

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

A blue partition for marine and freshwater
environments



Peter Woollard

Data Standards biocurator, ENA, EMBL-EBI

Acknowledgements: **Stéphane Pesant**, Josie Burgin, Suran Jayathilaka and Guy Cochrane



Blue-Cloud2026

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

EMBL-EBI



Challenge: determining the provenance of nucleotide data from 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?

Blue Partition



The European Nucleotide Archive(ENA)

What is the ENA?

- A repository of the world's nucleotide data
- European node of INSDC
- Creators of tools for submission and retrieval
- Checklists are used to capture metadata, using minimum requirement and more detailed optional fields



Key Sample Provenance Attributes at ENA

Sample site fields	Control	Total samples with these	Comment
Taxonomy “ <i>taxon_id</i> ”	NCBI taxonomy	100% (23,724,884)	
Geography “ <i>latitude</i> ” and “ <i>longitude</i> ”	[0-9]+.[0-9]+	15% (3,588,320)	1,918 samples with start and end coordinates
Keywords “ <i>Environment_biome</i> ”	Free text + ENVO terms advised	6% (1,523,513)	Often allows determination if marine/freshwater biome
country	Controlled “country” list	62% (14,782,362)	(9,045 samples are classified as sea or ocean)

*ENA snapshot on 2023-03-10; 23,724,884 samples in total

Environment by Geography

Calculate if latitude and longitude are within geographic shapes

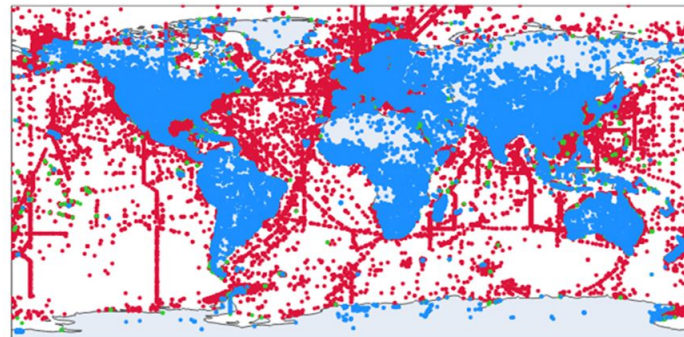
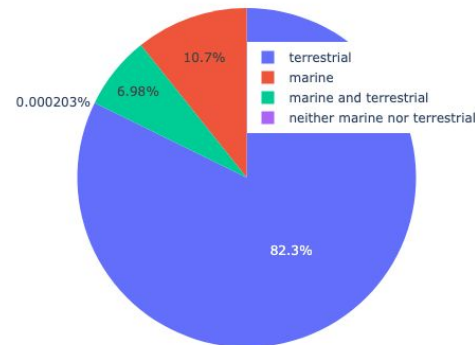
Examples of shapefile (multiple polygons with coordinates)

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).

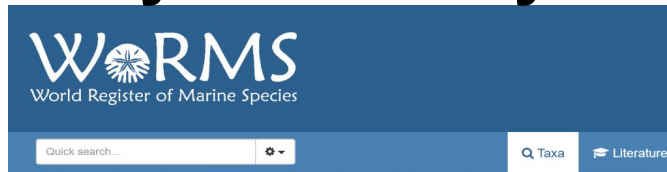
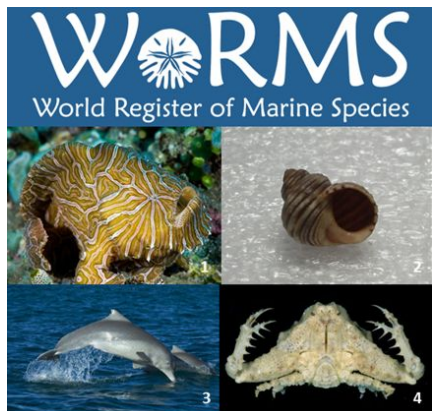
Longhurst_v4_2010.zip: Citation: Flanders Marine Institute (2009). Longhurst Provinces. Online at <https://www.marineregions.org/>. Unknown resolution, but likely coarse

