

Data Management for Synthesis

Ben Leinfelder
Matthew B. Jones

National Center for Ecological Analysis and Synthesis (NCEAS)
University of California Santa Barbara

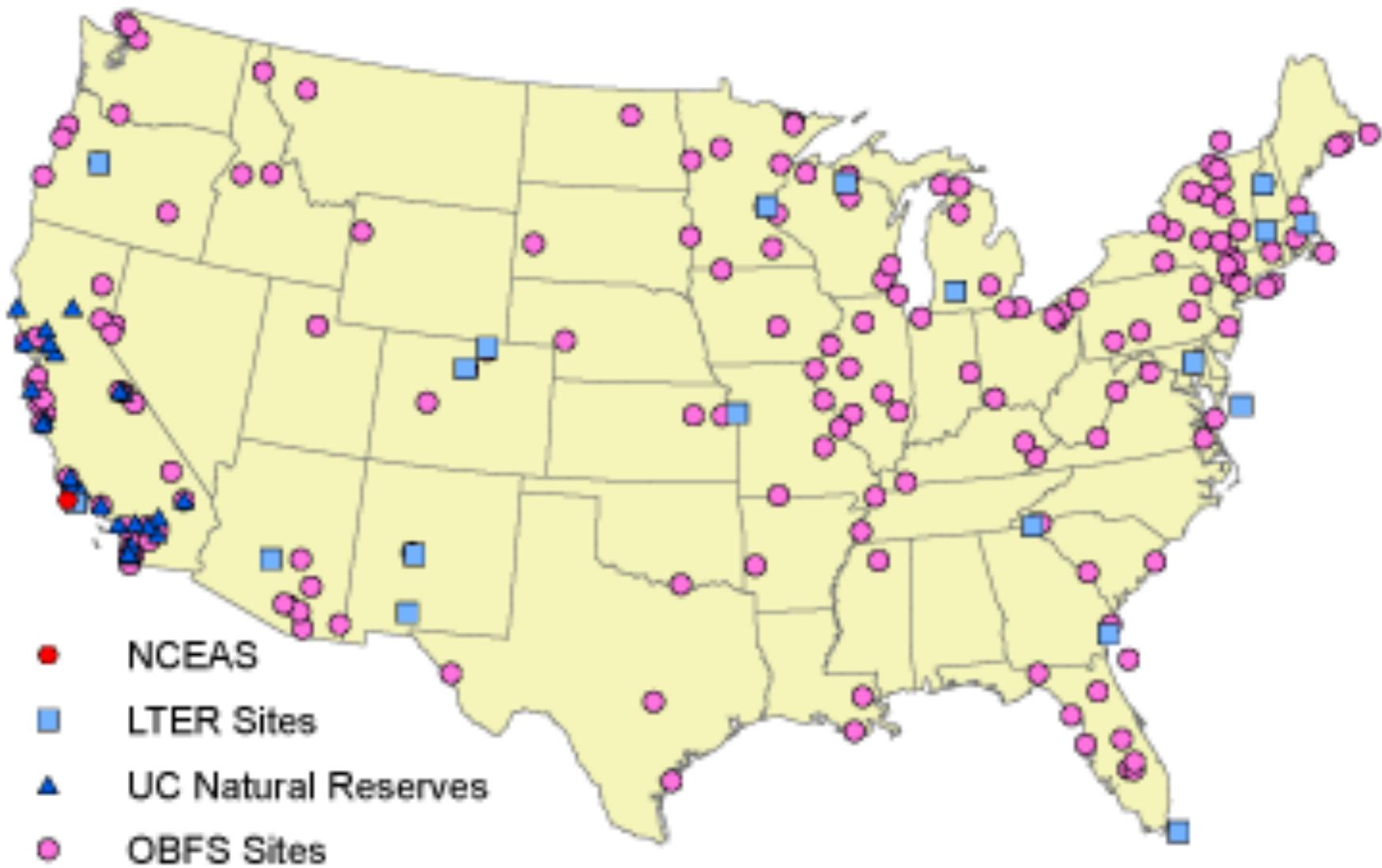


GLEON 16 Data workshop
October 27, 2014

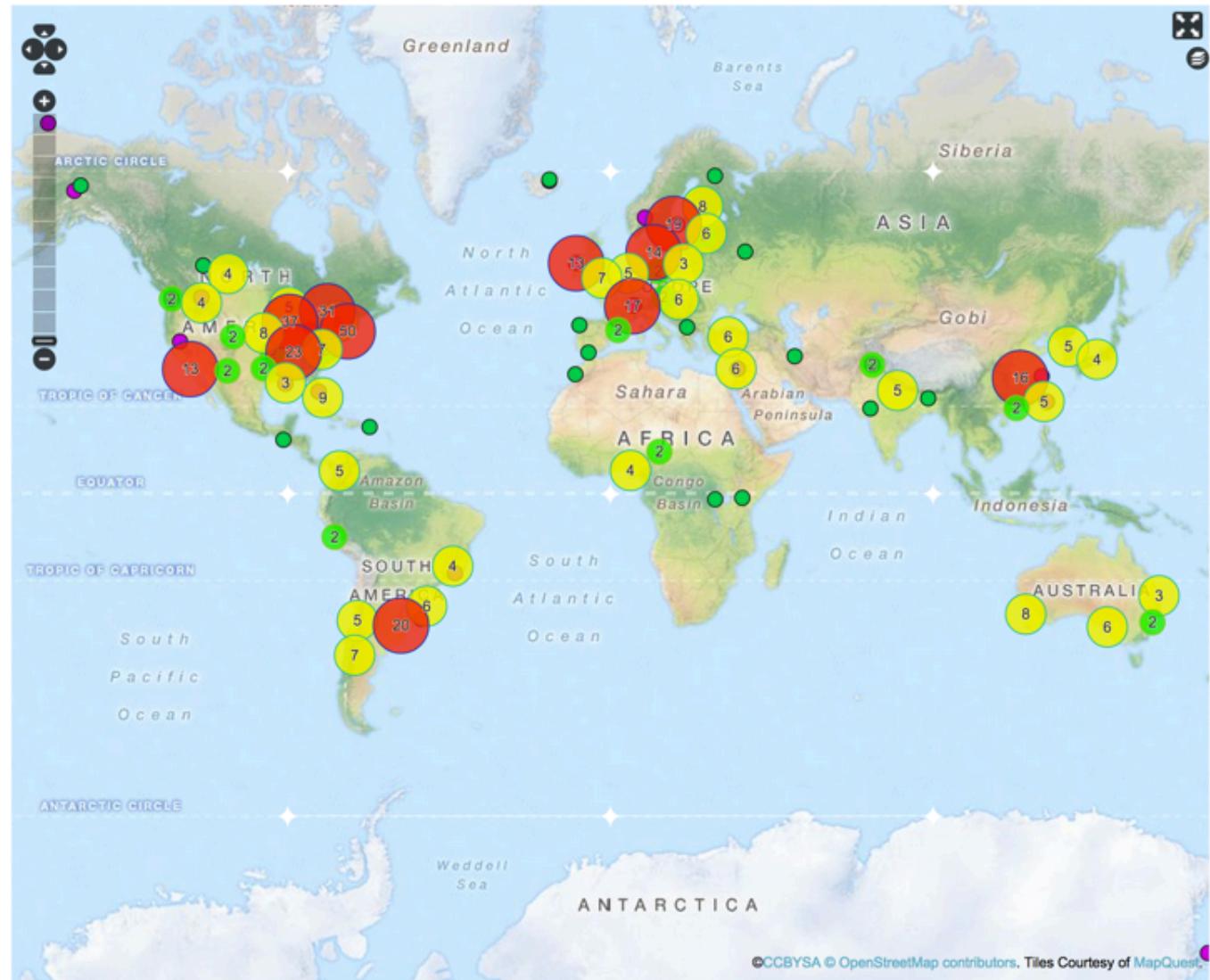
Barriers to Synthesis

- Data not preserved
 - Tiny proportion of ecological data are readily available
- Dispersed, isolated repositories
 - Each community has its own; disconnected; underutilized
- Lack of software interoperability
 - Metacat, DSpace, Mercury, iRODS, XMCat, OPeNDAP, ...
- Heterogeneous data
 - Many data formats, metadata formats, and varying semantics

