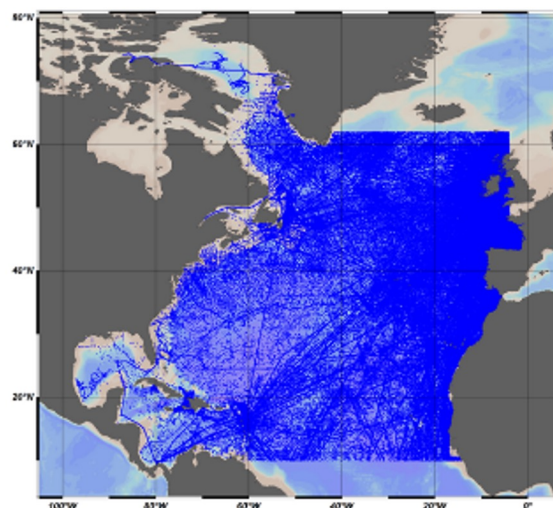
Current constraints on a PC of 16GB memory

→ North Atlantic ocean

→ period: 1950-2019

→ SeaDataNet + Copernicus data collections

~ **5 GB**

→ **Limit of data volume to work on PC**

→ **Difficult to make analysis on such volume**

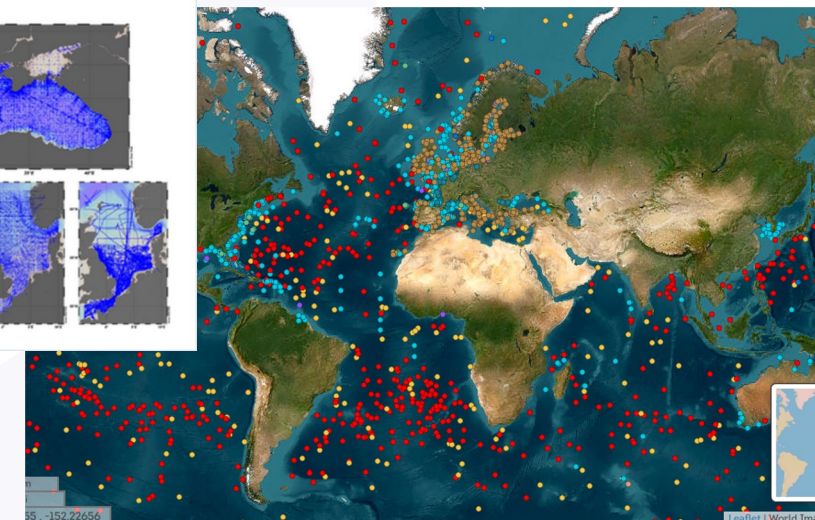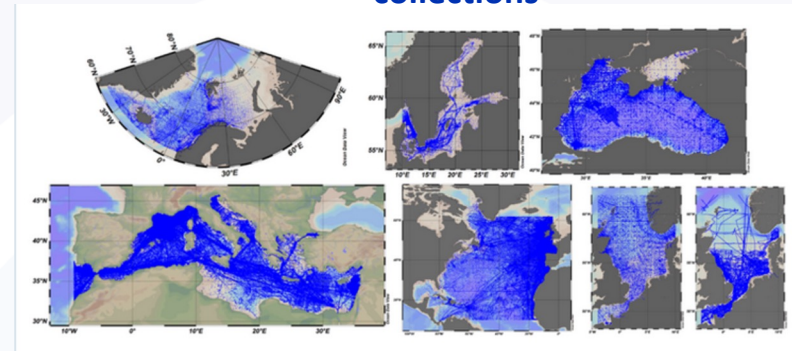**What we want to do (final objective)**

Global ocean

period: 1950-current time

SeaDataNet + Copernicus + WOD2018 datasets

= estimate **>20 GB**

→ **Possible thanks to Blue Cloud IT resources (Data lakes, processing services ...)**

→ **That will offer the possibility to have the data at one place and make the analysis at the same place for large and important data collections**

**Why?** **→ BENEFIT BLUE CLOUD 2026**

→ Large data collections homogenised with a good quality control
  ❖ Not possible up to now
  ❖ Thanks to Blue Cloud IT including big data techno

→ Workbenches or pipelines developed for several variables (EOVs)
  ❖ Reusable for other / updated datasets
  ❖ Depending on specific validation wanted (e.g. open sea vs coastal)
  ❖ Portable on other e-infrastructure than Blue Cloud after the project

→ Workbenches are prototypes
  ❖ Than can be used in operational oceanography
  ❖ That will contribute to the building of the Digital Twin of the Ocean

## Who?

**13 partners altogether**

Leader: IFREMER

Partners:  AWI, CMCC, SSBE, INGV, Oceanscope, Pokapok, HCMR, OGS, SMHI, ETHZ, SU, EMBL

Through 4 tasks

For 209 months of people

During the whole project

**T3.1 coordination, definition and design**

*Lead: Ifremer* Dominique Obaton

*Partners: AWI* Reiner Schlitzer*, CMCC* Massimiliano Drudi*, SSBE* Julia Vera

**T3.2 EOV workbench for physics: temperature and salinity**

*Lead: INGV* Simona Simoncelli

*Partners: Ifremer* Christine Coatanoan*, Oceanscope* Tanguy Szekely*, Pokapok* Jérôme Gourrion*, HCMR* Sissi Iona