GIS_hs_snapped.zip: Citation: Freshwater Ecoregions of the World (FEOW) hydrosheds (2008) <https://www.feow.org/download>

For samples with geographic provenances:



Environment by Taxonomy



WoRMS taxon details

★ *Calanus finmarchicus* (Gunnerus, 1770)

AphiaID	104464 (urn:lsid:marinespecies.org:taxname:104464)														
Classification	Biota > ★ Animalia (Kingdom) > ★ Arthropoda (Phylum) > ★ Crustacea (Subphylum) > ★ Multicrustacea (Superclass) > ★ Copepoda (Class) > ★ Neocopepoda (Infraclass) > ★ Cyclopoida (Superorder) > ★ Calanoida (Order) > ★ Calanidae (Family) > ★ Calanus (Genus) > ★ <i>Calanus finmarchicus</i> (Species)														
Status	accepted														
Rank	Species														
Typetaxon of	★ <i>Calanus</i> Leach, 1816														
Parent	★ <i>Calanus</i> Leach, 1816														
Orig. name	★ <i>Monoculus finmarchicus</i> Gunnerus, 1770														
Direct children (4)	Subspecies ★ <i>Calanus finmarchicus helgolandicus</i> Tanaka, 1907 Subspecies ★ <i>Calanus finmarchicus pallasii</i> Karavaev, 1894 Subspecies ★ <i>Calanus finmarchicus szekensis</i> Stalberg, 1961 Subspecies ★ <i>Calanus finmarchicus finmarchicus</i> (Gunnerus, 1770)														
Environment	marine, brackish, fresh, terrestrial														
Original description	(of ★ <i>Monoculus finmarchicus</i> Gunnerus, 1770) Gunnerus, 1770. <i>Kiøbenhavnske Selskab.</i> 10: 166-176., available online at https://www.biodidac.dk/monoculus-fimmaricus														
Descriptive notes	★ Distribution <i>C. helgolandicus</i> , a congener species of <i>C. finmarchicus</i> . ★ Distribution Arctic to West Indies														
Taxonomic citation	Walter, T.C.; Boxshall, G. (2023). World of Copepods Database https://marinespecies.org/aphia.php?p=taxdetails&id=104464														
Taxonomic edit history	<table> <thead> <tr> <th>Date</th> <th>action</th> </tr> </thead> <tbody> <tr> <td>2004-12-21 15:54:05Z</td> <td>created</td> </tr> <tr> <td>2008-08-06 13:32:15Z</td> <td>changed</td> </tr> <tr> <td>2009-11-29 21:58:33Z</td> <td>changed</td> </tr> <tr> <td>2010-06-28 13:28:09Z</td> <td>changed</td> </tr> <tr> <td>2020-08-10 10:49:04Z</td> <td>changed</td> </tr> <tr> <td>2020-10-06 15:25:25Z</td> <td>changed</td> </tr> </tbody> </table>	Date	action	2004-12-21 15:54:05Z	created	2008-08-06 13:32:15Z	changed	2009-11-29 21:58:33Z	changed	2010-06-28 13:28:09Z	changed	2020-08-10 10:49:04Z	changed	2020-10-06 15:25:25Z	changed
Date	action														
2004-12-21 15:54:05Z	created														
2008-08-06 13:32:15Z	changed														
2009-11-29 21:58:33Z	changed														
2010-06-28 13:28:09Z	changed														
2020-08-10 10:49:04Z	changed														
2020-10-06 15:25:25Z	changed														
Licensing	The webpage text is licensed under a Creative Commons Attribution 4.0 License														

Environment flag by taxonomy experts

- marine
- brackish
- freshwater
- terrestrial



★ *Calanus finmarchicus*...

