# Augmenting the Environmental Context in Biological Samples using Geographic and Taxonomic Metadata

P. Woollard, S. Pesant, J. Burgin, S. Jayathilaka and G. Cochrane

ORCID: 0000-0002-7654-6902, 0000-0002-4936-5209, 0000-0002-9818-1094, 0000-0002-0154-8807, 0000-0001-7954-7057

## Challenge: determining the provenance of nucleotide data for marine and freshwater environments



- Is a sample marine?
- Which Exclusive Economic Zone?
- Which ocean or sea?
- Is a sample from freshwater?
- What are the predicted light levels at the time and location of sampling?

## European Nucleotides Archive(ENA)

- A repository of the world's nucleotide data
- European node of INSDC
- Creators of tools for submission and retrieval

### Provide Checklists: with limited mandatory metadata

- balancing the needs of the different users
- science and technology all evolve.

### Key Sample Provenance Attributes at ENA

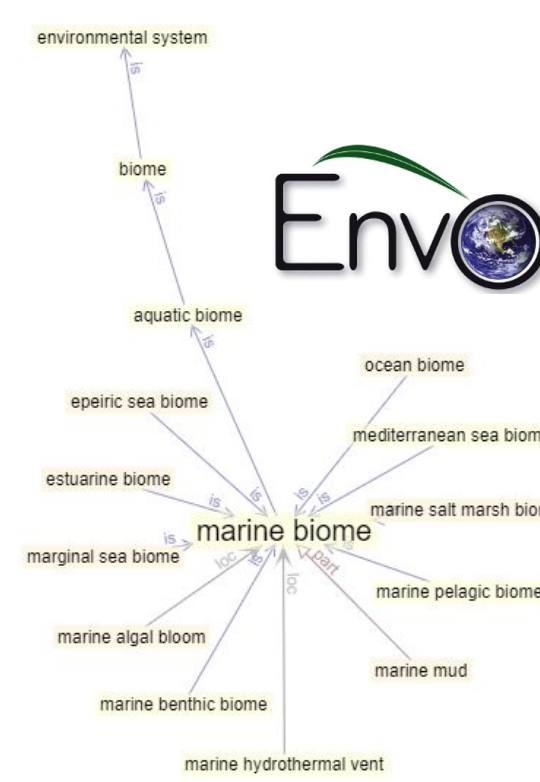| Sample site fields | Control | Total samples with these | Comment |
|---|---|---|---|
| Taxonomy *"taxon_id"* | NCBI taxonomy | 100% (23,724,884) | - mandatory! |
| Geography *"latitude"* and *"longitude"* | [0-9]+.[0-9]+ | 15% (3,588,320) | - much automated cleaning by ENA technical team to provide this cleanly.<br>- 1,918 samples with start and end coordinates |
| Keywords *"Environment_biome"* | Free text + sometimes ENVO terms | 6% (1,523,513) | - Often allows determination if marine/freshwater biome<br>- free text = harmonisation needed |
| country | Controlled "country" list | 62% (14,782,362) | - (9,045 samples are classified as sea or ocean) |

## Implementation

- External resources inputs:
  - Species environment labels from the World Register of Marine Species(WoRMS)
  - Geography as polygons in shapefiles (WWF, OpenStreetMap.de, marinerege etc.)
- A mainly python pandas based workflow was used to cope with 23 million samples worth of metadata



Environment flag by taxonomy experts
- marine
- brackish
- freshwater
- terrestrial

**Example of a shapefile (multiple polygons with coordinates)**
*Intersect_EEZ_IHO_v4_2020.zip:* Citation: Flanders Marine Institute (2020). The intersect of the Exclusive Economic Zones and IHO sea areas, version 4. Available online at https://www.marineregions.org/ https://doi.org/10.14284/402 The IHO polygon layer has a low resolution, but the EEZ polygon layer has a high resolution coastline (GSHHS).

### Mapping Latitude and Longitude Coordinates to Geographic Features



~20 seconds mapping all coords vs a shapefile

Coordinates of samples coloured by EEZ/IHO labels

### Environment by keywords

environment_biome examples (ENVO or free text)
cropland biome
freshwater biome
Human gut
village biome
ENVO:dense settlement biome
ENVO:00009003
coral reef
ENVO:human-associated habitat
forest
ENVO:2100002
Human
area of cropland
grassland
mouse gut
soil
marine biome
terrestrial biome [ENVO:00000446]
lentic water body

Unclassified "metagenome" taxonomic terms

## Merging Evidence

**Workflow: Integrating Evidence of "Blue" Domain**



**Combined coordinates, taxonomy and env_biome**
- Just counts



With the taxonomy info. in particular the marine domain massively increased in terms of sample coverage

### Serving relevant marine & freshwater data



Informing

### Essential annotations for science & society

www.ebi.ac.uk

Blue-Cloud2026
A federated European FAIR and Open Research Ecosystem for oceans, seas, coastal and inland waters

EMBL's European Bioinformatics Institute (EMBL-EBI)
Wellcome Genome Campus, Hinxton,
Cambridgeshire. CB10 1SD. UK.
T +44(0)1223 494 444