# Data life cycle for Life Science

**Stéphane PESANT**

**Senior marine biocurator**

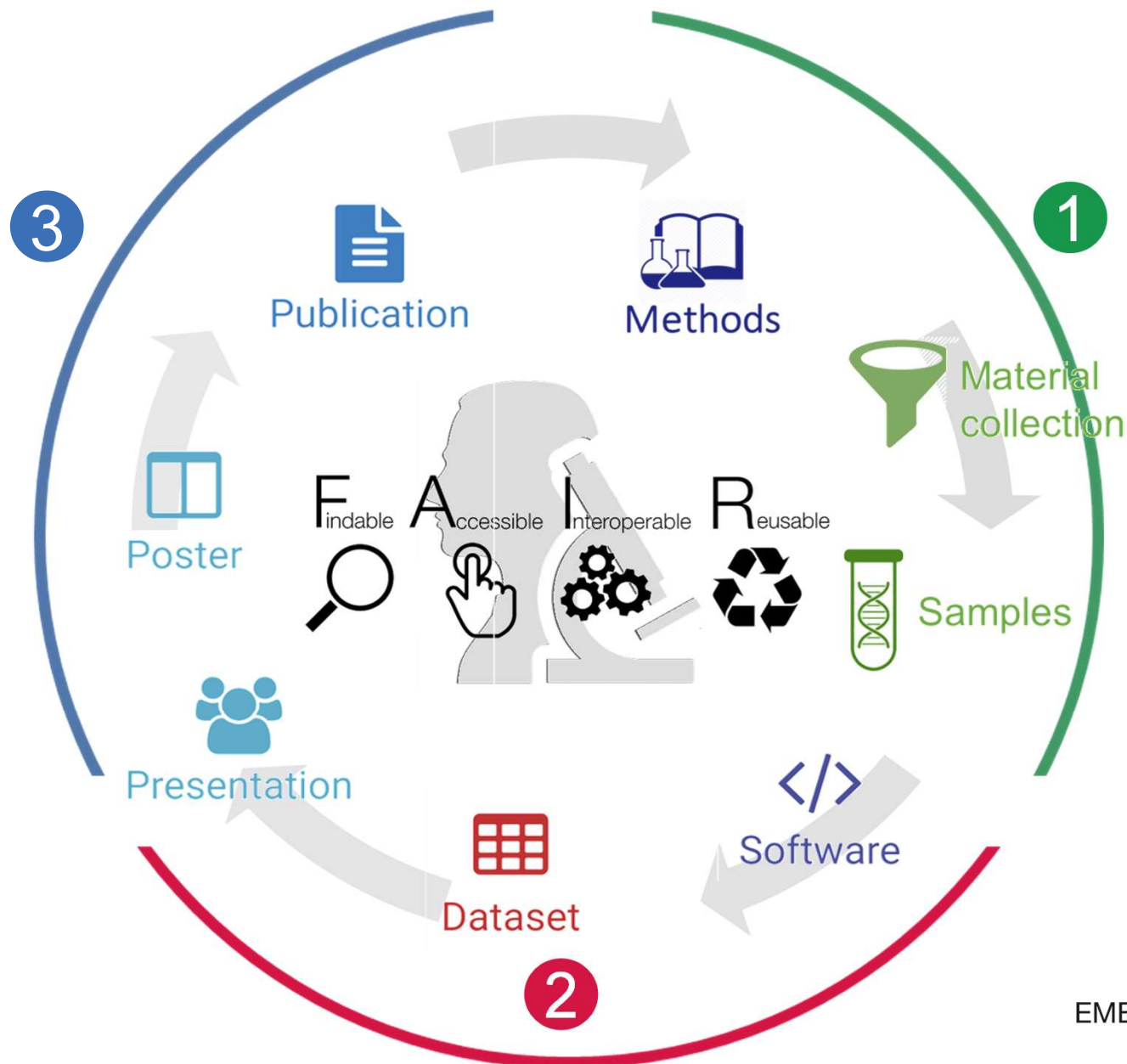SZN Seminar, Naples, 2023-02-09

EMBL-EBI

# Data life cycle

Sampling
Best Practices

Data Sharing
Best Practices

Publication
Best Practices

# Sampling best practices

Examplar initiatives - developing and experimenting with best practices

- Tara Oceans

- MicroB3

- EMO-BON

- AtlantECO

EMBL-EBI

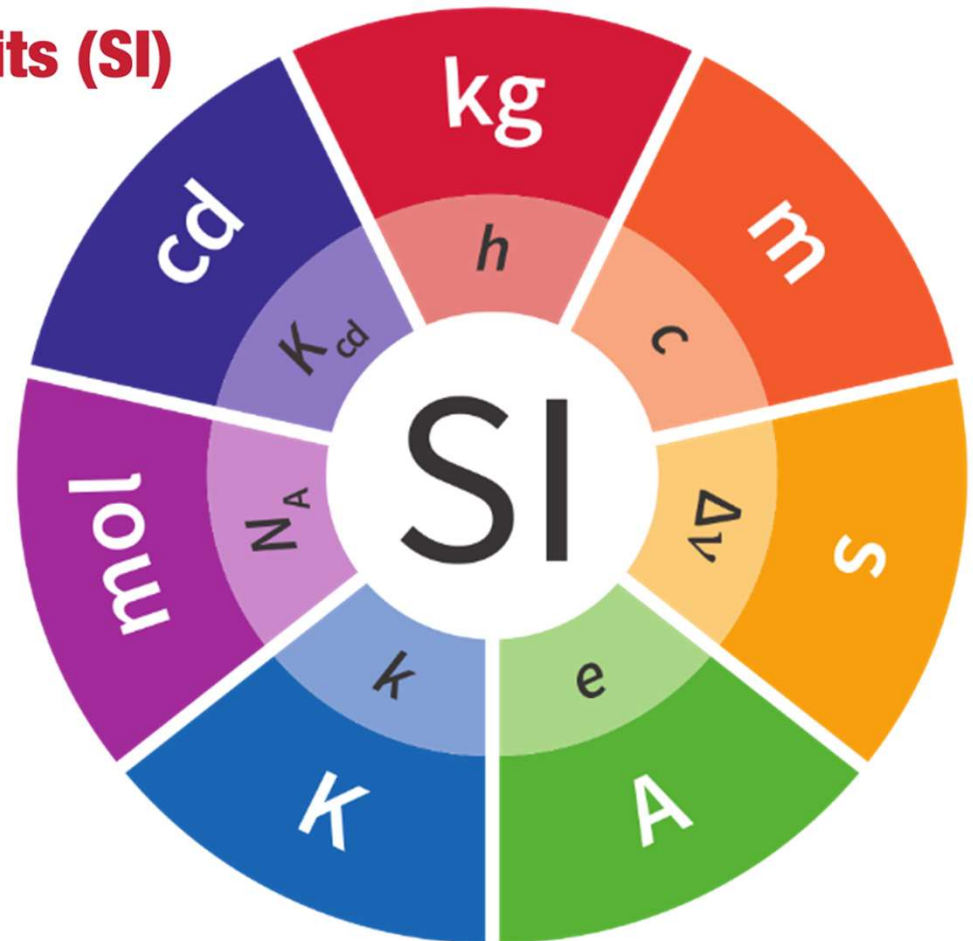# System of protocols



## International System of Units (SI)

### SI Base Units

| Base Quantity | Name | Symbol |
|---|---|---|
| Length | meter | m |
| Mass | kilogram | kg |
| Time | second | s |