~189K taxa at NCBI

~78K (41%) taxa matched NCBI:WoRMS

~11K (14%) with data at ENA

Samples

~625K (3%) from matched taxa NCBI:WoRMS

~91K (15%) with geographic provenance

Walter, T. Chad

Walter, T. Chad

<https://creativecommons.org/licenses/by/4.0/>

[taxonomic tree]

Sources (26)

Documented distribution (21)

Notes (3)

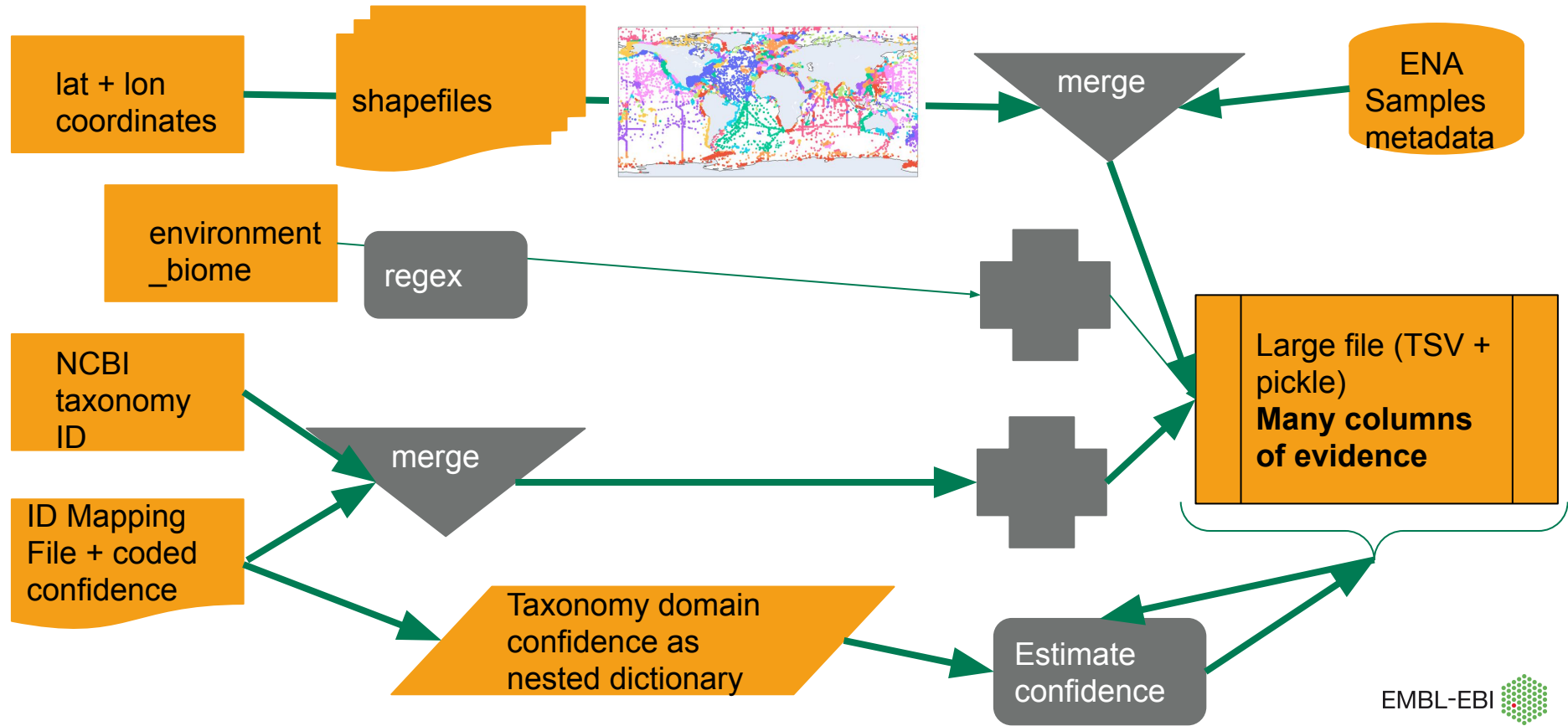
Attributes (149)

