



Structured Metadata at EDI

2018 March
Environmental Data Initiative (EDI)



Why Use a Metadata Standard?

A Standard provides a structure to describe data with:

- Common terms to allow consistency between records
- Common structure to quickly locate components

In search and retrieval, standards provide:

- Reliable, predictable format for computer interpretation
- A uniform summary description of the dataset



CC image by
ccarlstead on
Flickr



Metadata standards are optimized

- Content

Most metadata standards include:

Who – information on who to contact

What – a description of the available item

Some include:

Where – geographical locations

When – dates and times

- Examples

Darwin Core - organism occurrences (GBIF)

ISO 19115/19139 - geospatial data



XML



XML: a set of hierarchical custom elements for a particular community's use

Relatively verbose, and requires larger storage than some other formats

Does not have strong data typing or access control (on it's own)

Language requires some training, manual editing can be tedious

Common exchange format for web services

Easy programmatic access

Platform-independent, both human- and machine- readable

Editing tools are improving

XML Basics



XML Schema

- describes the structure of an XML document
- “XML Schema language” - also referred to as “XML Schema Definition” or XSD

So XML is (kind of) a language, and EML is a dialect

A community will write up its *specification* for a *standard* in XML *Schema*

XML Basics

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <eml:eml xmlns:eml="eml://ecoinformatics.org/eml-2.1.1"
3   xmlns:stmml="http://www.xml-cml.org/schema/stmml-1.1"
4   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
5   packageId="knb-lter-sbc.19.22" scope="system" system="knb"
6   xsi:schemaLocation="eml://ecoinformatics.org/eml-2.1.1 http://nis.lternet.edu/schemas/EML/eml-2.1.1/eml.xsd">
```

A little bit of XML vocabulary:

General:	Components:
Prolog "Root" element Validate (an XML doc)	XML element - contains text, other elements XML attribute - text only Namespace Schema, schemaLocation

Ecological Metadata Language - EML

- Based on Dublin Core, FGDC, STMML
- Rich, customized structures
- Widely used for ecological and environmental data
- Machine readable
 - Read into Matlab, SAS, R, SPSS
 - Parser and DB loader
- EML records number ~90,000



NCEAS

National Center for Ecological Analysis and Synthesis



EML Anatomy

- Access rules
- Dataset metadata
 - Identifier
 - *Title
 - *Creator(s)
 - Other metadata contributors
 - Associated parties
 - Publication date
 - Language
 - Abstract
 - Keywords
 - Intellectual rights statement
 - Geographic, temporal and taxonomic Coverage
 - *Contact
 - Publisher
 - Methods
 - Project
 - Data table (with its own sub-hierarchy)
- Other metadata