| Electric current | ampere | A |
| Thermodynamic temperature | kelvin | K |
| Amount of substance | mole | mol |
| Luminous intensity | candela | cd |

### SI Derived Units

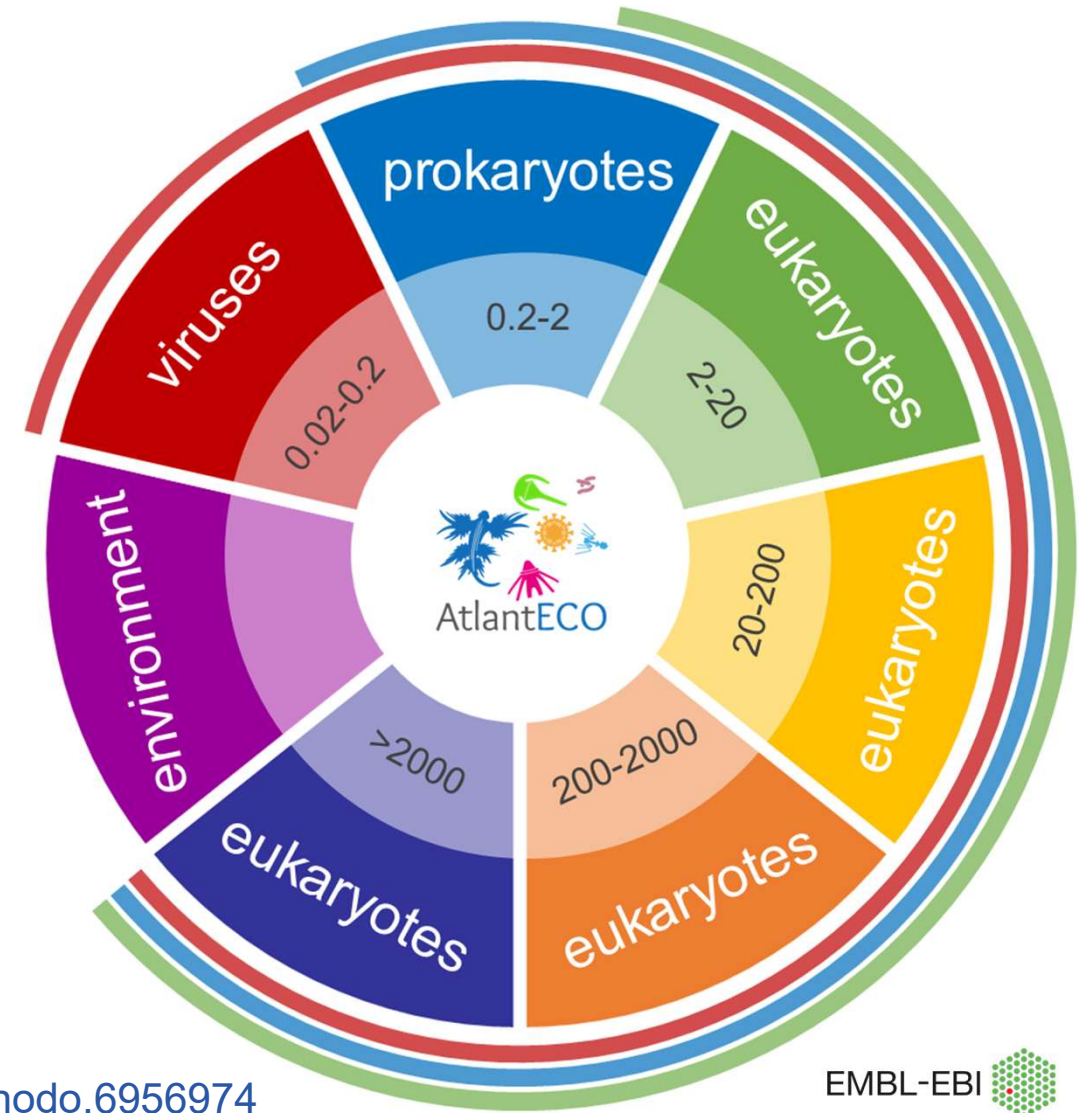| Derived Quantity | Name | Symbol | Equivalent SI units |
|---|---|---|---|
| Frequency | hertz | Hz | $s^{-1}$ |
| Force | newton | N | $m \cdot kg \cdot s^{-2}$ |
| Pressure | pascal | Pa | $N/m^2$ |
| Energy | joule | J | N·m |
| Power | watt | W | J/s |
| Electric charge | coulomb | C | s·A |
| Electric potential | volt | V | W/A |
| Electric resistance | ohm | Ω | V/A |
| Celsius temperature | degree Celsius | °C | K* |

*Unit degree Celsius is equal in magnitude to unit kelvin.