**T3.3 EOV workbench for eutrophication: chlorophyll, nutrients, oxygen.**

*Lead: OGS* Alessandra Giorgetti

*Partners: Ifremer* Julie Gatti*, Pokapok* Virginie Racapé*, SMHI* Lindh Markus

**T3.4 EOV workbench for ecosystems (genomic)**

*Lead: ETHZ* Meike Vogt

*Partners: SU* Jean-Olivier Irisson*, EMBL* Stéphane Pesant

**T3.2 EOV workbench for physics: temperature and salinity**

*Lead: INGV* Simona Simoncelli

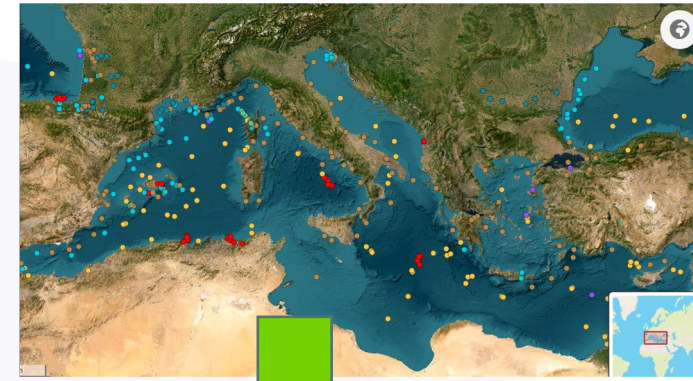*Partners: Ifremer* Christine Coatanoan, *Oceanscope* Tanguy Szekely, *Pokapok* Jérôme Gourrion, *HCMR* Sissi Iona

→ Datasets: SeaDataNet + Copernicus physics + WOD2018 datasets

→ Variables: Temperature & salinity

→ Methods: ODV, ISAS, MinMax

→ Start with the Mediterranean Sea, go to the global ocean



Mediterranean Sea Plot: 10 days, 400 data.
All database, 1 million data

Plot: 10 days, 16500 data. All database, 30 millions data

**T3.3 EOV workbench for eutrophication: chlorophyll, nutrients, oxygen**
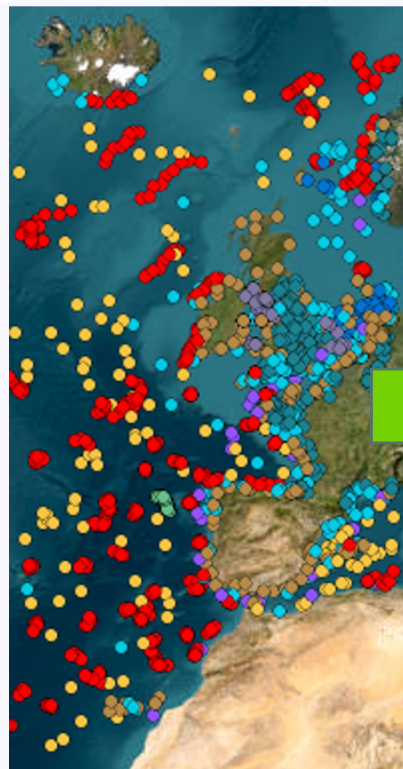
*Lead: OGS* Alessandra Giorgetti

*Partners: Ifremer* Julie Gatti*, Pokapok* Virginie Racapé*, SMHI* Lindh Markus

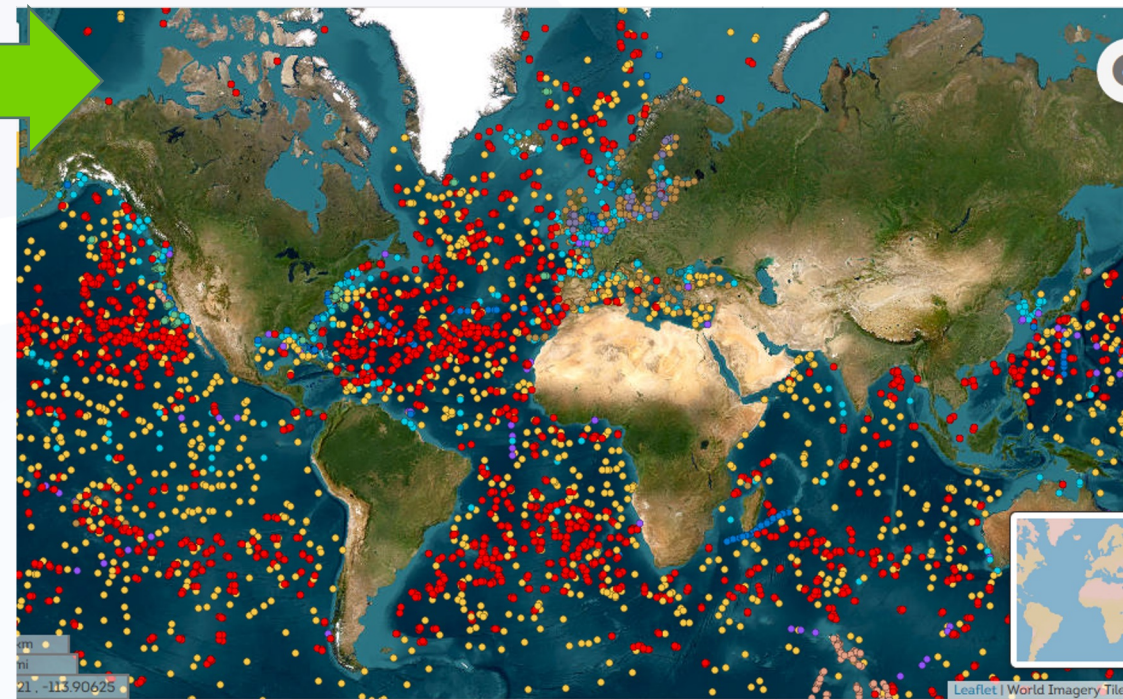→ Datasets: EMODnet chemistry + Copernicus +WOD2018 datasets