```
1 <?xml version="1.0" encoding="UTF-8"?>
2 <eml:eml xmlns:eml="http://ecoinformatics.org/eml-2.1.1"
3   xmlns:stnml="http://www.xml-cml.org/schema/stnml-1.1"
4   xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance"
5   packageId="knb-lter-sbc.19.22" scope="system" system="knb"
6   xsi:schemaLocation="eml://ecoinformatics.org/eml-2.1.1 http://nis.lternet.edu/schemas/EML/eml-2.1.1/eml.xsd">
7   <access authSystem="knb" order="allowFirst">
8     <allow> [3 lines]
12    <allow> [3 lines]
16  </access>
17  <dataset>
18    <alternateIdentifier system="http://doi.org">10.6073/pasta/62803c95783c4e771695d1c6cc3d23ac</alternateIdentifier>
19    <alternateIdentifier>knb-lter-sbc.19</alternateIdentifier>
20    <shortName>KFCF Reef quad-swath counts</shortName>
21    <title>SBC LTER: Reef: Kelp Forest Community Dynamics: Invertebrate and algal density</title>
22    <creator id="sbclter"> [11 lines]
34    <creator id="dreed"> [17 lines]
52    <pubDate>2016-09-06</pubDate>
53    <language>english</language>
54    <abstract> [2 lines]
57    <keywordSet> [2 lines]
60    <keywordSet> [13 lines]
74    <keywordSet> [3 lines]
78    <intellectualRights> [30 lines]
109   <distribution> [4 lines]
114   <coverage>
115     <geographicCoverage id="ABUR"> [9 lines]
125     <geographicCoverage id="NAPL"> [8 lines]
134     <temporalCoverage> [9 lines]
144     <taxonomicCoverage> [6 lines]
151   </coverage>
152   <contact> [13 lines]
166   <publisher> [2 lines]
169   <methods> [25 lines]
195   <project> [95 lines]
291   <dataTable id="ent01">
292     <entityName>Benthic community survey, inverts and understory algae, all years</entityName>
293     <entityDescription>abundance and size of selected species of benthic invertebrates and
294       understory algae in fixed plots along permanent transects.</entityDescription>
295     <physical> [29 lines]
325     <attributeList> [629 lines]
955     <numberOfRecords>295879</numberOfRecords>
956   </dataTable>
957 </dataset>
958 <additionalMetadata>
959 <metadata>
960 <stnml:unitList> [7 lines]
968 </metadata>
969 </additionalMetadata>
970 </eml:eml>
```

Metadata Template

❖ EDI Metadata Template (2017)¹

- Dataset Title
 - (be descriptive, more than 5 words)
- Short name or nickname you use to refer to this dataset:
- Abstract
 - (include what, why, where, when, and how)
- Investigators
 - (list in order as for a paper with e-mail addresses, organization and preferably ORCID ID, if you don't have one, get it, it's easy and free: <http://orcid.org/>) add table rows as needed

- Other personnel names and roles
 - (field crew, data entry etc. with e-mail addresses, organization and ORCID ID)

- Keywords
 - (list and separate by comma, please check out these resources <http://vocab.teresa.edu/>) Please determine one or two keywords that best describe your lab, station, and/or project (e.g., Trout Lake Station, NTL LTER, UW Center for Limnology).
- Funding of this work:
 - Add rows to table if several grants were involved, list only the main PI, start with main grant first:

- Timeframe
 - Begin date
 - End date
 - Data collection ongoing/completed
- Geographic location
 - Verbal description:
 - North bounding coordinates (decimals)
 - South bounding coordinates (decimals)
 - East bounding coordinates (decimals)
 - West bounding coordinates (decimals)
- Taxonomic species or groups
- Methods
 - (please be specific, include instrument descriptions, or point to a protocol online, if this is a data compilation please specify datasets used, preferably their DOI or URL plus general citation information)