EMBL-EBI

# System of protocols

- Size

- Taxonomy

- Across size-fractions

 

- Genomics

- Transcriptomics

- Proteomics

- Metabolomics

- Phenomics



https://doi.org/10.5281/zenodo.6956974
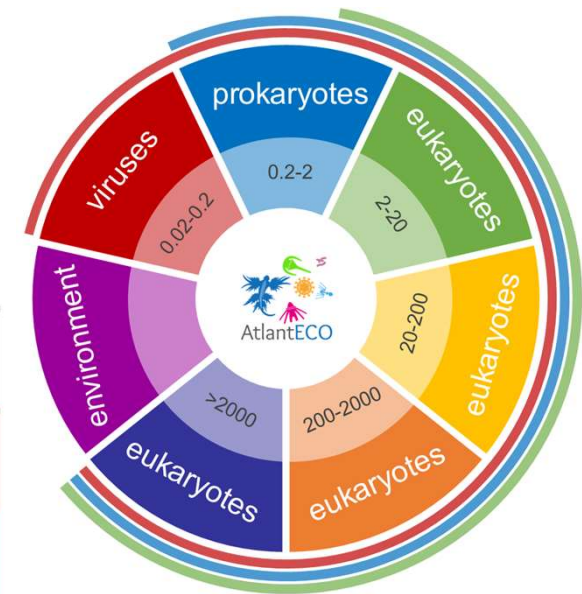
EMBL-EBI

# System of protocols

## Base phenomics protocols

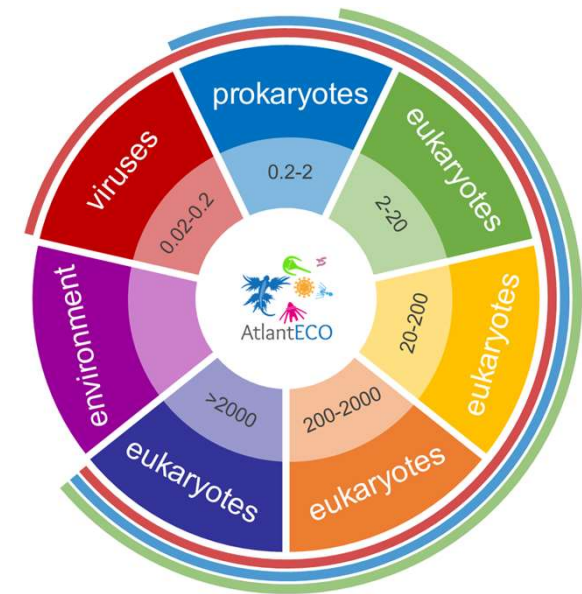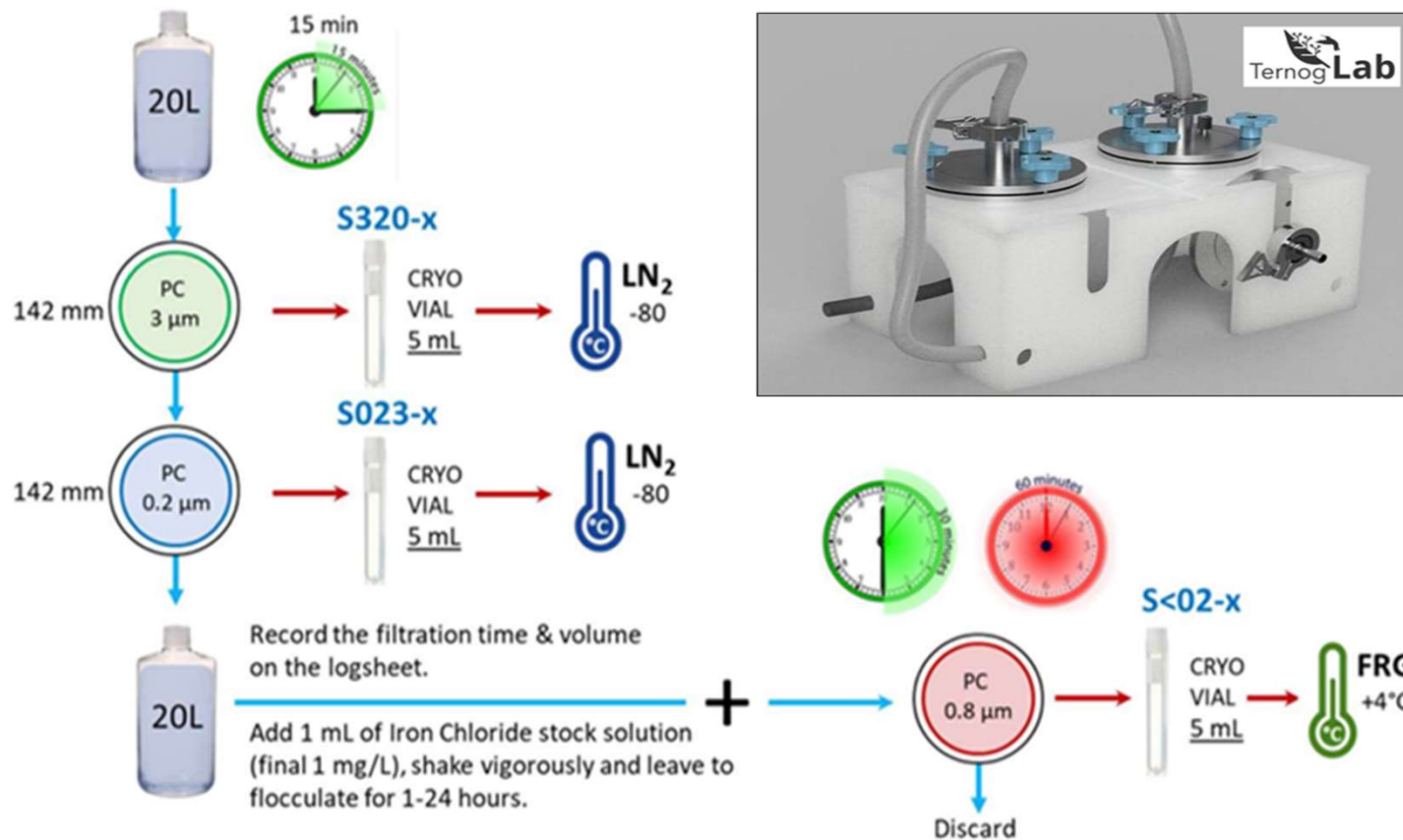| Base Protocol | Target Size fraction | Target Analysis | Target Volume (L) | Target Concentration | Preservation |
|---|---|---|---|---|---|
| I002 | 0.02-0.2 | Flow cytometry | $10^{-3}$ | | $LN_2$ or -80°C |
| I02 | 0.2-2 | Flow cytometry | $10^{-3}$ | | $LN_2$ or -80°C |
| I2 | 2-20 | Fluorescence microscopy | $10^2$ | | +4°C |
| I20 | 20-200 | Flow imaging microscopy | $10^2$-$10^4$ | | live |
| I200 | 200-2000 | Flatbed scan imaging | $10^3$-$10^5$ | | formaldehyde |
| I2000 | >2000 | Flatbed scan imaging | $10^3$-$10^5$ | | formaldehyde |
| environment | multiple | multiple | multiple | | |