Dispersed data from field stations



GLEON Sites



“...understand, predict
and communicate the
role and response of
lakes in a changing
global environment”

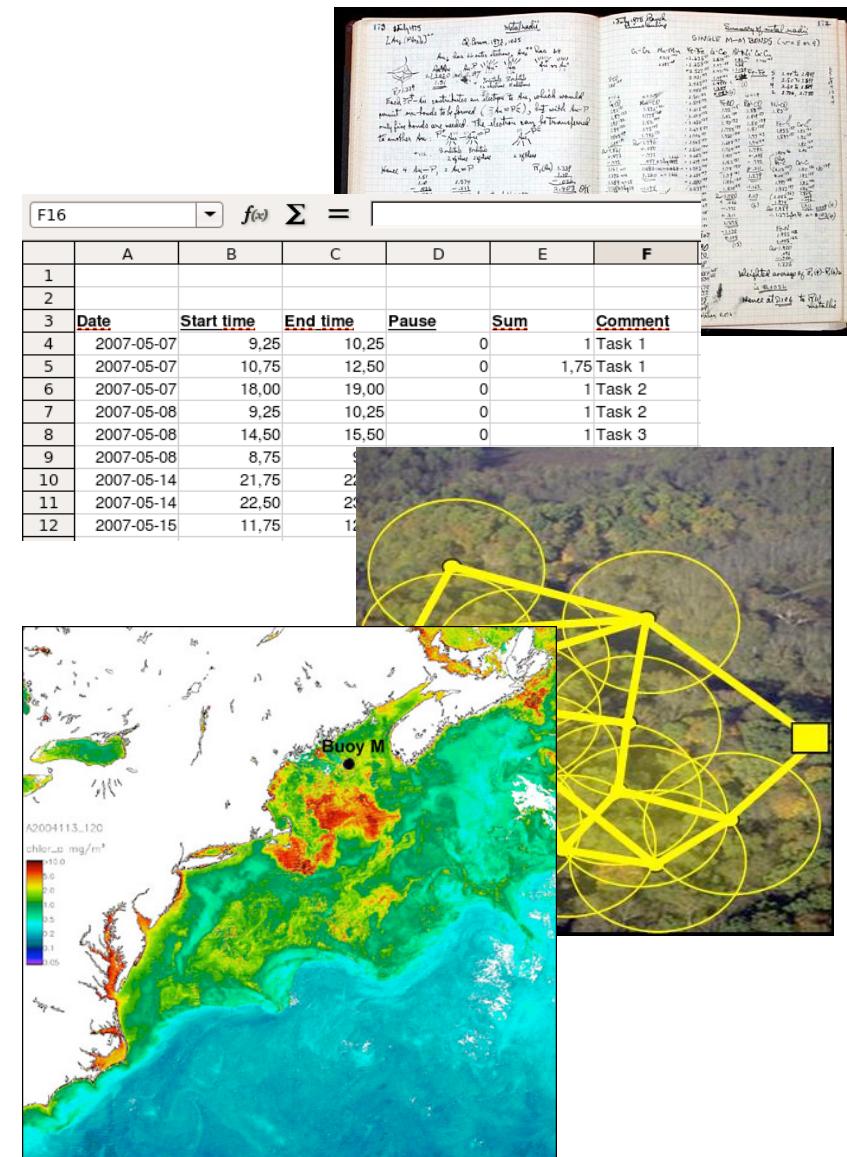
What data are in scope?

