



Ecological Metadata

John Porter University of Virginia



DO UNTO DATA BEFORE IT DOES UNTO YOU.





Slide from Susan Stafford



The Ecoinformatics Challenge:



- Can we make information available to ecologists:
 - In ways they can locate the information they need?
 - With information in forms they can readily use?
- How can we assure that the information is current and accurate?





Metadata is needed to reach our goals:

- Enabling new scientific approaches by making data available
 - Regional and Global Analyses
 - Multi-disciplinary syntheses
 - Multi-factorial analyses
- Creating data that is still useful 20-years or more into the future



Slide from William Michener & James Brunt

V C R LTER



Increasing value of data





Evolution of Data Sharing - Traditional Model









- Data Discovery
 - Scientists need to be able to identify important data sets
- Data Retrieval
 - Scientists need know how and where to access data
- Data Use
 - Scientists need to know enough details about how the how the data were collected and stored
- Data Archiving
 - Ecological data can grow more valuable with time, but only if the critical information required to retrieve and interpret the data remains available





	Levels of secondary data utilization and associated metadata content			
Metadata descriptor classes	Level I: exchange with expert colleague	Level II: searchable and third party data reuse	Level III: publishable and auditable	
I. Data set descriptors II. Research origin descriptors III. Data set status and accessibility	Х	X X X	X X X	
 IV. Data structural descriptors V. Supplemental descriptors 	х	x	X X	

Michener et al. 1997 – NONGEOSPATIAL METADATA FOR THE ECOLOGICAL SCIENCES, Ecological Applications Vol. 7, No. 1, pp. 330–342.





- Entropy is not easily conquered! The natural tendency is towards increasing disorder (2nd law of thermodynamics)
- The diversity of ecological data – Genome data is much simpler!
- The diversity of users and uses to which data will be put







- General Standards

 Dublin Core
 designed for publications

 Spatial Metadata

 ISO 19115 International Standard
 FGDC (Federal Geographic Data Committee)
- Taxonomic & Collection Data
 Darwin Core
- Catalog
 - Directory Interchange Format (DIF)





- Each of the existing metadata standards includes:
 - Who information on who to contact
 - What a description of the available item
- Some include:
 - Where geographical locations
 - When dates and times
- Standards vary widely in the degree of structure from free text to XML and database schema

The wonderful thing about standards is that there are so many of them to choose from..... - Anonymous



Metadata – Two examples



Benson, Barbara Trout Lake Temperature Water temperature data was collected hourly at Trout Lake **Trout Lake Temperature** Water temperature data was collected hourly at Trout Lake from January 1, 2005 to December 31, 2005 Collected by Barbara

Benson

Both metadata documents are readable by humans, can't be processed by a computer into new forms because both the <u>content</u> and <u>structure</u> are different



Metadata – Content Standardized



Originator: Benson, Barbara Trout Lake Temperature Water temperature data was collected hourly at Trout Lake Start Date: January 1, 2005 End Date: December 31, 2005 Trout Lake Temperature Water temperature data was collected hourly at Trout Lake Time period: January 1, 2005 to December 31, 2005

Originator: Benson, Barbara

Now both metadata documents have the same content – title, originator, etc. but it still can't be automatically processed by a computer because the <u>structure</u> is different

Metadata – Structure Standardized



- <title>Trout Lake Temperature</title> <originator>
- <namelast>Benson</namelast> <namefirst>Barbara</namefirst>
- </originator>
- <abstract>Water temperature data was collected hourly at Trout Lake</abstract>

<date>

<start>January 1, 2005</start> <end>December 31, 2005</end> </date>

<title>Trout Lake Temperature</title> <originator> <namelast>Benson</namelast> <namefirst>Barbara</namefirst> </originator> <abstract>Water temperature data was collected hourly at Trout Lake </abstract> <date> <start>January 1, 2005</start> <end>December 31, 2005</end> </date>

With standardized content and structure, computers can automatically extract information from the metadata!



Why Ecological Metadata Language?