→ Variables: chlorophyll, nutrients, oxygen

→ Methods: ODV, Copernicus procedures

→ Start with the North East Atlantic, go to the global ocean



North East Atlantic
Plot: 10 days, 1790 data.
All database, 2,4 millions data



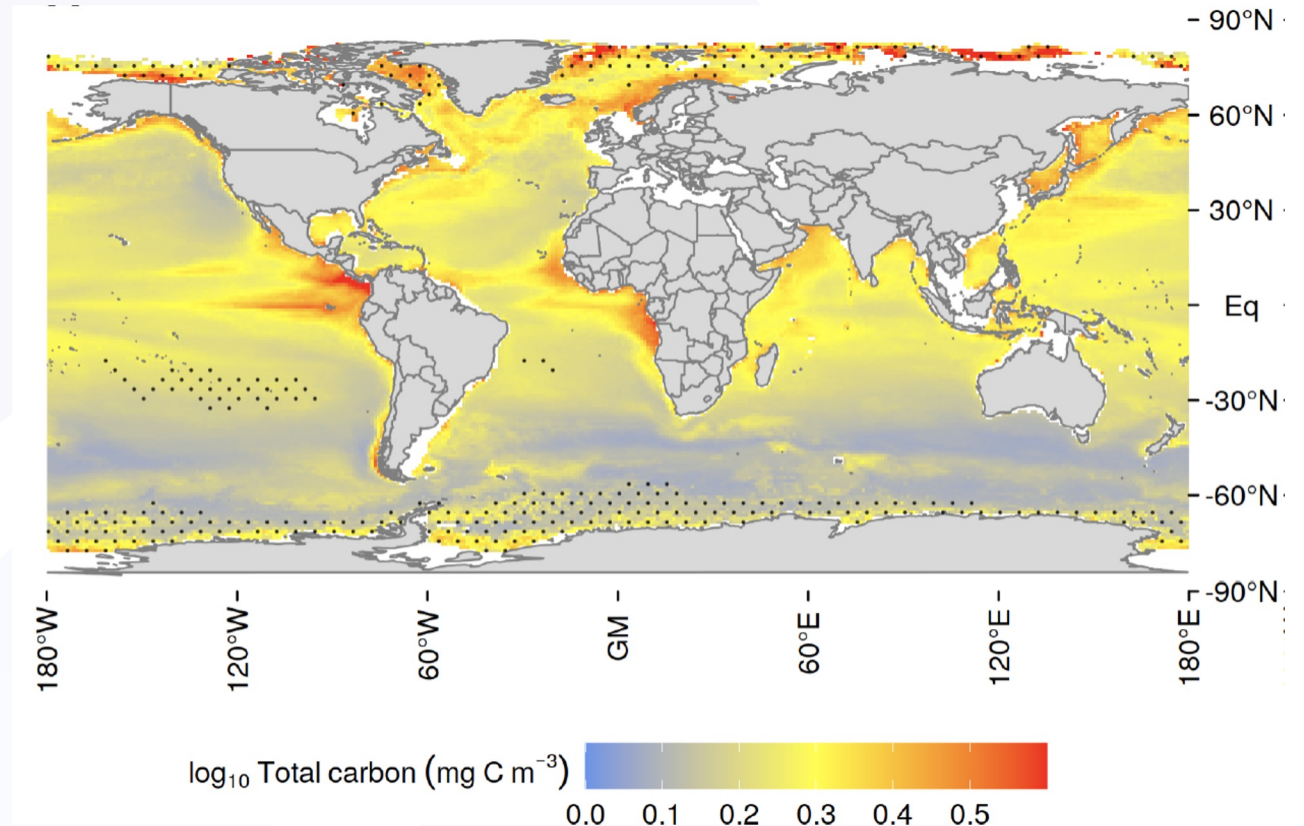Plot: 10 days, 16500 data. All database, 30 millions data

**T3.4 EOV workbench for ecosystems (genomic)**
*Lead: ETHZ Meike Vogt*
*Partners: SU Jean-Olivier Irisson, EMBL Stéphane Pesant*

→ Datasets: large collections of plankton observations EMODnet biology/EurOBIS, ELIXIR

→ Variables: plankton diversity, biogeography, biomass and relative abundance

→ Methods: Mgnify, EcoTaxa

→ Global ocean

Aggregation and harmonization of large collections of plankton observations based on traditional counts, quantitative imaging and genomic methods
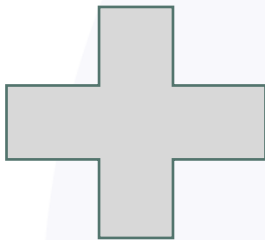
**T3.1 coordination, definition and design**