Links (10)

Images (15)



Workflow: Integrating Evidence of “Blue” Domain



Essential annotations for science & society

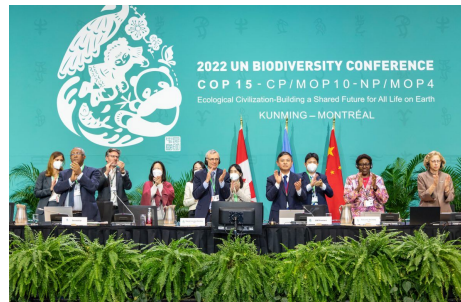


Clearinghouse

curl "https://www.ebi.ac.uk/ena/clearinghouse/api/curations/16f53cc2-532c-486e-b8b1-b07ddb05f93b"

2>/dev/null | jq

```
{
  "curations": [
    {
      "id": "16f53cc2-532c-486e-b8b1-b07ddb05f93b",
      "recordType": "sample",
      "recordId": "SAM00003553",
      "attributePost": "EEZ-name",
      "valuePost": "Japanese Exclusive Economic Zone (mrgid:8487)",
      "attributeDelete": false,
      "assertionMethod": "automatic assertion",
      "assertionEvidences": [
        { "identifier": "ECO:0000203", "shortForm": "ECO_0000203", "label": "automatic assertion" },
        { "identifier": "ECO:0000366", "shortForm": "ECO_0000366", "label": "evidence based on logical inference from automatic annotation used in automatic assertion" }
      ],
      "assertionSource": "",
      "assertionAdditionalInfo": "confidence:high; evidence:sample coordinates within EEZ shapetile",
      "providerName": "European Nucleotide Archive",
      "providerUrl": "https://www.ebi.ac.uk/ena/browser/home",
      "submittedTimestamp": "2023-04-05T07:09:55.109+0000",
      "updatedTimestamp": "2023-04-05T07:09:55.109+0000",
      "suppressed": false
    }
  ],
  "totalAttributes": 1
}
```



Biosample: SAM00003553

HSJLFDJ01_Monbetsu3V7V9_V7V9

Organism: marine metagenome
Accession: SAM00003553
Sample Title: HSJLFDJ01_Monbetsu3V7V9_V7V9
Location: 44.337 N 143.38083 E
Sample Alias: SAM00003553
Secondary Accession: DRS012359
Sample Name: DRS012359
Target Gene: 16S rRNA
Geo Loc Name: Japan:Monbetsu-shi,Kayakoiom
Seq Meth: pyrosequencing

Show More

3rd Party Curations

Operation	Attribute Name	Attribute Value	Assertion Method	Assertion Evidences	Provider Name	Status
Add	EEZ-name	Japanese Exclusive Economic Zone (mrgid:8487)	automatic assertion	automatic assertion; evidence based on logical inference from automatic annotation used in automatic assertion	European Nucleotide Archive	
Add	EEZ-sovereign-level-1	Japan (mrgid:2121; ISO3166-1 alpha-3:JPN; ISO3166-1 num:392)	automatic assertion	automatic assertion; evidence based on logical inference from automatic annotation used in automatic annotation	European Nucleotide Archive	



Thank you



The Blue-Cloud 2026 project has received funding from the European Union's Horizon Europe research and innovation programme under grant agreement No 101094227. The H2020 project Blue-Cloud received funding from the European Union's Horizon programme call BG-07-2019-2020, topic: [A] 2019 - Blue Cloud services, Grant Agreement No. 862409.



BACKUPS

(REMOVE THIS SLIDE BEFORE

Abstract - 10mins lightning talk 7mins taking 3mins Q+A

POSTING)

Title: Providing Expanded Contextual Metadata for Biological Samples using Both Geographic and Taxonomic Factors

Authors: All for us.