- Other existing standards do not support the development of interoperable systems
 - Lack needed structure
 - Lack needed content to support automated access to common ecological data formats (e.g., delimited text, formatted text, SQL databases)



- Currently EML version 2.01 is being implemented by US LTER sites
- EML uses eXtensible Markup Language (XML) schema to define content and structure
- The EML standard is made up of a number of different modules





- eml-resource -- bibliographic info
- eml-party -- people and organizations
- eml-entity -- file/object info
- eml-attribute -- variable/attribute info
- eml-access -- access control
- eml-distribution -- distribution info



 EML -Ecological Metadata
 Language – different EML
 modules are
 nested within
 other modules



Full list of EML 2.01 Modules



- eml
- eml-access
- eml-attribute
- eml-constraint
- eml-coverage
- eml-dataset
- eml-dataTable
- eml-entity
- eml-literature
- eml-methods
- eml-party

- eml-physical
- eml-project
- eml-protocol
- eml-resource
- eml-software
- eml-spatialRaster
- eml-spatialReference
- eml-spatialVector
- eml-storedProcedure
- eml-text
- eml-view

TIOLLOGGGLOD

Overview of EML modules and their use

2.1. Module Overview Foreword

2.2. Root-level structure

- 2.2.1. The emi module A metadata container
- 2.2.2. The emi-resource module Base information for all resources

2.3. Top-level resources

- 2.3.1. The eml-dataset module Dataset specific information
- 2.3.2. The emi-literature module Citation specific information
- 2.3.3. The eml-software module Software specific information
- 2.3.4. The eml-protocol module Research protocol specific information
- 2.4. Supporting Modules Adding detail to top-level resources
 - 2.4.1. The emi-access module Access control rules for resources
 - 2.4.2. The emi-physical module Physical file format
 - 2.4.3. The eml-party module People and organization information
 - 2.4.4. The emi-coverage module Geographic, temporal, and taxonomic extents of resources
 - 2.4.5. The eml-project module Research context information for resources
 - 2.4.6. The eml-methods module Methodological information for resources
- 2.5. Data organization Modules describing dataset structures
 - 2.5.1. The eml-entity module Entity level information within datasets
 - 2.5.2. The eml-attribute module Attribute level information within dataset entities
 - 2.5.3. The emi-constraint module Relationships among and within dataset entities
- 2.6. Entity types Detailed information for discipline specific entities
 - 2.6.1. The emi-dataTable module Logical information about data table entities
 - 2.6.2. The emi-spatialRaster module Logical information about regularly gridded geospatial image data
 - 2.6.3. The emi-spatialVector module Logical information about non-gridded geospatial image data
 - 2.6.4. Schema for validating spatial referencing descriptions.
 - 2.6.5. The emi-storedProcedure module Data tables resulting from procedures stored in a database
 - 2.6.6. The emi-view module Data tables resulting from a database query
- 2.7. Utility modules Metadata documentation enhancements
 - 2.7.1. The emi-text module Text field formatting
 - 2.7.2. Dependency Chart
-

EML Documentation describes each module and how it is used

http://www.ecoinformatics.org

V

C

~

software Type			
eml-dataset Do	cumentation - Nets	scape Browser	Eile Edit 🛛 🗹 🥓 Microphone 👰 Tools 😰 📮
🖳 🕑 💽 🔅	o 🕋 🗔	est Practices LTER 🗢 SE	ARCH) 🕒 http://knb.ecoinformatics.org/software/eml/eml-2.0.1/eml-dat 😳 😢 🗾
Personal V G 74° C	🕑 🖾 Setup E-m	ail 오 🛛 🚁 🕫 Net	scape.com 🛅 Inside Netscape
🕑 🙉 🚊 Index of /bj.	🗡 💽 2005 LTER	R I 🗡 📑 ESA > N	neet × 🎦 http:/7.html × 🎦 ecoinformat × 🔂 eml-dataset Documentation 🔤 🏹 🗂
Content of this field:			Description of this field:
Elements:	Use:	How many:	DatasetType is the base type for the dataset element. The dataset field
A choice of (encompasses all information about a single dataset. A dataset is defined as all of the information describing a data collection event. This event may take place
A sequence of (over some period of time and include many actual collections (a time series or
res:ResourceGroup			remote sensing application) or it could be just one actual collection (a day in the
purpose	optional		field).
maintenance	optional		
contact	required	unbounded	
publisher	optional		
pubPlace	optional		
methods	optional	FMI c	locuments consist of nested
project	optional		
access	optional	modu	les each of which have OPTIONAL
A choice of (an extine d	mouu	ICS CACH OF WHICH HAVE OF FIONAL
	requirea	and D	FOULPED alamants
CRatialBactor	required		
	requireu		
snatialVector	required		
OR	requirea	Horo	is the EML specification for a
storedProcedure	required		is the live specification for a
OR	. equilea	datas	ot
view	reauired	ualdS	
OR			
otherEntity	required		
Eind: EMI		Eind Previous 🗐 Hid	abliabt 🔲 Match case