*Lead: Ifremer* Dominique Obaton

*Partners: AWI* Reiner Schiltzer*, CMCC* Massimiliano Drudi*, SSBE* Julia Vera

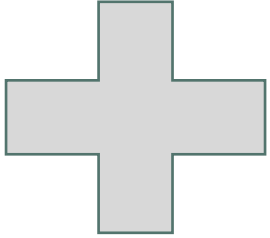Within the workpackage

→ Coordination

→ Common points, developments between tasks, especially methods in workbenches 1 and 2

With other workpackages

→ Strong dependency of WP2 (data lakes, Virtual Data Analysis Platform -VDAP) and WP5 (Virtual Research Environment -VRE)

→ Highly validated datasets from workbenches + some applications/use for/within the use cases of developed in WP4

→ Follow and liaise with project coordination (WP1)

→ Contribution to the communication – outreach & engagement (WP6)

→ Contribute to the exploitation and sustainability (WP7) *Big challenge*

| Deliverable (number) | Deliverable name & description | WP number | Delivery date |
|---|---|---|---|
| D3.1 | **1st release of aggregated and harmonised EOV datasets**: preliminary aggregated and harmonised EOV datasets obtained thanks to each workbench. | WP3 | M20 |
| D7.1 | **Individual Exploitation Plans of Workbenches**: A compilation of the individual Exploitation Plans of each of the Workbenches set up in WP3 | WP7 | M24 |
| D3.2 | **Final release of open aggregated and harmonised EOV datasets**: final and open aggregated and harmonised EOV datasets stamped BC2026 obtained thanks to the workbenches. | WP3 | M42 |
| D3.3 | **Workbenches and tools (notebooks, model, containers with instructions)**: Portable prototypes for operational services and data infrastructures. | WP3 | M42 |

**Task 3.2 workbench for physics: temperature and salinity**

**Simona Simoncelli, INGV**

First task meeting (remote) before Technical Scientific Committee: 28 February at 14:00 CET

**Task 3.3 workbench for eutrophication: chlorophyll, nutrients, oxygen**
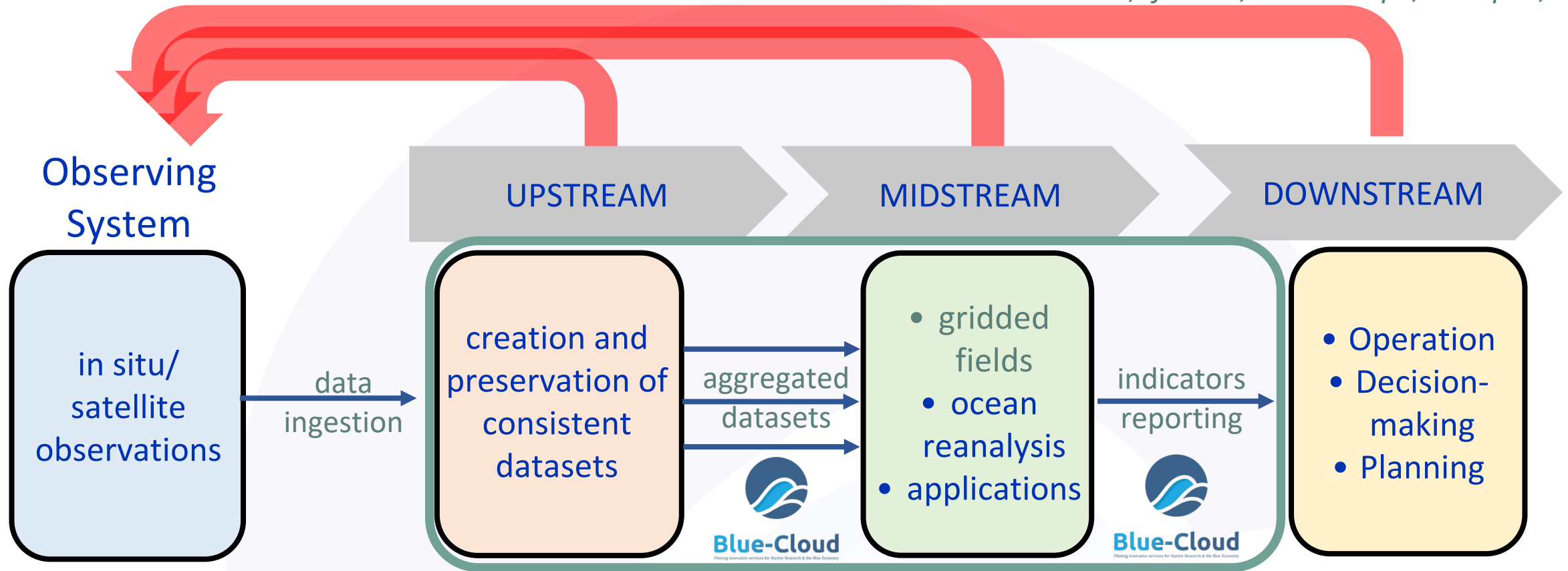
**Alessandra Giorgetti, OGS**