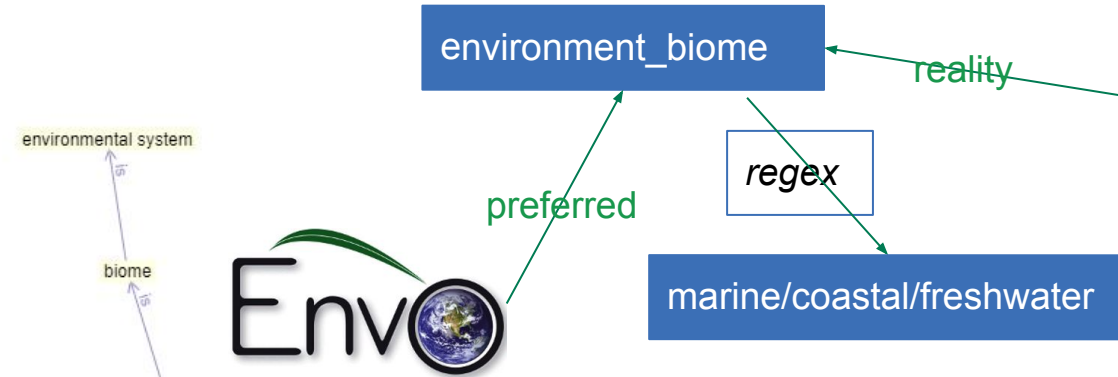
The European Nucleotide Archive (ENA; <https://www.ebi.ac.uk/ena>) is a long-standing database of record for nucleotide sequence data and associated metadata. The ENA has minimal required metadata standards for submitted records to balance the needs of the data generators/submitters and making the metadata as FAIR as possible for downstream users though recommended standards are not always utilised to their full potential and details can be left out.

There are nearly 200,000 marine samples alone within ENA and as part of the BlueCloud project (<https://blue-cloud.org/>) it was identified that there was a need to enhance the available specific metadata for marine and freshwater samples. By utilising user-provided geographic metadata, we can assert additional contextual metadata to enhance the existing sample records. Approximately 17% of all ENA samples have GPS coordinates. We have used the GPS coordinates to determine additional metadata, for example, the geographic political regions (e.g. countries and EEZs) and environment types (e.g. land and sea), via computational geometry. These were compared to existing submitter metadata provided with these samples. Additionally organism taxonomies were categorised with their likely marine or freshwater environment. The submitter, GPS and taxonomy insights were merged and compared. As expected much of the time there is clear cut metadata agreement, sometimes explainable differences and occasionally harder to explain or understand differences.

For the ENA and similar archives, submitter entered data is the record and so metadata cannot be changed substantively on the primary record without the approval of data owners. The extra contextual metadata is being added to the ELIXIR Contextual Data Clearinghouse see

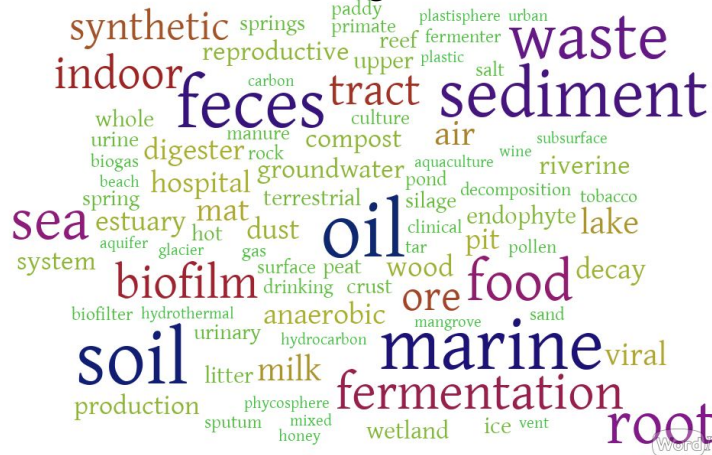
<https://elixir-europe.org/internal-projects/commissioned-services/establishment-data-clearinghouse>; the metadata will be programmatically available from <https://www.ebi.ac.uk/ena/clearinghouse/api/>. It will thus be straightforward to