- Biological
 - e.g., Ecosystem, Organism, Population, Species, Community, Biome, Gene

- Environmental
 - e.g., Atmospheric, Chemical, Ecological, Hydrological, Oceanographic, Physical

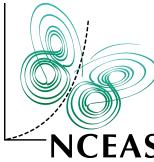
- Social
 - e.g., Land use, human population

- Economic
 - e.g., trade, ecosystem services, resource extraction



Metadata and data heterogeneity

- Every community has
 - many data schemas
 - one for each project and person
 - many data formats
 - ASCII, NetCDF, HDF, GeoTiff, ...
 - many metadata schemas
 - Biological Data Profile, Darwin Core, Dublin Core, Ecological Metadata Language, Open GIS schemas, ISO Schemas, ...
- Accepting this heterogeneity is critical



Biodiversity data heterogeneity

Space

Time

Taxa

(a) Georgia Coastal Ecosystems LTER (EML)

(b) *Mephitis mephitis* specimen record (DarwinCore)

(c) NOAA Ocean Buoy Data Station 46069 - South Santa Rosa Island, CA

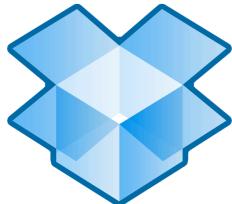
YYYY	MM	DD	hh	mm	WD	WSPD	GST	WVHT	DPD	APD	MWD	BAR	ATMP	WTMP	DEWP	VIS	TIDE
2006	01	01	00	00	284	6.9	8.2	3.50	13.79	8.11	275	1013.2	15.5	14.3	999.0	99.0	99.00
2006	01	01	01	00	287	6.0	7.6	3.13	13.79	8.08	271	1013.5	15.3	14.3	999.0	99.0	99.00
2006	01	01	02	00	276	4.3	5.6	3.20	13.79	8.49	272	1013.7	14.9	14.4	999.0	99.0	99.00

(d) Macroecological data for fossil occurrences (Paleobiology Database)

Synthesis requires sharing data

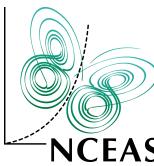
- NCEAS' approach to data sharing
 - Deal with data heterogeneity
 - Distributed data management
 - Centralized search
 - Semi-automated analysis tools
- A grass-roots network with global partners
 - NCEAS, LTER, iLTER, PISCO, ESA, SanParks, SAEON, TERN, TEAM, ...

Software diversity



GMN

python
powered



Solutions

- Preserve data
- Adopt standards

<EML/>



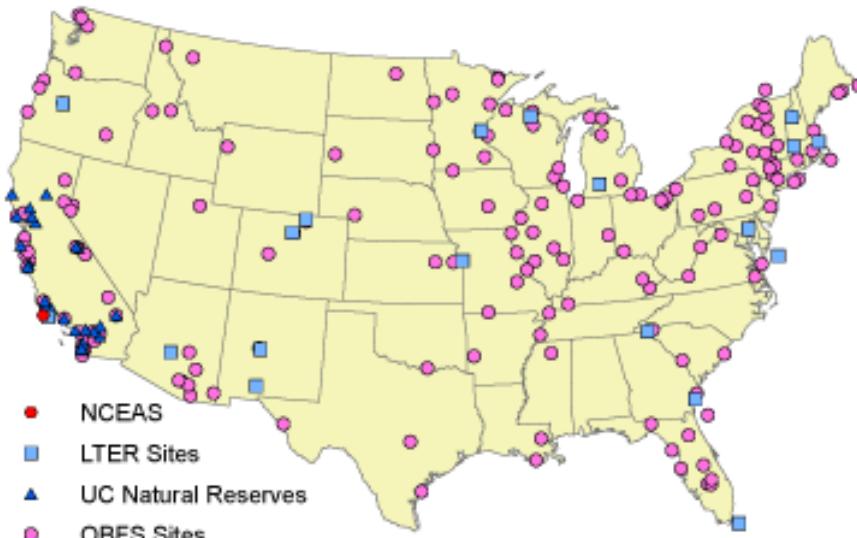
- Create networks
- Create interoperable software