First task meeting (remote) before TSC: 3rd March at 11.00 am CET

**Task 3.4 workbench for ecosystems**

**Meike Vogt, ETHZ**

First in person meetings (a) between ETHZ-EMBL: January 10-13, 2023; (b) between ETHZ-SU: February 6-9, 2023
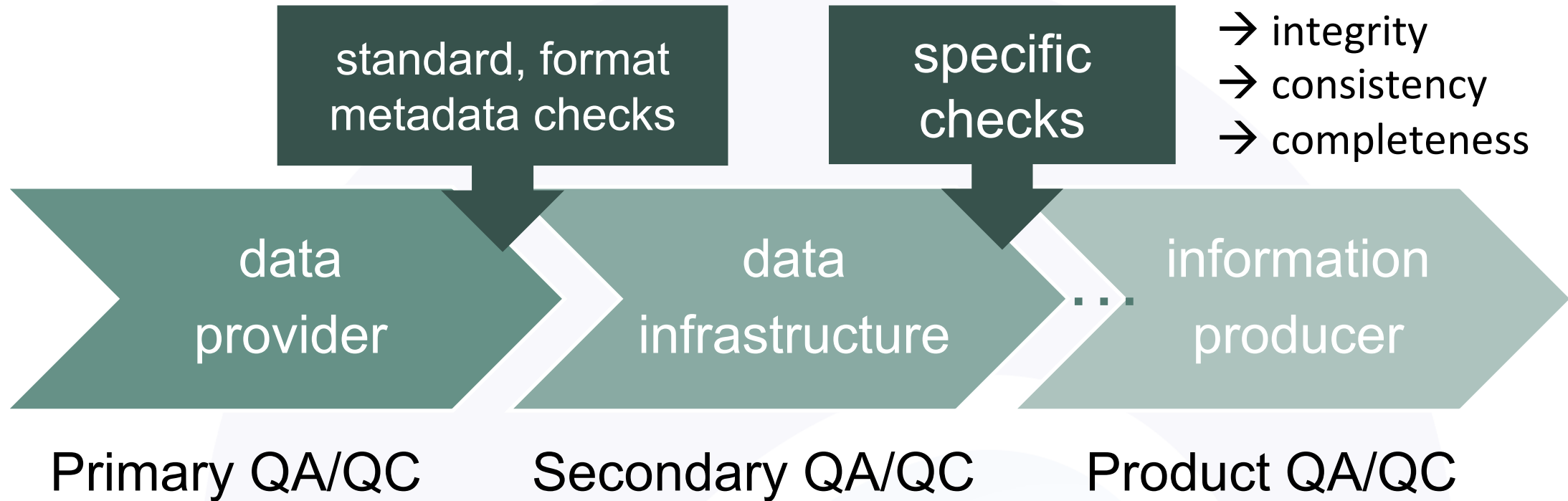
# Integration of SeaDataNet T&S data with other data sources (Copernicus, NOAA-WOD2018) → **duplicates removal**

*INGV, Ifremer, Oceanscope, Pokapok, HCMR*
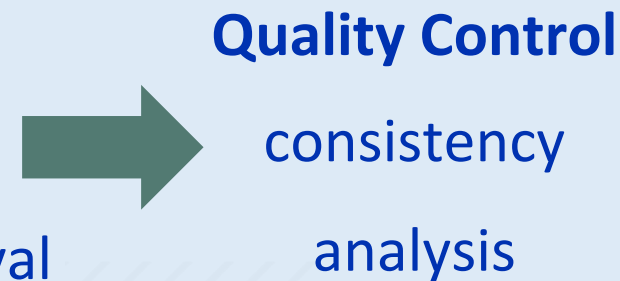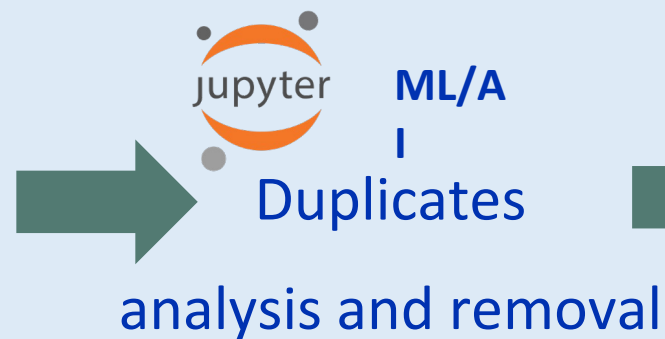
**Blue Cloud T&S dataset**
and
input datasets

**gridded fields**
climatologies

**WP4**

**Vlab 4
Ocean Indicators**
Ocean Heat Content

**WP7**

*OGS, Ifremer, Pokapok, SMHI*

**Ambition and scope:**
- Providing a toolbox, configurable by the user, to **build customizable datasets from different combinations of collections** or QC procedures and to compute various statistical parameters to assess the consistency of observed and derived quantities.

**How:**
- Definition and implementation of an efficient production **workflow that will merge multi-source datasets** to obtain an integrated and most complete dataset for the North East Atlantic.

- **Comparison** of EMODnet chemistry and Copernicus marine in-situ **QC methods** to provide a set of procedures depending on regions and/or user's aims.

- **Testing** the statistical parameters to i) evaluate the consistency of the initial input dataset and ii) to compare the data selection obtained after applying different QC strategies.

*OGS, Ifremer, Pokapok, SMHI*

- **Currently** several datasets in different data infrastructures with their own QC/QA are synchronised manually. These need to be further integrated within the Blue-Cloud 2026 data infrastructure.