```
knb-lter-sbc.19.xml X
eml:eml dataset
17 <dataset>
18 <alternateIdentifier>knb-lter-sbc.19</alternateIdentifier>
19 <shortName>KFCD Reef quad-swath counts</shortName>
20 <title>SBC LTER: Reef: Kelp Forest Community Dynamics: Invertebrate and algal density</title>
21 <creator id="sbclter"> [11 lines]
22 <creator id="dreed"> [17 lines]
23 <associatedParty id="sharrer"> [18 lines]
24 <pubDate>2016-09-06</pubDate>
25 <language>english</language>
26 <abstract>
27 <para>These data describe the abundance and size of a select group of about 50 species of [8 lines]
28 </para>
29 </abstract>
30 <keywordSet> [2 lines]
31 <keywordSet> [13 lines]
32 <keywordSet> [3 lines]
33 <keywordSet> [5 lines]
34 <keywordSet> [3 lines]
35 <keywordSet> [3 lines]
36 <intellectualRights>
37 <para> [39 lines]
38 </para>
39 </intellectualRights>
40 <coverage>
41 <geographicCoverage id="ABUR"> [9 lines]
42 <geographicCoverage id="AHND"> [8 lines]
43 <geographicCoverage id="AQUE"> [8 lines]
44 <geographicCoverage id="BULL"> [8 lines]
45 <geographicCoverage id="CARP"> [8 lines]
46 <geographicCoverage id="GOLB"> [8 lines]
47 <geographicCoverage id="IVEE"> [8 lines]
48 <geographicCoverage id="MOHK"> [8 lines]
49 <geographicCoverage id="NAPL"> [8 lines]
50 <geographicCoverage id="SCDI"> [8 lines]
51 <geographicCoverage id="SCTW"> [8 lines]
52 </geographicCoverage>
53 <temporalCoverage>
54 <rangeOfDates> [7 lines]
55 </rangeOfDates>
56 </temporalCoverage>
57 </coverage>
58 <contact> [13 lines]
59 <publisher> [2 lines]
60 <methods>
61 <methodStep> [24 lines]
62 </methodStep>
63 </methods>
64 <project> [95 lines]
65 </project>
```

REQUEST DATA:

Available Online:

Download data after acceptance of SBC LTER Data Use Agreement:

DOWNLOAD DATA: LTE Benthic inverts and understory algae, all years**Name:** LTE Benthic inverts and understory algae, all years**Description:** Abundance and size of selected species of benthic invertebrates and understory algae in fixed plots along permanent transects in kelp removal experiment**Number of Records:** 271982**Number of Columns:** 29**Table Structure:****Object Name:** LTE_Quad_Swath_All_Years_20160315.csv**Size:** 55 megabyte

Text Format:	Number of Header Lines:	1					
	Record Delimiter:	\r\n					
	Orientation:	column					
	Simple Delimited:	<table border="1"> <tr> <td>Field Delimiter:</td> <td>,</td> </tr> <tr> <td>Collapse Delimiters:</td> <td>no</td> </tr> <tr> <td>Quote Character:</td> <td>"</td> </tr> </table>	Field Delimiter:	,	Collapse Delimiters:	no	Quote Character:
Field Delimiter:	,						
Collapse Delimiters:	no						
Quote Character:	"						

Table Column Descriptions

	Year	Month	Date	Site code	Transect name	Treatment	Quadrat Name	Transect side	Species code	Count	Observer code
Column Name	YEAR	MONTH	DATE	SITE	TRANSECT	TREATMENT	QUAD	SIDE	SP_CODE	COUNT	OBS_CODE
Definition	Calendar year	Month of data collection	Date of data collection	Code for reef site. See list of codes for meaning	40m transects defined by six permanent markers (stainless steel eyebolts or rebar stakes) at 0. 8. 16.	Experimental treatment, frequency of removal	Data collected in quadrats of 20mx1m called 20 or 40 (0-20, or 20-40m, also known as swaths), or 1mx1m quadrats at	Side of the transect the data was taken on. Either 'I' (inshore) or 'O' (offshore).	2-4 letter code used by SBC LTER. Numbers are used to differentiate size classes for Macrocyctis.	Number of individuals in the area surveyed.	Numeric code indicating the SBCLTER data collector

Resources



Dataset in examples, “SBC LTER Time-Series of Kelp Forest Invertebrate and Algal Density”
<https://portal.edirepository.org/nis/mapbrowse?scope=knb-lter-sbc&identifier=19>

Best Practices for EML datasets (originated by LTER IMC)

EML Project

Download the current release, read basic documentation: <https://www.nceas.ucsb.edu/eoinfo/tools>

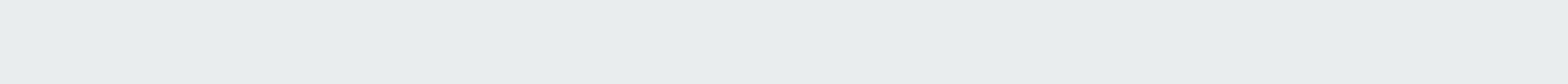
Maintained in GitHub: <https://github.com/NCEAS/eml>

Find resources at:

<https://environmentaldatainitiative.org> > Resources > 5 Phases of data publishing > Phase 3. create EML metadata

<https://www.liquid-technologies.com/xml-tutorial>

<https://www.w3schools.com/xml/>





Data Table metadata

Some components can be automated
Attribute list: your knowledge of the data is essential

Data Table

- Column name: exactly as it appears in the dataset. Please avoid special characters, dashes and spaces.
- Description: please be specific, it can be lengthy
- Unit: please avoid special characters and describe units in this pattern:
e.g. microSiemenPerCentimeter, microgramsPerLiter, absorptionPerMolePerCentimeter
- Code explanation: if you use codes in your column, please explain in this way: e.g. LR=Little Rock Lake, A=Sample suspect, J=Nonstandard routine followed
- Data format: please tell us exactly how the date and time is formatted: e.g. mm/dd/yyyy hh:mm:ss plus the time zone and whether or not daylight savings was observed.
- If a code for 'no data' is used, please specify: e.g. -99999
- Please add rows as needed

Notes and Comments

-
-
-

