

Going beyond the spreadsheet - developing Best Practices in ‘long-tail’ environmental data curation and publishing

Corinna Gries¹, Renée Brown², Mary Gastil-Buhl³, Sarah Elmendorf⁴, Hap Garritt⁵, Mary Martin⁶, Greg Maurer⁷, An Nguyen⁸, John Porter⁹, and Timothy Whiteaker¹⁰

¹University of Wisconsin Madison

²University of New Mexico Main Campus

³Moorea Coral Reef LTER

⁴University of Colorado Boulder

⁵Woods Hole Marine Biological Lab

⁶University of New Hampshire Main Campus

⁷New Mexico State University

⁸University of Texas at Austin

⁹University of Virginia

¹⁰CRWR

November 26, 2022

Abstract

The research data repository of the Environmental Data Initiative (EDI) is a signatory of the FAIR Data Principles. Building on over 30 years of data curation research and experience in the NSF-funded US Long-Term Ecological Research program (LTER), it provides mature functionalities, well established workflows, and support for ‘long-tail’ environmental data publication. High quality scientific metadata are enforced through automatic checks against community developed rules and the Ecological Metadata Language (EML) standard. Although the EDI repository is far along the continuum of making its data FAIR, representatives from EDI and the LTER Information Management community have recently been developing best practices for the edge cases in environmental data publishing. Here we discuss and seek feedback on how to best handle the publication of these ‘long-tail’ data when extensive additional data are available along with e.g., genomics data, physical specimens, or flux tower data. While these latter data are better handled in other discipline-specific repositories such as NCBI, iDigBio, and AmeriFlux, they are frequently associated with other data collected at the same time and location, or even from the same samples. This is particularly relevant across the LTER Network, where sites represent integrative research projects. Questions we address (and seek community input from) include: How to archive documents and images when they are data, e.g., field notebooks, or time-lapse photographs of plant phenology? How to deal with data from Unmanned Vehicle (e.g., drones and underwater gliders), acoustic data, or model outputs, which may be several terabytes in size? How should processing scripts or modeling code be associated with data? Overall, these best practices address issues of Findability and Accessibility of data as well as greater transparency of the research process.

Going Beyond the Spreadsheet

Developing Best Practices in 'long-tail' environmental data
curation and publishing



C. Gries, R.F. Brown, G.Gastil-Buhl, S. Elmendorf, H. Garritt,
M. Martin, G. Maurer, A. Nguyen, J.H. Porter, T. Whiteaker



Introduction

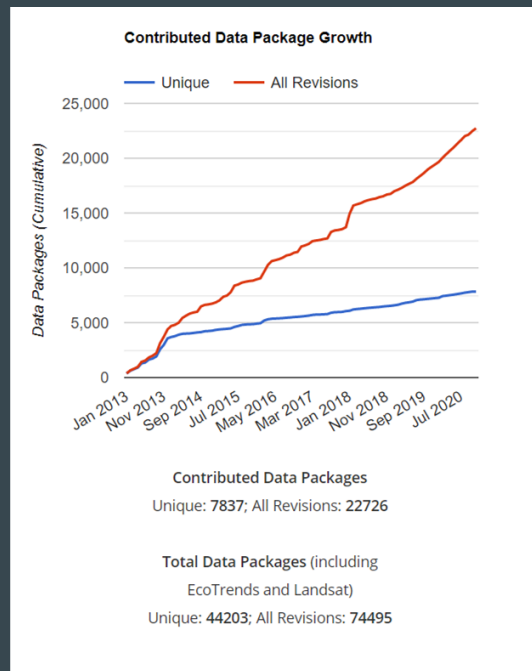
Environmental Data Initiative repository