**Methods:**
- ODV, DIVAnd, Copernicus procedures.

**Description:**
- The task will aim at **generating harmonised and validated EOV data collections** for Chlorophyll, nutrients and oxygen integrating several datasets released from different EU and non-EU data infrastructures.

- The information associated with the observations, metadata, will be **mapped** into a common schema and the quality control procedures will be **exchanged**. A dedicated protocol will be jointly set out to identify and handle the potential duplicate observations.

- The Workbench will be **developed and tested for the Northeast Atlantic Sea** with the aim of further **extending** it to the global ocean and will be made available for implementation by other data infrastructures.

EOV workbench for eutrophication: *chlorophyll, nutrients, oxygen*

**eosc** | Blue-Cloud2026
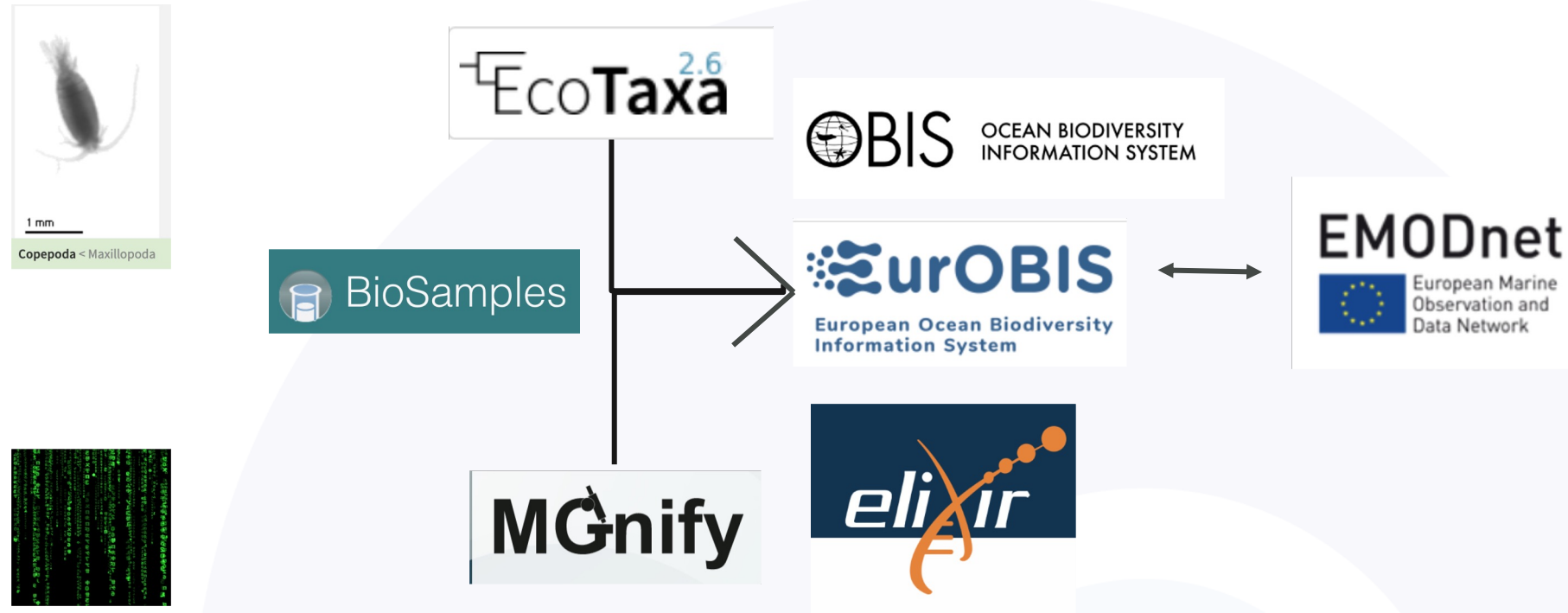
*OGS, Ifremer, Pokapok, SMHI*

- ODV, DIVAnd, Copernicus procedures **will need further developments** to support the interaction with chosen technologies to handle the large volumes of data the WorkBenches will have to deal with (Task 3.1)

- IT functionalities and FAIRness of **BlueCloud services** (WP2 & WP5) will need to be used for dataset integration and analysis.

- **Interoperability services** of the data infrastructures, common vocabularies and brokering services will be needed to allow dataset aggregation and harmonisation.

- To be further discussed at the **task meeting planned for March 3rd, 2023 at 11.00 am CET** & in Technical and Scientific Committee (TSC) meeting in Amsterdam fixed for 28-29 March 2023

eosc | Blue-Cloud2026

*ETHZ*, *SU*, *EMBL*

**Aims:**

● Improve **the availability, quality and interoperability of large collections of plankton observations** based on traditional counts, quantitative imaging and genomic methods available from the EMODnet/EurOBIS and ELIXIR data infrastructures

● Develop **a generic machine learning modelling pipeline** in Blue-Cloud 2026 to generate high-quality interpolated maps of the global distribution of plankton biogeography and diversity (binary data; presence-absence), biomass (quantitative data; e.g. biomass) and community composition (relative abundance; e.g. percent reads)