EMBL-EBI

# System of protocols

## Base genomics protocols

| Base Protocol | Target Size fraction | Target Analysis | Target Volume (L) | Target Time (min) | Preservation |
|---|---|---|---|---|---|
| S002 | 0.02-0.2 | MetaG, MetaT | 20 | Flocculated 4-24h | +4°C |
| S02 | 0.2-2 | MetaB, MetaG, MetaT | 20 | <15 | LN$_2$ or -80°C |
| S2 | 2-20 | MetaB, MetaG, MetaT | 20 | <15 | LN$_2$ or -80°C |
| S20 | 20-200 | MetaB, MetaG, MetaT | $10^2$-$10^4$ | <15 | LN$_2$ or -80°C |
| S200 | 200-2000 | MetaB, MetaG, MetaT | $10^3$-$10^5$ | <15 | LN$_2$ or -80°C |
| S2000 | >2000 | MetaB, MetaG, MetaT | $10^3$-$10^5$ | <15 | LN$_2$ or -80°C |
| eDNA | >0.2 | MetaB | 2 | | LN$_2$ or -80°C |



Derived Protocols

- Size fractions
- Filtration volume
- Filtration time
- Preservation method
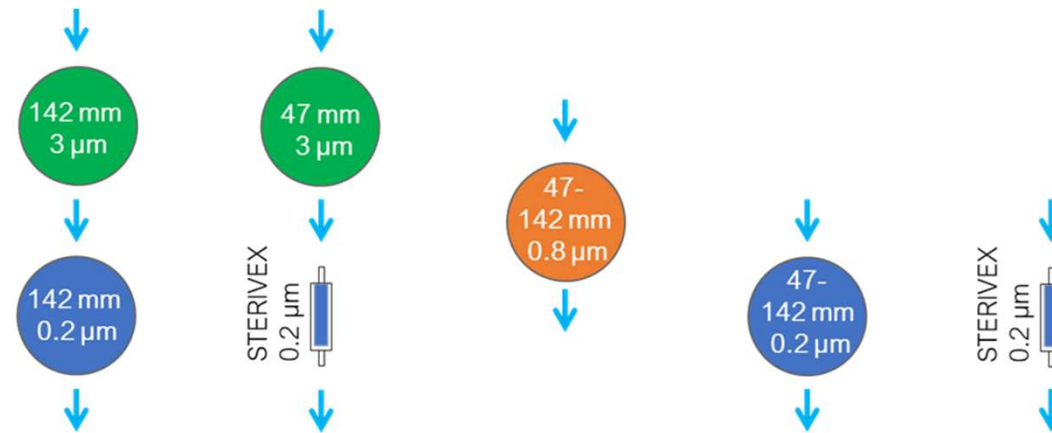
EMBL-EBI