FAIR data repository

Ecological Metadata Language

Data package: data table(s), metadata

Metadata generation tools

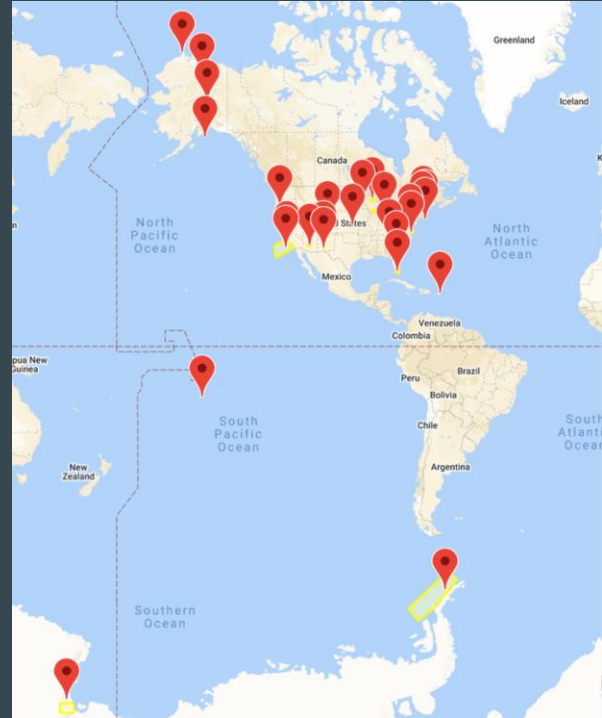


Introduction

Long-Term Ecological Research network

40 years of data management

Professional data managers at each site



Archiving other data types

Modeled Data

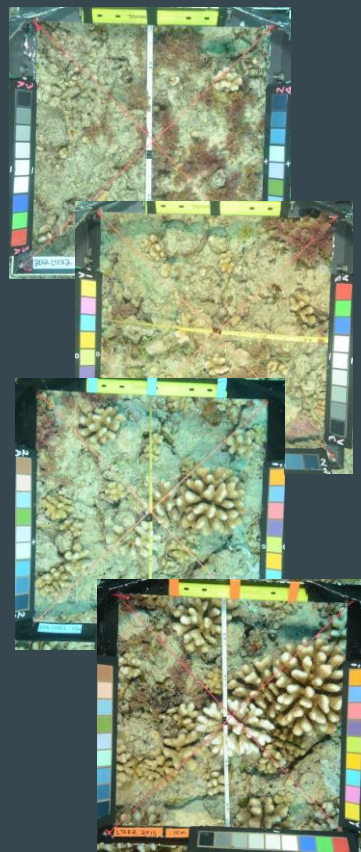
Code and Scripts

Data in More than one Repository

Images and Documents as Data

Spatial Data

Data from Small, Moving Platforms



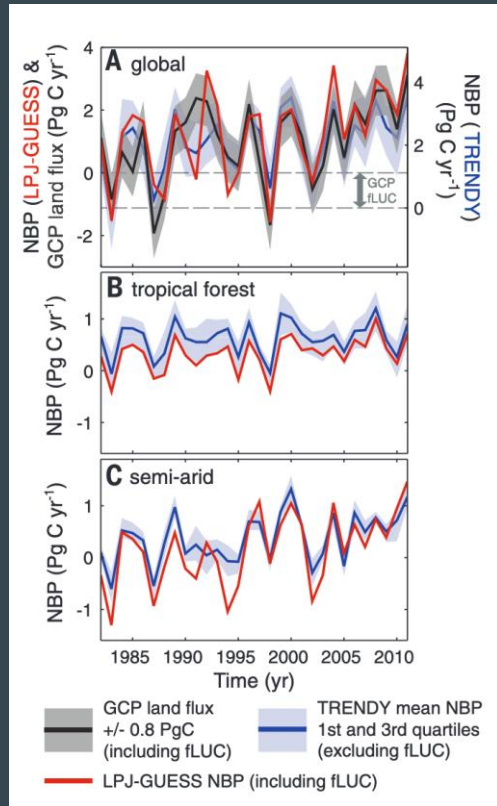
Modeled Data

Transparency / Reproducibility

Parameter settings, input data, output data

Model code

Size



From Ahstrom et al 2015



Code and Scripts

```
6  
7 # Fix any interval or ratio columns mistakenly read in as nominal and nominal columns read as numer  
8  
9 if (class(dt1$SEQ)=="factor") dt1$SEQ <-as.numeric(levels(dt1$SEQ))[as.integer(dt1$SEQ) ]  
0 if (class(dt1$YEAR)=="factor") dt1$YEAR <-as.numeric(levels(dt1$YEAR))[as.integer(dt1$YEAR) ]  
1 if (class(dt1$MONTH)=="factor") dt1$MONTH <-as.numeric(levels(dt1$MONTH))[as.integer(dt1$MONTH) ]  
2 if (class(dt1$LONG)=="factor") dt1$LONG <-as.numeric(levels(dt1$LONG))[as.integer(dt1$LONG) ]  
3 if (class(dt1$LAT55)=="factor") dt1$LAT55 <-as.numeric(levels(dt1$LAT55))[as.integer(dt1$LAT55) ]  
4 if (class(dt1$LAT525)=="factor") dt1$LAT525 <-as.numeric(levels(dt1$LAT525))[as.integer(dt1$LAT525) ]  
5 if (class(dt1$LAT50)=="factor") dt1$LAT50 <-as.numeric(levels(dt1$LAT50))[as.integer(dt1$LAT50) ]  
6 if (class(dt1$LAT50)=="factor") dt1$LAT50 <-as.numeric(levels(dt1$LAT50))[as.integer(dt1$LAT50) ]
```

Model code, data manipulation scripts (not binary software packages)