```
knb-ltr-sbc.19.xml X
em:eml dataset dataTable attributeList
416 <dataTable id="ent01">
417 <entityName>Benthic community survey, inverts and understory algae, all years</entityName>
418 <entityDescription>abundance and size of selected species of benthic invertebrates
419 and understory algae in fixed plots along permanent transects.</entityDescription>
440 <physical> [20 lines]
441 <attributeList>
442 <attribute id="ent1.att3"> [15 lines]
443 <attribute id="ent1.att4"> [96 lines]
444 <attribute id="ent1.att8">
445 <attributeName>sp_code</attributeName>
446 <attributeLabel>Species code</attributeLabel>
447 <attributeDefinition>2-4 letter code generally composed of the first letter of the
448 genus (G) name and the first three letters of the species (S) name (format=GSSS)</attributeDefinition>
449 <storageType typeSystem="http://www.w3.org/2001/XMLSchema-datatypes">string</storageType>
450 <measurementScale>
451 <nominal>
452 <nonNumericDomain>
453 <textDomain>
454 <definition>any text</definition>
455 </textDomain>
456 </nonNumericDomain>
457 </nominal>
458 </measurementScale>
459 <missingValueCode>
460 <code>-99999</code>
461 <codeExplanation>no information available</codeExplanation>
462 </missingValueCode>
463 </attribute>
464 <attribute id="ent1.att9">
465 <attributeName>count</attributeName>
466 <attributeLabel>Count</attributeLabel>
467 <attributeDefinition>number of individuals in area</attributeDefinition>
468 <storageType typeSystem="http://www.w3.org/2001/XMLSchema-datatypes">integer</storageType>
469 <measurementScale>
470 <ratio>
471 <unit>
472 <customUnit>number</customUnit>
473 </unit>
474 <precision>1</precision>
475 <numericDomain>
476 <numberType>real</numberType>
477 </numericDomain>
478 </ratio>
479 </measurementScale>
480 <missingValueCode>
481 <code>-99999</code>
482 <codeExplanation>no information available</codeExplanation>
483 </missingValueCode>
484 </attribute>
485 </attributeList>
486 <numberOfRecords>295879</numberOfRecords>
487 </dataTable>
```

Attributes & Units



- Attribute: a “property” of an object (data table)
 - In databases, a table column is called an “attribute”
 - Often referred to as “variables”, “parameters”, “columns” or “field names”
- Unit: a particular physical quantity
 - Defined and adopted by convention
 - Comparable

To describe a data table you need a moderate understanding of

- how to define the table’s attributes,*
- when and how to define a unit, and*
- the relationship between the two.*

EML Attribute components



1 Attribute Name:

Usually the name you would give that column in a script

2 Attribute Label:

Longer, for display, Use whole words, capitals, etc.

3 Attribute Definition:

As complete and unambiguous as you need them to be, for the data to be understood

4 Measurement Scale:

- Nominal - *attribute can be considered a category*
- Ordinal - *categories that have a logical or ordered relationship to one another*
- Interval - *the magnitude between the steps is known; equidistant points*
- Ratio - *have a meaningful zero, which allows ratios between values to have meaning*
- Datetime - *Gregorian dates and times*

5 Unit:

Interval and Ratio measurements only

Choose from: <http://unit.lternet.edu>

EML Attributes - Measurement Scale



Nominal	<i>Values are members of a category</i>	string	Place and taxon names, coded values (eg, 1=male, 2=female), text comments
Ordinal	<i>Nominal categories that have a logical or ordered relationship to one another</i>	string	Academic grades, quality rankings (eg, 1=high, 2=medium, 3=low)
Interval	<i>Ordinal, but the magnitude between the steps is known; equidistant points</i>	numeric	Celsius scale, pH
Ratio	<i>Interval, with a meaningful zero, so ratios between values to have meaning</i>	numeric	Temperature in Kelvin, lengths, concentrations, organism densities
Datetime	<i>Gregorian dates and times</i>	datetime	Points in time, e.g., with formats like YYYY-MM-DD, hh:mm:ss.s

EML Attributes - code lists

Nominal	<i>Values are members of a category</i>	Place and taxon names, coded values (eg, 1=male, 2=female)	<pre><codeDefinition> <code>ABUR</code> <definition>Arroyo Burro Reef</definition> </codeDefinition> <codeDefinition> <code>NAPL</code> <definition>Naples Reef</definition> </codeDefinition> ...</pre>
Ordinal	<i>Nominal categories that have a logical or ordered relationship to one another</i>	Academic grades, quality rankings (eg, 1=high, 2=medium, 3=low)	<pre><codeDefinition> <code>A</code> <definition>scored higher than 90%</definition> </codeDefinition> <codeDefinition> <code>B</code> <definition>score 80 - 89%</definition> </codeDefinition> <codeDefinition> <code>C</code> <definition>score 70 - 79%</definition> </codeDefinition> ...</pre>