ds:DatasetType

cit:CitationType

dataset 🗗 –∕ 🗗 🖽

citation ⊟+(~i+)⊞

emi -----









EML documents can be created in a variety of ways. These include:

- General XML Editors (e.g., oXygen, XMLspy)
- Customized XML Editors (e.g., Morpho)
- Output from metadata stored in databases
- Metadata stored in specialized spreadsheets



	l xls_eml_01_SAMPLE.xls							
1	2	В	C	D				
	17	I. LTER Dataset Information		· · · · · · · · · · · · · · · · · · ·				
	18	I TER site acronym	TECE					
	19	Metacat nackage ID	UT ND Grahl 001.1					
	20	Dataset LTER Identification Number	LT ND Grahl 001					
	21	Dataset Title	Taylor Slough Water Quality Data for 2000-2001					
Г	· 22	Dataset Creator Salutation	Dr.					
Ш	· 23	Dataset Creator First Name	Daniel					
Ш	· 24	Dataset Creator Last Name	Childers					
Ш	· 25	Dataset Creator Organization Name		Florida Coastal Everglades LTER Program				
Ш	· 26	Dataset Creator Position Name	4					
Ш	· 27	Dataset Creator Mail Street Address	SERCIOE 148 Florida International University University Park	ECS 253 Florida International University U				
Ш	· 28	Dataset Creator Mail City	Miami	Miami				
Ш	· 29	Dataset Creator Mail State	FL	FL				
Ш	· 30	Dataset Creator Mail Zip Code	33199	33199				
Ш	· 31	Dataset Creator Mail Country	USA	USA				
Ш	· 32	Dataset Creator Voice Telephone	305-348-3101	305-348-6054				
Ш	· 33	Dataset Creator Facsimile Telephone	305-348-1986	305-348-4096				
Ш	· 34	Dataset Creator Electronic Mail Address	<u>childers@fiu.edu</u>	fcelter@fiu.edu				
	: 35	Dataset Creator URL	http://www.fiu.edu/~ecosyst/	http://fcelter.fiu.edu/				
E	36	+-Dataset Creator						
	37	Dataset Abstract	Water quality samples are being collected using ISCO					
	38	Dataset Keywords	inorganic nutrients	water quality				
۱.	39	Dataset Keyword I hesaurus						
Ш	· 40	Dataset Geographic Description	Shark River Slough Transect within Everglades National Park					
Ш	· 41	Dataset West Bounding Coordinate	-81.07794623					
Ш	• 42	Dataset East Bounding Coordinate	-80.72742805 Non 20145474					
Ш	43	Dataset North Bounding Coordinate	25.76145171					
Ш	44	Dataset South Bounding Coordinate	25.36462994					
Ш	40	Dataset Beginning Temporal Coverage Date	2000-10-31					
Ш	. 40	Dataset Ending Temporal Coverage Date	2001-11-13					
	. 48	Dataset Taxon Rank Value						
Ш	. 40	Dataset Common Taxon Names	Evcal snraadshaat	for				
Ľ	49	+ Dataset Coverage (Coographic Temperal Tayonomic)	LACCI Spicausiicet					
-	51	Patasat Intellectual Pights						
Г	. 57	Dataset Intellectual Rights	+ MI - Processed by a r	program to				
Ш	. 53	Dataset Offine Medium Name						
Ш	. 50	Dataset Offine Medium Name	croate EML document					
	. 54	Dataset Offline Medium Density Unite						
	. 50	Dataset Online Medium Velume						
	. 50	Dataset Offline Medium Format						
	1 59	+ Dataset Distribution						
Г	. 50	Dataset Associated Party First Name	Tim					
	. 60	Dataset Accordated Party Last Name	Grahl					
H	- + +	M General Metadata / MethodsCitation / MethodsProtocol / ResearchProjects /	DataTable / References / IM Use Only / Units IM Use Only /					