Environment by Keywords

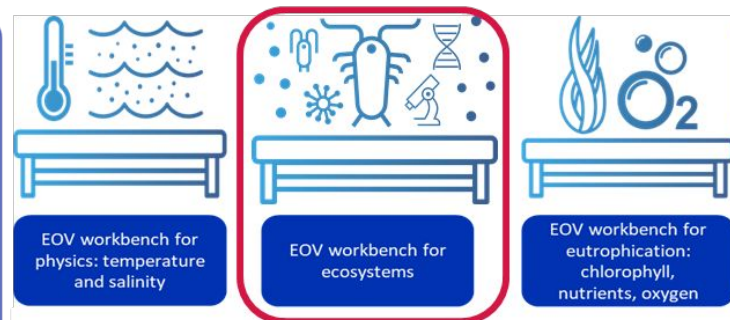
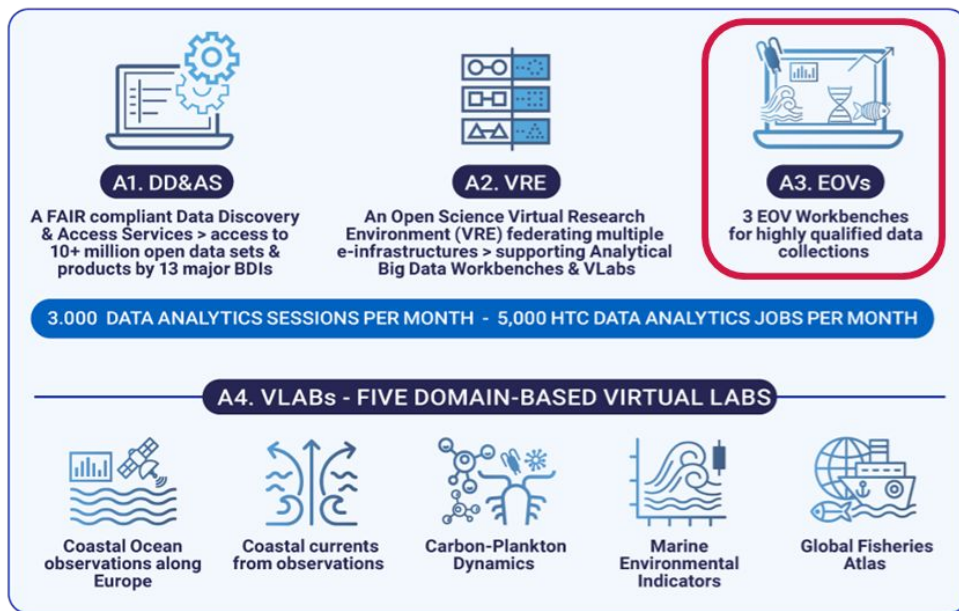
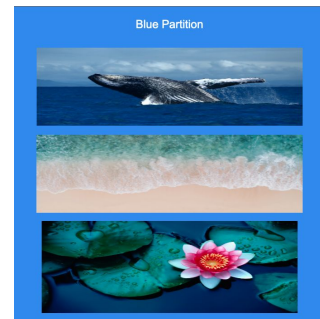


marine benthic biome
cropland biome
freshwater biome
coral reef
ENVO:human-associated habitat
forest
Bos taurus
ENVO:2100002
area of cropland
grassland
mouse gut
soil
marine biome
terrestrial biome [ENVO:00000446]
lentic water body

Unclassified “metagenome” taxonomic terms



Serving relevant marine & freshwater data



Essential Ocean Variables (EOV) Workbenches

“Developing and testing analytical Blue Cloud WorkBenches to generate high quality data products”

Mapping Latitude and Longitude Coordinates to Geographic Features