In EDI / code repository

Inside data package / referenced

CodeMeta file



Data in More than one Repository

Related datasets managed in different repositories

E.g., GenBank, FluxNet, IDigBio

Data inventory table (accession #)

```
ACGGTAGCTAATACCGCATAACGTCGCAAG/  
TTACTAGCGATTCCAAC TTCATAAGGTCGAG  
CGTATTCACCGCGGCATGCTGATCCGCGATT/  
GGGGAAAGATTTATCGCCAAAAGATTGGCC
```



Images and Documents as Data

Images, digital field sheets, lab notebooks, protocols

Reanalysis, transparency, documentation

Design data packages for usability

Data inventory table

LTER BIOLOGICAL FIELD DATA

Sequence # 6301 A

Lake SP Station 1 Date 7 SEP 2020 Page 1 of 2

Perm. Anchor? ☒ N
Observers P. SP, P. M
On Station Time 9:25
Off Station Time 11:05

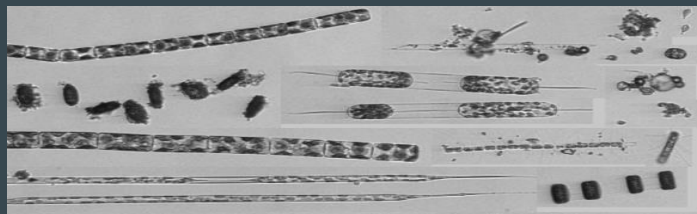
Air Temp (C) 8
Wind Sp. (mph) 4
Wind Direction NW
Wave Size (cm) 5
% Cloud Cover 70

Bottom at (m) 2.1
with (gear) Siphon
Secchi Disk Depth (m) 7.0

Equip: 2m Sch-P
83µ mesh
Start: 9:40
End: 10:30

PERISTALTIC PUMP

DEPTH (m)	SAMPLE TAKEN	CHLOROPHYLL		VOLUME FILTERED	PHYTOPLANKTON	
		FILTER LETTER	MAX PSI		SPAN (m)	STRATUM, VOLUME RUMPED
Surf.		D	15	2700		
1	X					
2						
3	X	R	15	2700		
4						
5	X	L	15	2800		
6						
7	X					
8		P	15	1670		



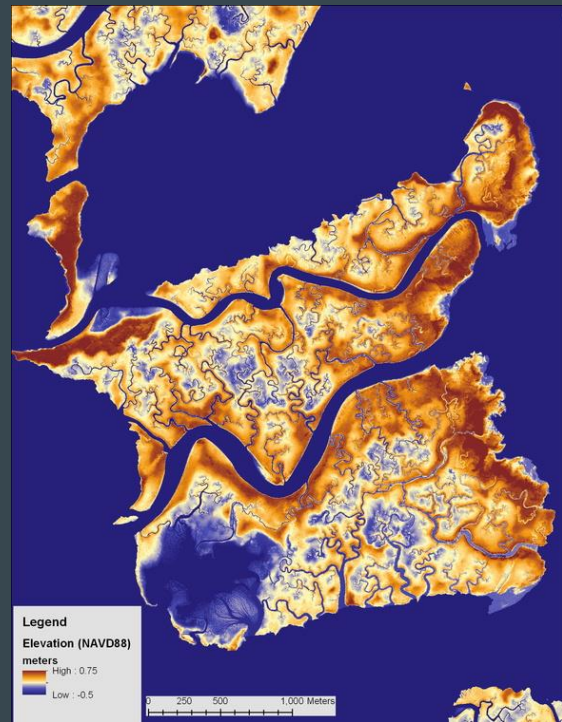
Spatial Data

Widely readable format (shapefile, geoTIFF)

Include metadata files in other formats

EML specific recommendations

Metadata	Geography	Table
Road GIS Layer for the Eastern Shore of Virginia 2013		
Type	Shapefile	
Tags	transportation, Roads, Street Centerline	
Summary		
<p>The TIGER/Line shapefiles and related database files (.dbf) are an extract of selected geographic and cartographic information from the U.S. Census Bureau's Master Address File / Topologically Integrated Geographic Encoding and Referencing (MAF/TIGER) Database (MTDB). The MTDB represents a seamless national file with no overlaps or gaps between parts, however, each TIGER/Line shapefile is designed to stand alone as an independent data set, or they can be combined to cover the entire nation. The All Roads Shapefile includes all features within the MTDB Super Class "Road/Path Features" distinguished where the MAF/TIGER Feature Classification Code (MTFCC) for the feature in MTDB that begins with "S". This includes all primary, secondary, local neighborhood, and rural roads, city streets, vehicular trails (4wd), ramps,</p>		



Data from Small, Moving Platforms

Sensors mounted on uncrewed aerial or underwater systems

Data Inventory Table

Publish raw data, processing code, derived data



What's new

Data inventory table

with additional metadata and accession numbers/GUIDs

Publish raw data and processing code

File with metadata in other formats

as part of the data package