PRESERVE DATA

Data in the KNB



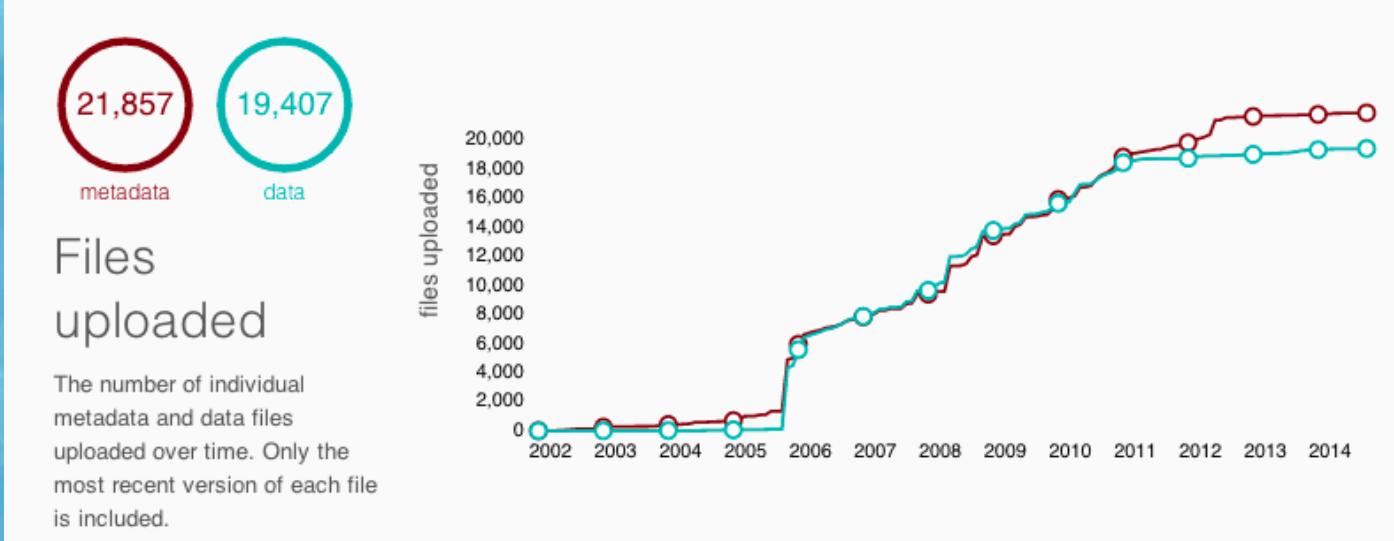
- Diverse Contributors
- Individual investigators
- Field stations and networks
- Government agencies
- Non-profit partnerships
- Scientific Societies
- Synthesis centers

Data Types

- Ecological
- Environmental
- Demographic
- Social/Legal/Economic



KNB Data Distribution



Metacat Data Server

- Data and metadata management
- Store, search, and document data
- Customizable web-based search interface
- Web metadata entry tool
- DOI Support

- Runs on Linux, Windows, MacOS
- Replication capabilities
- Postgres or Oracle backend
- OAI-PMH harvester
- GPL open source license



The screenshot shows the knb website interface. At the top, there are navigation links for Map (selected), Satellite, ABOUT, DATA, SHARE, TOOLS, a search bar "Search for data" with a magnifying glass icon, and a "SIGN IN" button. Below the header is a map of South America and Central America, with red location markers and numbers indicating data packages. A sidebar on the left contains a "FILTER" section with various search and filter options: "Clear all filters", "Anything" (with a search icon), "Data attribute (density, length, ...)" (with a search icon), "Only results with data" (with a checkbox and a help icon), "Creator" (with a search icon), a date range selector from "1900" to "2014" (with a help icon), "Data covers" (with a checkbox and a help icon), "Published between" (with a checkbox and a help icon), "Taxon" (with a search icon), "Location" (with a search icon), and a dropdown "Only results with all spatial coverage inside the map" (with a help icon). On the right, a large box displays the title "Mapping 1 to 112 of 112 data packages" and a sorting option "Sort by Most recent". Below this are six data package entries, each with an info icon, a location pin, and a folder icon:

- Amy Thissell, and Jennifer Balch. 2014. Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013). doi:10.5063/F1XW4GQH 78 views
- Amy Thissell, and Jennifer Balch. 2014. Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Flammability Metrics (2013). doi:10.5063/F12N5066 5 views
- Amy Thissell, and Jennifer Balch. 2014. Testing Amazon Transitional Forest Leaf Flammability: Leaf Traits (2013). doi:10.5063/F16D5QX4 9 views
- University of Technology, Sydney, and Edd Hammill. 2014. Mosquito distributions. knb.464.2 55 views
- University of Colorado, Boulder, and Amy Thissell. 2014. Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013). knb.460.2 15 views
- University of Colorado, Boulder, and Amy Thissell. 2014. Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013). knb.458.1 36 views
- Edmar da Silva Prado, Albanita de Jesus Rodrigues da Silva, and Carolina Volkmer de Castilho. 2014. Perfil químico de *Chaetocnemus schomburgkianus* (Kuntze). Pay & Hoffman 10 views



Data Citation



Data Set Citation:

When using this data, please cite the data package:

Kline T.