# System of protocols



Derived Protocols

- Size fractions
- Filtration volume
- Filtration time
- Preservation method

# Community Survey
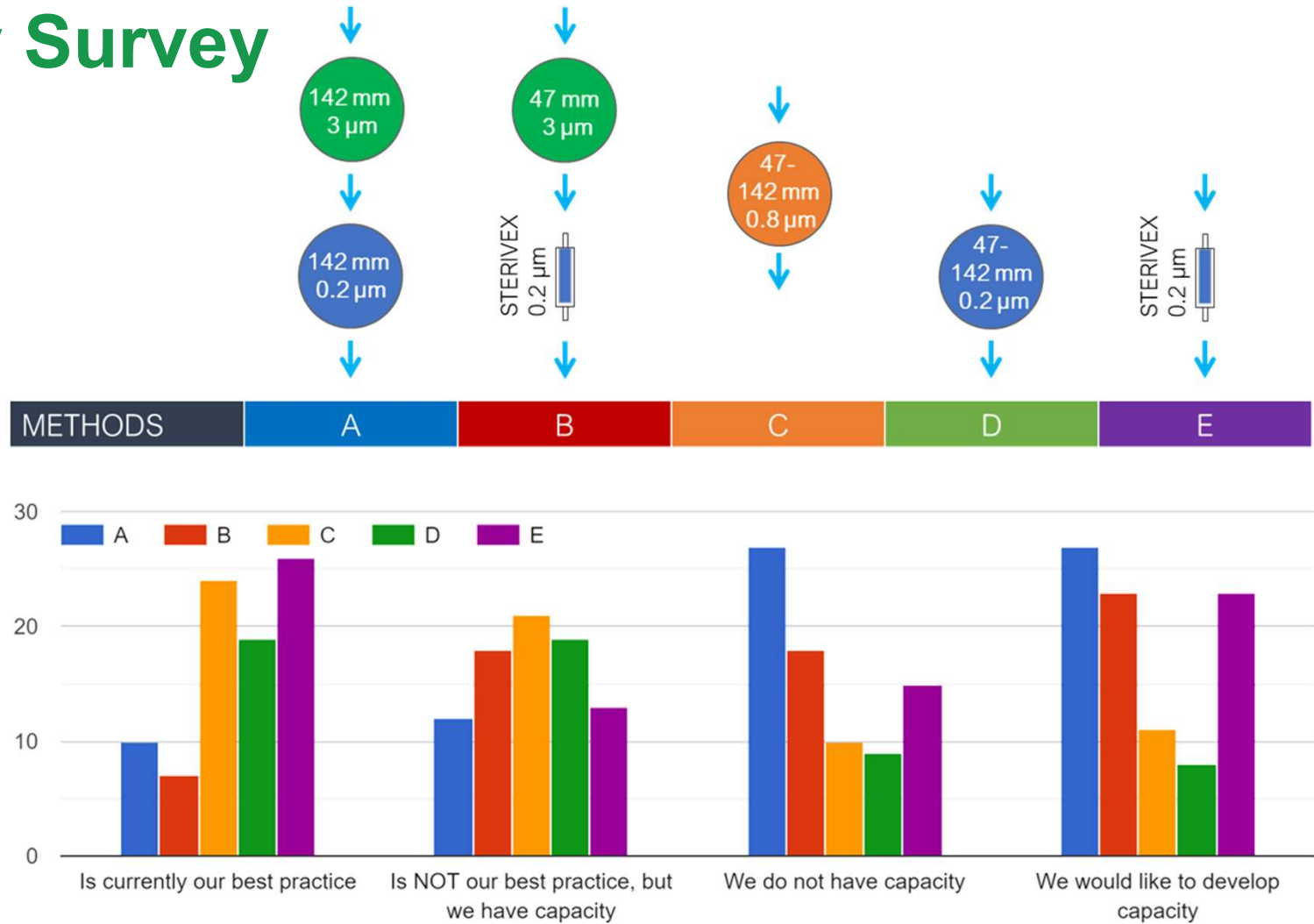
All Atlantic Ocean
Microbiome Sampling



| METHODS | A | B | C | D | E |
|---|---|---|---|---|---|
| Filters | 2x 142 mm membranes | 47 mm membr. + sterivex | 47 or 142 mm membrane | 47 or 142 mm membrane | sterivex |
| Vol. in 15 min. | 10-20 L | 2-5 L | 1-10L | 1-10L | 1-10 L |
| Pump system | large peristaltic | small peristaltic | various | various | various |
| **ADOPTED BY** | **A** | **B** | **C** | **D** | **E** |
| Ocean Sampling Day | | | (✓) | | ✓ |
| Bio-GO-SHIP* | | | | | ✓ |
| Tara Oceans | ✓ | | (✓) | (✓) | |
| EMO-BON** | ✓ | | | | |

*Bio-GO-SHIP Linking marine biodiversity and biogeochemistry (https://biogoship.org/)
**European Marine Omics Biodiversity Observatory Network (https://www.embrc.eu/emo-bon)

EMBL-EBI

# Collecting Metadata – unique identifiers



**Sample provenance metadata**

Write in left column of the logsheet:
1. Sampling depth (m)
2. Replicate # or Control #

Write on the large barcode sticker:
1. Station ID (e.g. 021)
2. Sampling depth (e.g. 200 m)
3. Protocol label (e.g. S320)
4. Replicate # or Control # (e.g R1)

Fix the large barcode sticker on the sample container

Fix the corresponding small barcode sticker on the logsheet, in the appropriate protocol column



EMBL-EBI

# Collecting Metadata – provenance & context

# Collecting Metadata – structured checklists



| attribute | Format / units | comment |
|---|---|---|
| sample id | SAMEA0000000 | |
| sample label | Alpha-numeric, e.g. "project_date-time_station_environment_size-fraction_method" | human readable and meaningful label |
| sampling design, label(s) | alpha-numeric | campaign, station, site, transect, etc. |
| sampling device | alpha-numeric | device name & specifications |
| operator | alpha-numeric | initials or full name |
| sampling date and time | yyyy-mm-dd T hh:mm | in UTC |
| latitude | N/S dd.dddddd or N/S dd mm.mmm | ultimately in decimal degree N |
| longitude | E/W ddd.dddddd or E/S ddd mm.mmm | ultimately in decimal degree E |
| elevation, depth below soil surface | cm | ultimately in metre |
| elevation, depth below sediment surface | cm | ultimately in metre |
| elevation, depth below water surface | m | |
| elevation, altitude above sea level | m | |
| methodological details | alpha-numeric | e.g. processing time & volume |

# Data sharing best practices

Permanent archives selected for the different data types:

- BioSamples for metadata (https://www.ebi.ac.uk/biosamples/)

- ENA for genomics data (http://www.ebi.ac.uk/ena)

- MGnify for metagenomic data (https://www.ebi.ac.uk/metagenomics/)

- PRIDE for proteomics data (https://www.ebi.ac.uk/pride/)

- Metabolights for metabolomics data (https://www.ebi.ac.uk/metabolights/)

- BioImage Archive for imaging data (https://www.ebi.ac.uk/bioimage-archive/)