Specialized Tool -Morpho



Sample Wizard Pages

Owner	Details

Salutation: First Name:

Organization:

Position Name: Address 1: Address 2: City:

Postal Code: Phone: Email:

You can pick from one of the earlier er

X

One of the

three required

ntries that you have made.		30	
		Enter a description	on of the geographic coverage. Enter a general description of the geographic area in which the data were be a simple place name (e.g., Santa Barbara) or a fuller description.
		Description:	
		Set the geograph box' containing the fractional degrees.	hic coordinates which bound the coverage: Latitude and longitude values are used to create a 'bounding region of interest. Drag or click on the map and then edit the text boxes if necessary. [Default entries are in To enter in degrees/minutes/seconds, simply type a space between the degrees, minutes, and seconds values]
Sta	ite:	Bounding Box:	S.S.W BOSE 45.0 S Zoom In Zoom Out O Box Tool O Point Tool
		Named Regions:	[Belize] Jaguar Creek [Namibia] Gobabeb Training and Research Centre Angelo Coast Range Reserve UCNRS Ano Nuevo Island Reserve UCNRS Antic LTER (ARC) Battimore Ecosystem Study LTER (BES) Roders Merine Reserve LICNRS

OK Cancel

Using EML - Search



LTER



Template

Title:

Date

 Depending upon the XSLT template used, the same EML document can be displayed many different ways







- For metadata to be used to automatically generate or configure analysis tools, it needs to contain the details required by those tools
 - Lists of attributes/variables
 - Information on storage or access formats
- Need to provide a reliable structure so that those details can be automatically extracted



Using EML

• FMI transformations are not restricted to documents for viewing. They can also be used to create programs for processing the associated data







- Researchers use a wide variety of tools to analyze data
 - Spreadsheets
 - Statistical Packages
- Many ecological datasets are available as text files
- Therefore, we need a way to facilitate the ingestion of data in text files into different tools used by researchers



We need to be able to translate an EML Document into a Statistical Program



EML Document

<?xml version="1.0" encoding="UTF-8"?> <eml:eml packageld="knb-lter-vcr.76.1" system="VCR" xmlns:ds="eml://ecoinformatics.org/dataset-2.0.1" xmlns:eml="eml://ecoinformatics.org/eml-2.0.1" xmlns:stmml="http://www.xmlcml.org/schema/stmml" xmlns:xsi="http://www.w3.org/2001/XMLSchemainstance" xsi:schemaLocation="eml://ecoinformatics.org/e ml-2.0.1 http://gce-Iter.marsci.uga.edu/Iter/files/schemas/eml-201/eml.xsd"> <dataset id="76" system="VCR"> <alternateIdentifier>VCR00073</alternateIdentifier <title>Water Quality - physical data</title> <creator> <individualName> <salutation>Dr.</salutation> <givenName>Robert</givenName> <surName>Christian</surName> </individualName> <address> <deliveryPoint>East Carolina University, Department of Biology</deliveryPoint>

Statistical Program



/ FILE="PUT-LOCAL-PATH-TO-DATA-FILE-HERE" /ARRANGEMENT=Delimited

/DELIMITERS=","

/QUALIFIER=""

/VARIABLES=

SITE A DATE A TIME F10.2 SCTTEMP F SCTSAL F SCTCOND F10.2 TEMP F REFRSAL F DOTEMP F DO F SECCHI F DEPTH F WIND A .

execute.