● Provide objective, multi-criteria quality assessments of output layers of added value based on based on multi-model ensemble projections
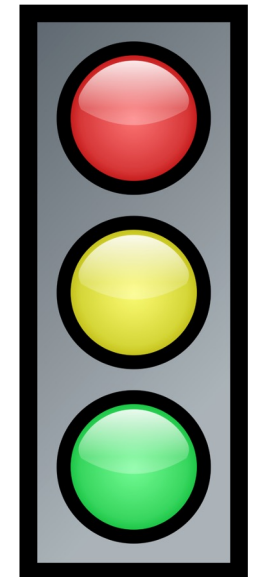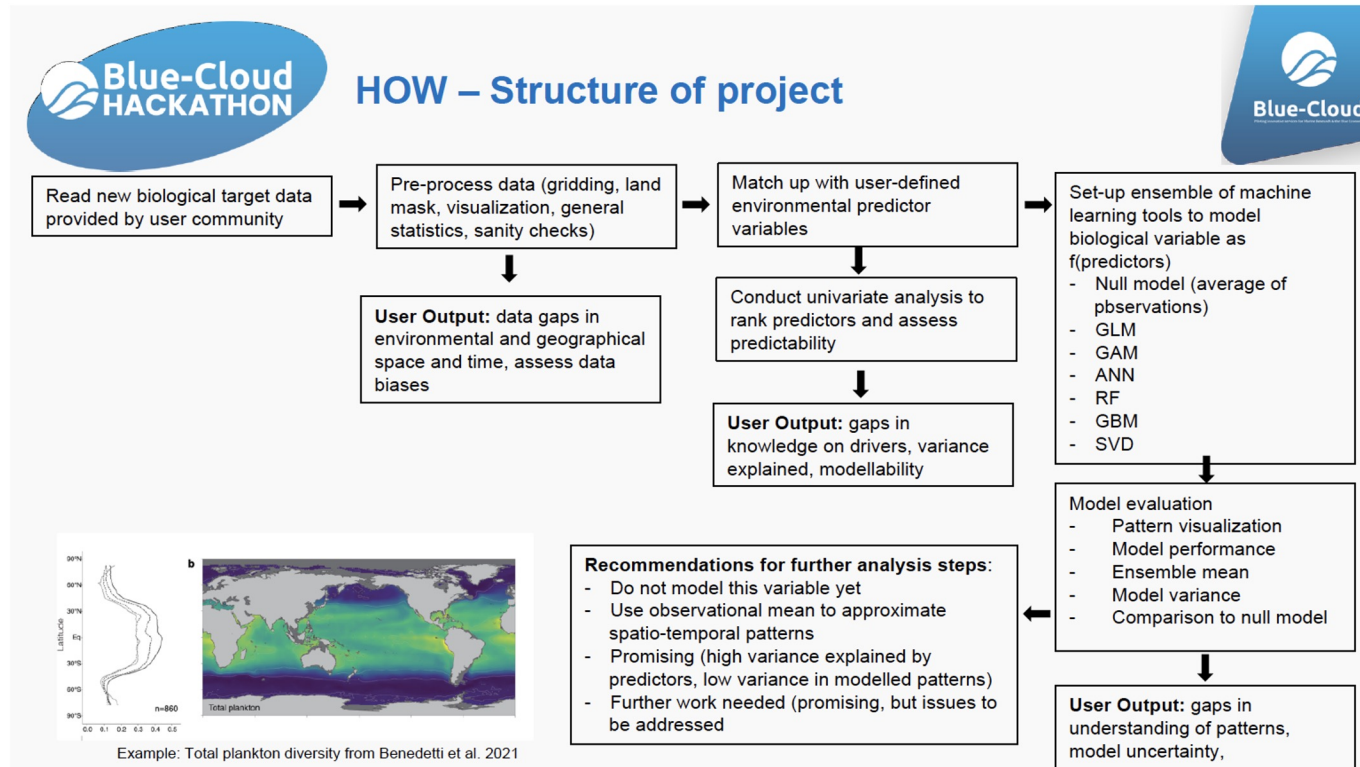
**Data and Methods:**