EMBL-EBI

# Data workflow

**primary data**
curation & archiving

**annotated data**
analysis & archiving

**Environmental**
EMODnet, Copernicus, etc.

**Provenance & Envir. Context**
BioSamples

**Genomics**
ENA, Mgnify & Ensembl

**Proteomics**
PRIDE

**Metabolomics**
Metabolights

**Imaging**
BioImage archive & EcoTaxa

EMBL-EBI

# Data workflow

**primary data**
curation & archiving

**annotated data**
analysis & archiving

**Environmental**
EMODnet, Copernicus, etc.

**Provenance & Envir. Context**
BioSamples

**Genomics**
ENA, Mgnify & Ensembl

**Proteomics**
PRIDE

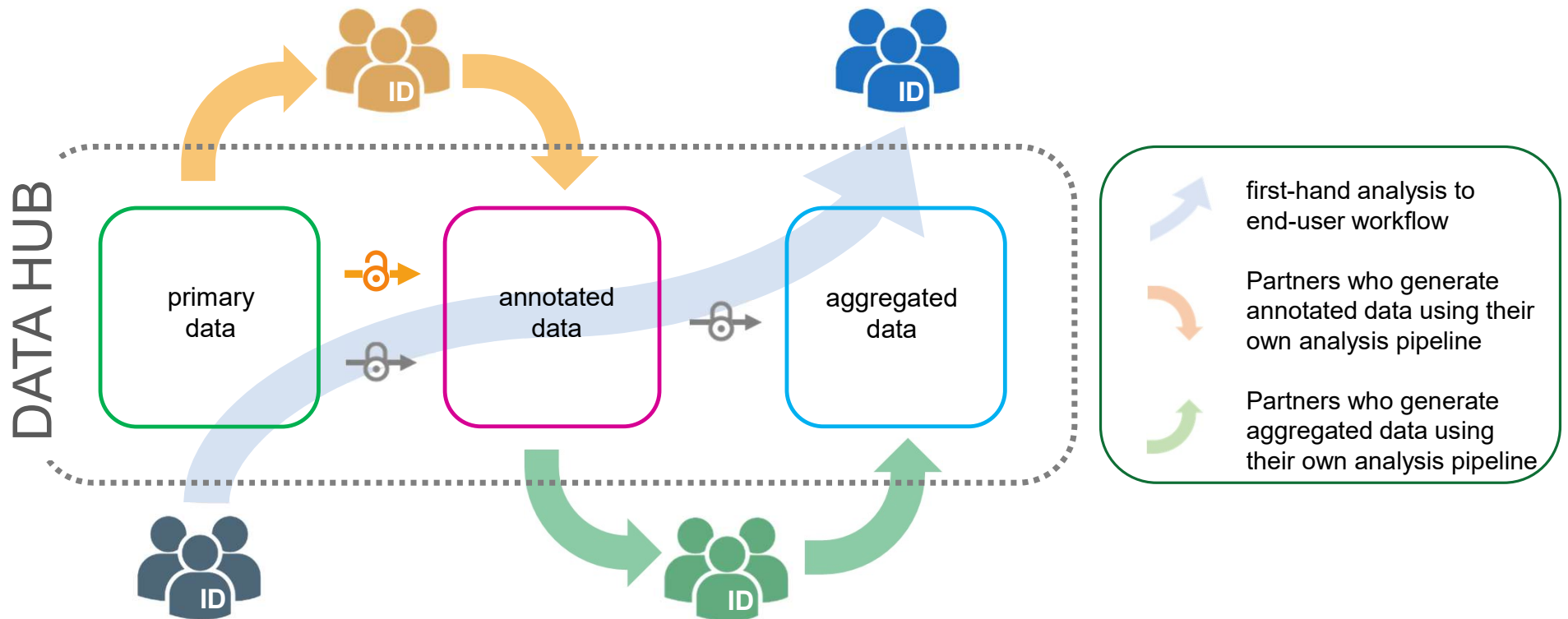**Metabolomics**
Metabolights

**Imaging**
BioImage archive & EcoTaxa

DATA HUB

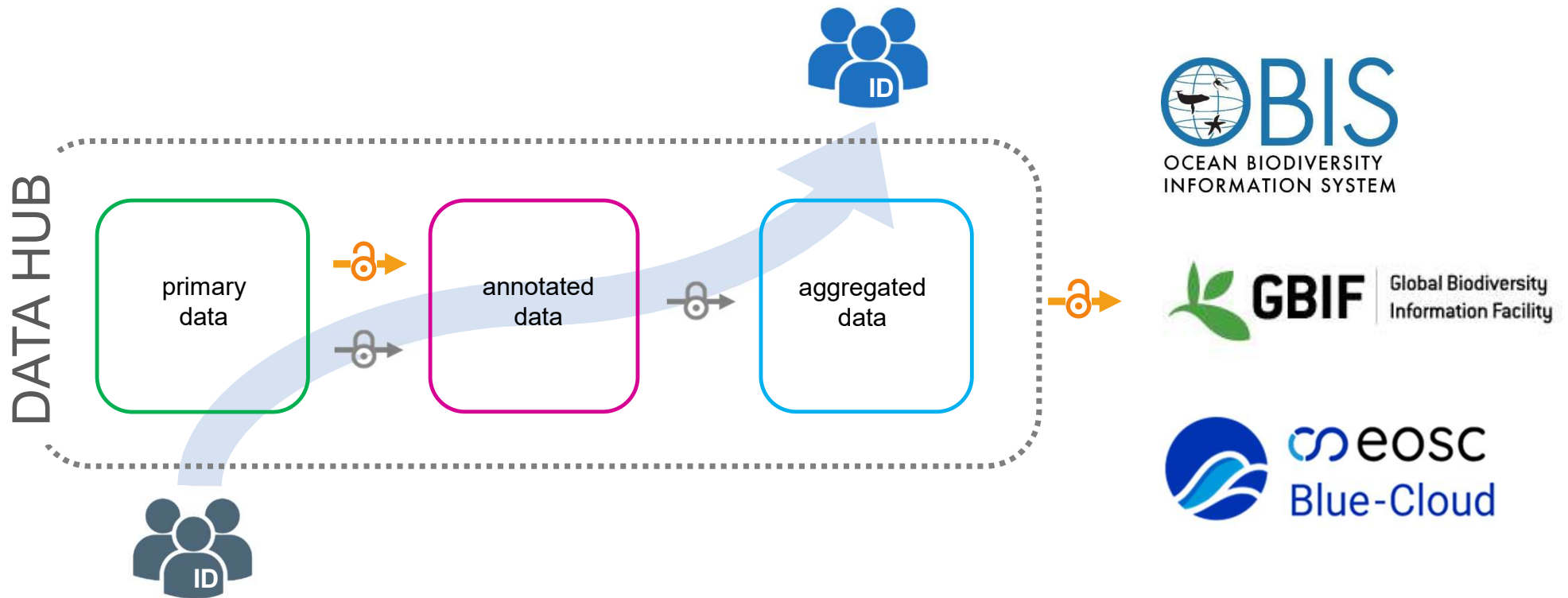**aggregated data**
analysis & archiving

< sample IDs >
< raw data file IDs >
< annotation file IDs >

< entity IDs >

EMBL-EBI

# Data Hub workflow

# Data Hub workflow

Discovery & Access

Virtual Labs & Workbenches
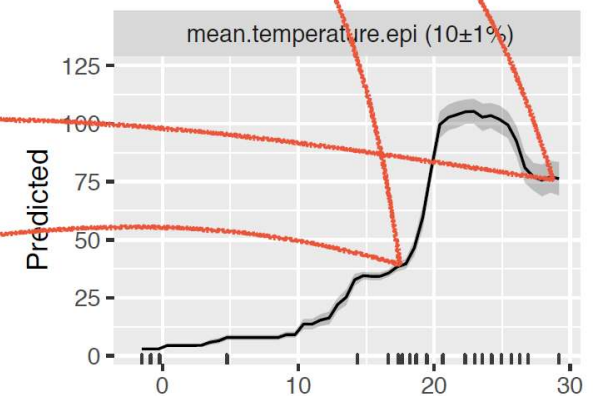
mean.temperature.epi

mean.temperature.epi (10±1%)

EMBL-EBI

eosc | Blue-Cloud

https://blue-cloud.org/

Discovery & Access

Virtual Environment

Publishing Services

EMODnet

zenodo

OpenAIRE

EMBL-EBI

# Publication best practices

Exemplar initiatives

- Tara Oceans
- Tara Pacific
- AtlantECO

Best Practices

- Jointly owned results
- Early notification of intent to use "jointly owned results"
- Co-authorship
- Open Access publication