VAR LABELS	SITE 'SAMPLE SITE OR STATION- ' .
VAR LABELS	DATE 'DATE SAMPLE COLLECTED- ' .
VAR LABELS	TIME 'TIME SAMPLE COLLECTED- none' .
VAR LABELS	SCTTEMP 'TEMPERATURE BY SCT- DEGREES C' .
VAR LABELS	SCTSAL 'SALINITY BY SCT- PPT' .
VAR LABELS mg/L'	DOTEMP 'TEMPERATURE BY DISSOLVED OXYGEN METER-
VAR LABELS mg/L'	DO 'DISSOLVED OXYGEN BY DISSOLVED OXYGEN METER-
VAR LABELS	SECCHI 'SECCHI DEPTH- cm' .
VAR LABELS	DEPTH 'WATER DEPTH- m' .
VAR LABELS	WIND 'WIND SPEED AND DIRECTION- ' .
Frequencies	variables=SITE /order=analysis.
Frequencies	variables=DATE /order=analysis.
Frequencies	variables=WIND /order=analysis.
Descriptives	variables=TIME .
Descriptives	variables=SCTTEMP.











ne -





- US LTER Network
- Organization of Biological Field Stations (OBFS)
- Ecological Society of America, Long-Term Studies Section
- Cyberinfrastructure Research Projects
 - Knowledge Network for Biocomplexity (KNB)
 - Science Environment for Ecological Knowledge (SEEK)

Long- Term	LTSS Data Registry				
Section	<u>LTSS Home</u>	<u>Registry Home</u>	<u>Register a</u> <u>New Data Set</u>	Search for Data	
Data Registry Use this form to subm Please have a look at your browser's Reload If you have any quest esaadmin@nceas.ucs *Denotes a required BASIC INFORMATION *First *Las	/ Form hit a new data set des the <u>Guide for Compl</u> d/Refresh function to r tions, comments or pro <u>sb.edu</u> . I field.	cription for inclusion in the eting the Data Registry Fo make sure you see the late oblems regarding this form	The Long-Te Section Data includes: •Identification •Data set own •Abstract •Keywords •Coverage (sp temporal, taxo •Data Collection	rm Studies a Registry information her batial, bnomic) on Methods	e
*Data S *Organization	et Title		•Distribution I	nformation	
PRINCIPAL DATA S	ET OWNER (What's t	<u>his?)</u>			<u>Hide</u>
*First	t Name				

- The U.S. LTER sites have developed a "Best Practices" guide for US LTER sites
 - It calls for a tiered approach with sites moving from EML suitable for data identification to full-content EML metadata that supports interoperability

oShapes 🔻 📐 🔌

1/33

At 1"

Ln 1 Col 1

http://cvs.lternet.edu/cgibin/viewcvs.cgi/emlbestpractices/ emlbestpractices-1.0/

LTER



_ <u>©</u>¥

REC TRK EXT OVR English (U.S.





- Minimum content for adequate data set discovery in a general cataloging system or repository
 - title
 - creator
 - contact
 - publisher
 - pubDate
 - keywords
 - abstract (recommended)
 - dataset/distribution (i.e. url for general dataset information)

EML Best Practices for LTER Sites - Oct. 2004





- Level 1 content, plus coverage information to support targeted searches, adding elements:
 - -Geographic Coverage
 - -Taxonomic Coverage
 - -Temporal Coverage





- Level 2 content, plus data set details to enable enduser evaluation of the methodology and data entities, adding elements:
 - Intellectual Rights
 - project
 - methods
 - dataTable/entityGroup
 - dataTable/attributes

EML Best Practices for LTER Sites – Oct. 2004





- Level 3 content plus data access details to support automated data retrieval, adding elements:
 - -access
 - -physical

EML Best Practices for LTER Sites – Oct. 2004





- Level 4 content plus complete attribute and quality control details to support computer-assisted data integration and re-sampling, adding elements:
 - -Attribute List (full descriptions)
 - -Constraint
 - –Quality Control



More Information on EML and EML-related Projects

http://www.ecoinformatics.org