● Repositories: EurOBIS, OBIS, GBIF - plankton data from traditional methods, imaging and 'omics
● Data resources/analysis pipelines/infrastructures: MGnify, Ecotaxa, BioSamples (EMODnet)
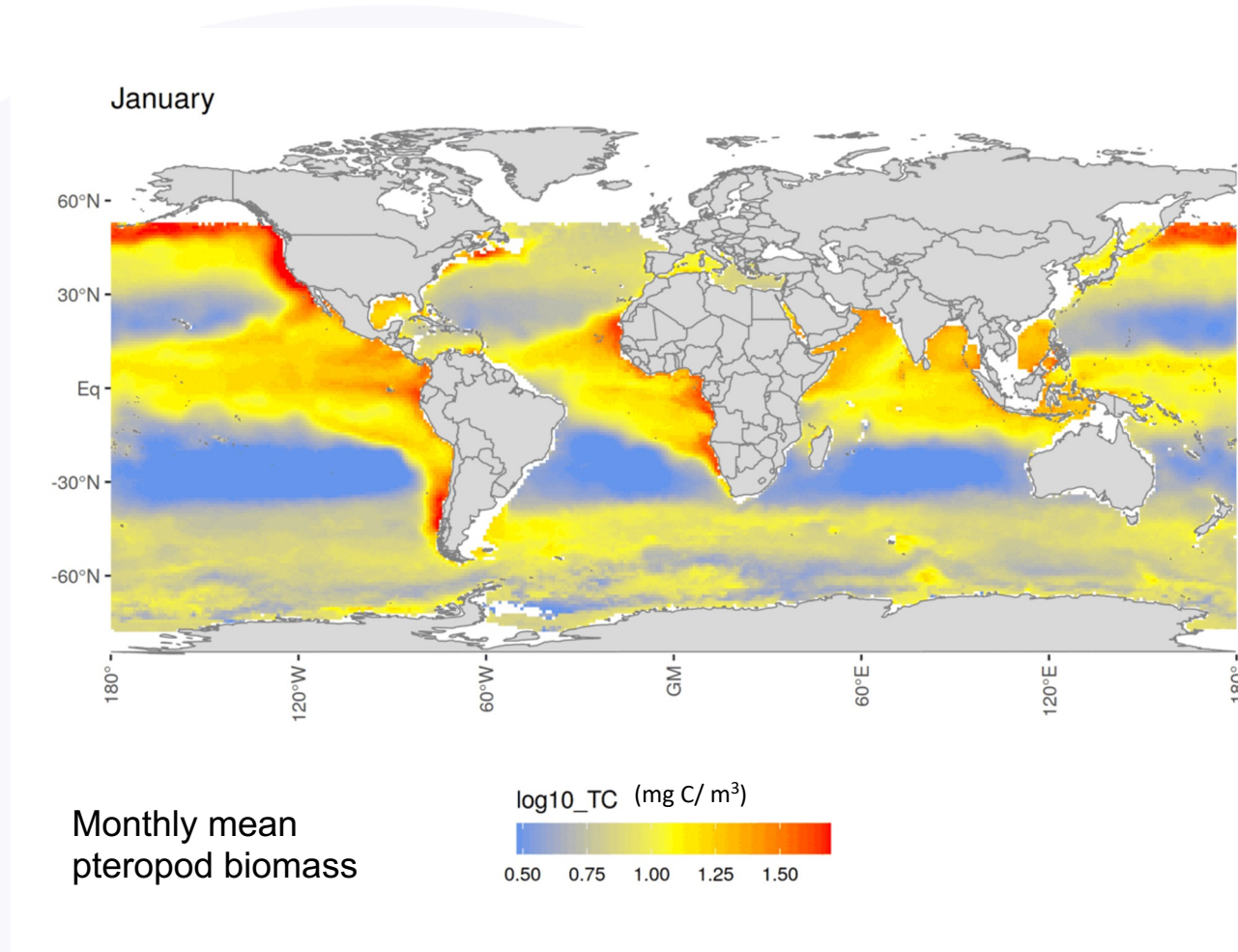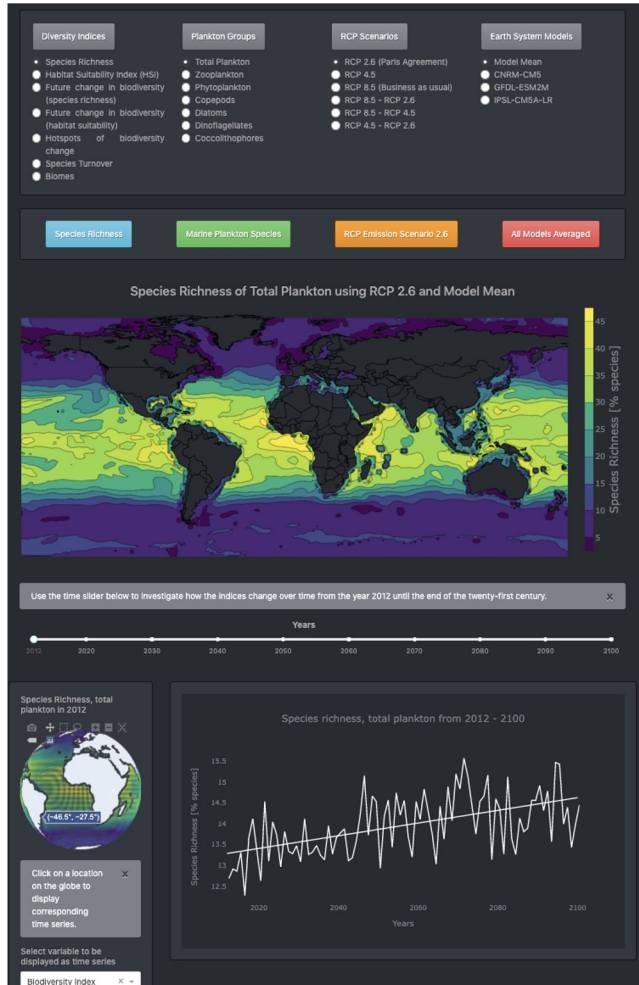● Methods: Species distribution modelling (common machine learning tools)

ETHZ, *SU, EMBL*



**Facilitated data access and increased biological data availability:**

- Improved **exchange of interoperable plankton data** between EcoTaxa/MGnify and key global data repositories (e.g. EurOBIS) and research infrastructures (e.g. EMODnet)

- **Provenance metadata curation pipeline** (BioSamples) for omics data

ETHZ, SU, EMBL

**Prototype modelling pipeline developed during BlueCloud hackathon 2021 (proof-of-concept):**

- **Automatized, ensemble-based modelling/machine learning pipeline** for the projections of binary (presence-absence), continuous quantitative (biomass) and proportional (percent genetic reads) data

- Rigorous quality-control of output products using **a multi-metric quality assessment**

**Explorable layers of added value for ecosystem monitoring, conservation and policy making**