EMBL-EBI

# Early notification

MMA Results

**Notification of intent**
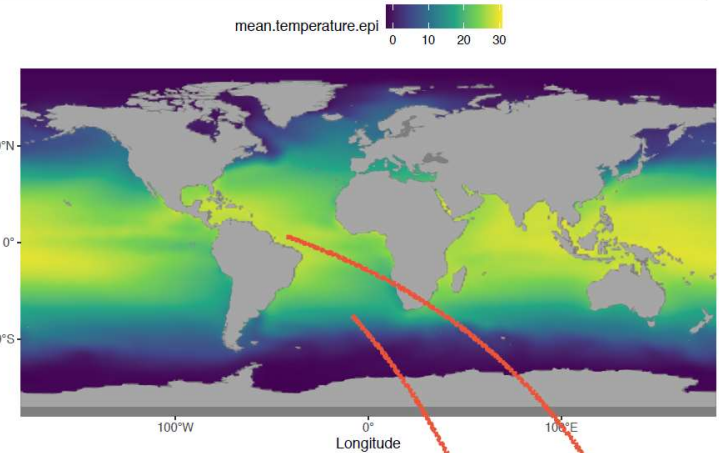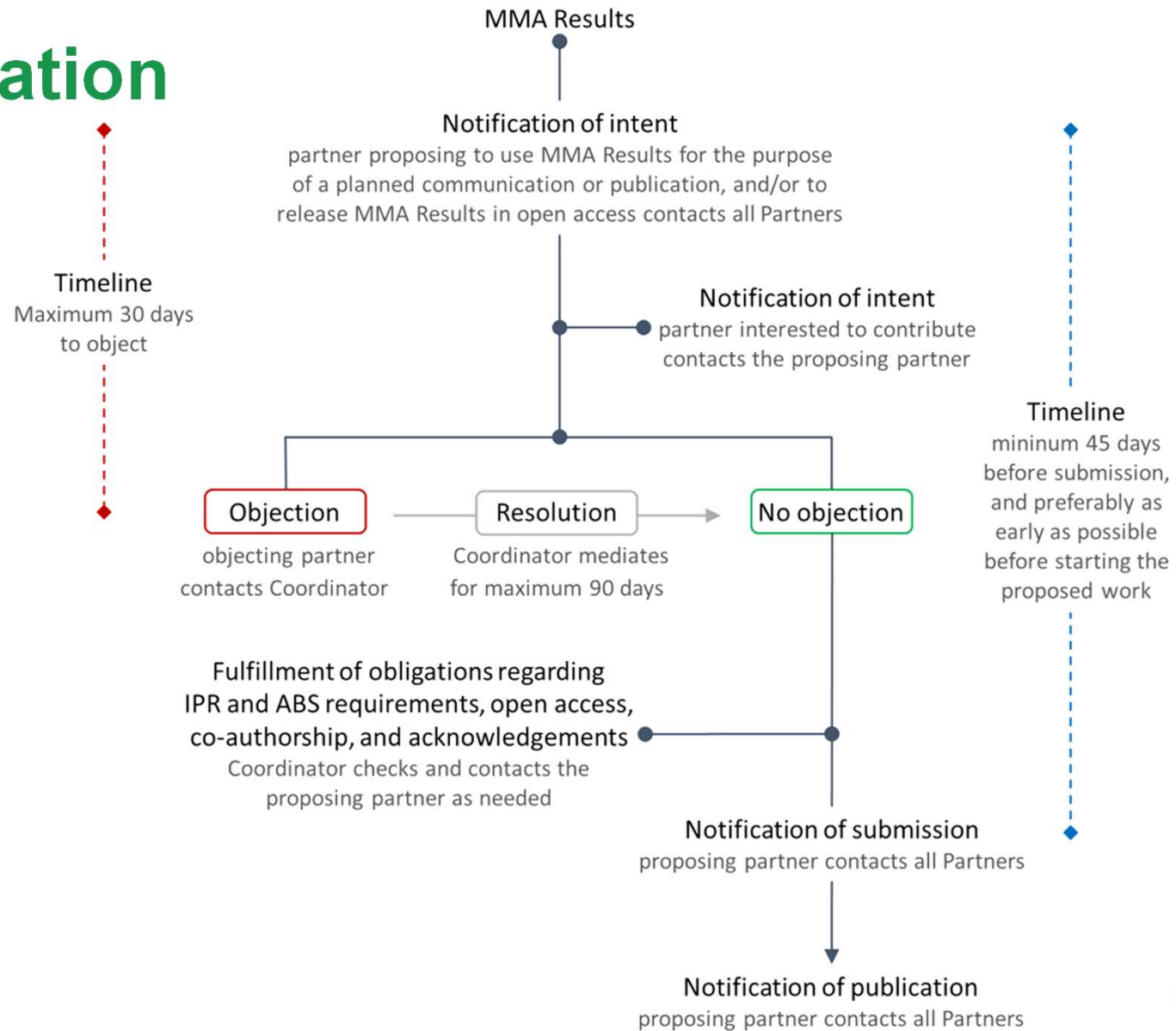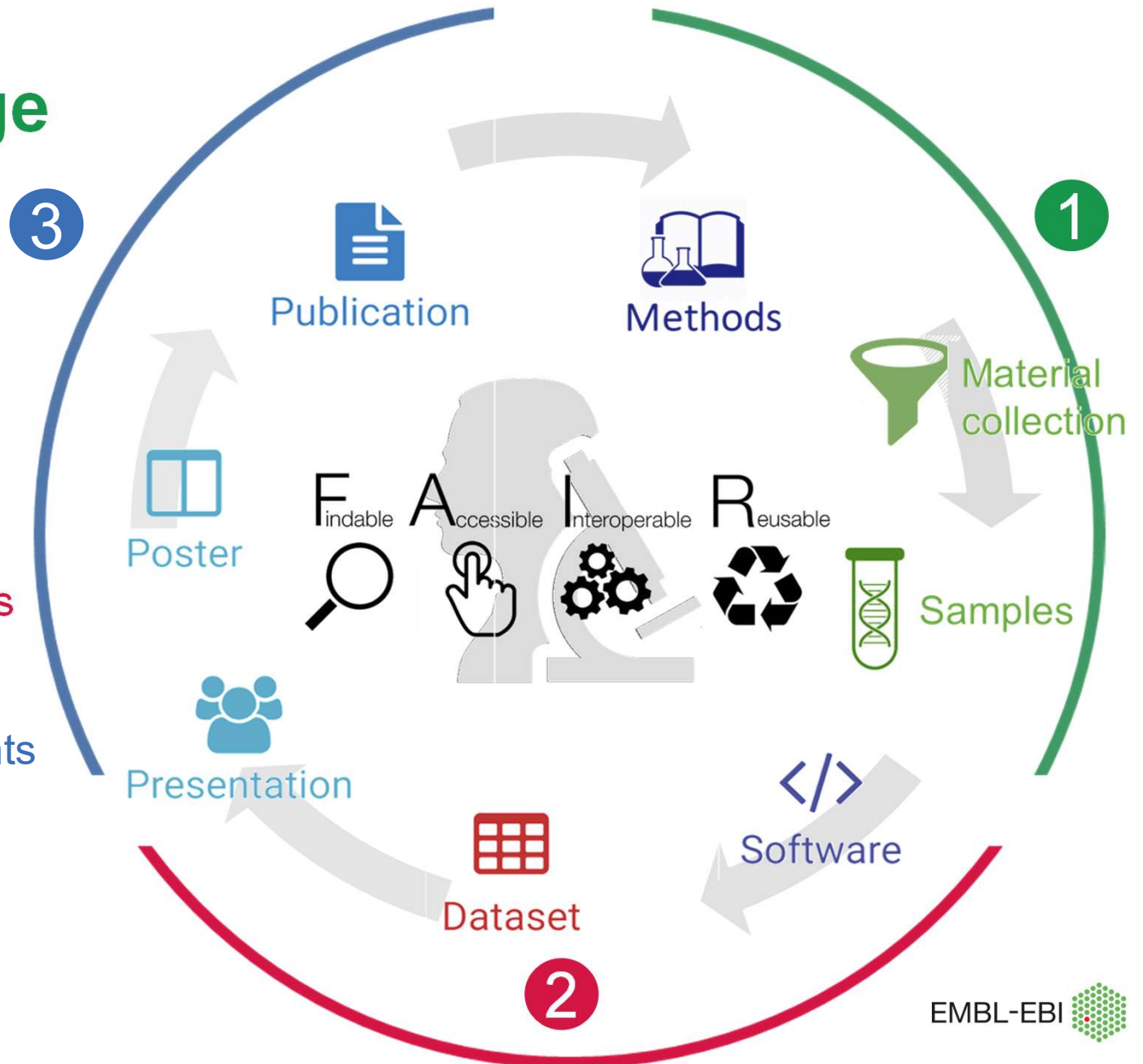partner proposing to use MMA Results for the purpose
of a planned communication or publication, and/or to
release MMA Results in open access contacts all Partners

**Notification of intent**
partner interested to contribute
contacts the proposing partner

**Timeline**
Maximum 30 days
to object

**Timeline**
mininum 45 days
before submission,
and preferably as
early as possible
before starting the
proposed work

| Objection | Resolution | No objection |

objecting partner
contacts Coordinator

Coordinator mediates
for maximum 90 days

**Fulfillment of obligations regarding
IPR and ABS requirements, open access,
co-authorship, and acknowledgements**
Coordinator checks and contacts the
proposing partner as needed

**Notification of submission**
proposing partner contacts all Partners

**Notification of publication**
proposing partner contacts all Partners

EMBL-EBI

# Thank you
# for your kind attention

**Stéphane PESANT**

**Senior marine biocurator**

pesant@ebi.ac.uk

EMBL-EBI