SEA: Confirming Food Web Dependencies with Stable Isotope Tracers: Food Webs of Fishes-Prince William Sound, Alaska (1994-1998)

couture.9.15 (<http://evos.nceas.ucsb.edu/evos/metacat/couture.9.15/default>)

- Facility to issue DOI identifiers for publicly archived data sets:

doi:10.xxxx/AA/gleon.9.15

- Always resolve to the data set
- Used in journals to cite data usage



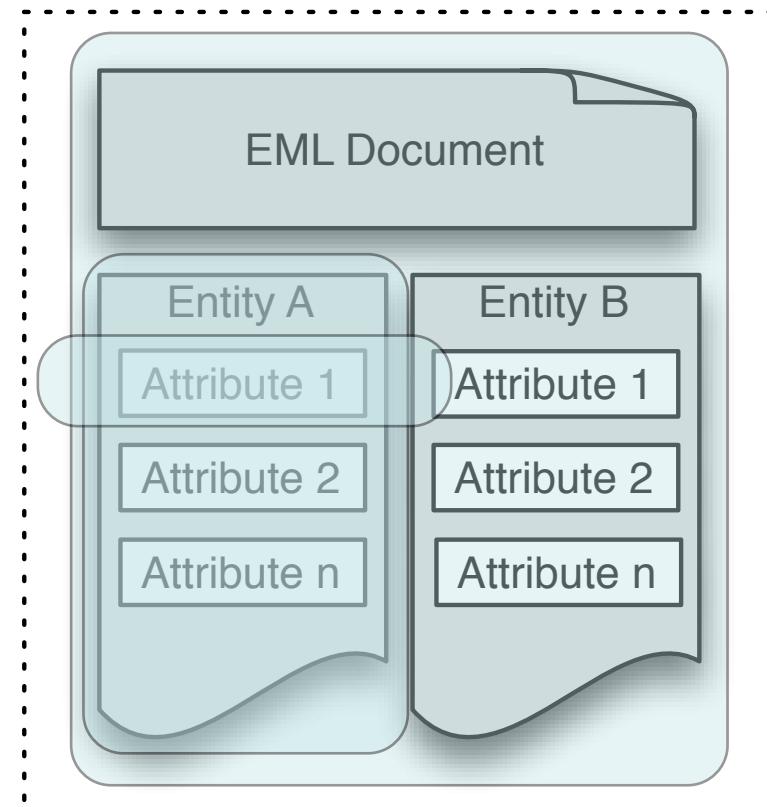
ADOPT STANDARDS

Metadata and data heterogeneity

- Every community has
 - many data schemas
 - one for each project and person
 - many data formats
 - ASCII, NetCDF, HDF, GeoTiff, ...
 - many metadata schemas
 - Biological Data Profile, Darwin Core, Dublin Core, **Ecological Metadata Language (EML)**, Open GIS schemas, ISO Schemas, ...
- Accepting this heterogeneity is critical

EML Data Package

- Data Package (eml-dataset)
 - collection of data Entities
- Entity (eml-dataTable)
 - tabular data
 - other
- Attribute
 - data column



Metadata



knB ABOUT DATA SHARE TOOLS Search for data SIGN IN

< Back to search

Thissell A and Balch J. Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013) (doi:10.5063/F1XW4GQH)

Files	Package	urn:uuid:19d860f4-7291-4cd0-a55c-5aba5f511662 .zip	14.79 MB	4 downloads
	Metadata	doi:10.5063/F1XW4GQH.xml (EML)	11.28 KB	79 views
	Data	jkbalch.14.1	text/csv	14.78 MB 6 downloads Details

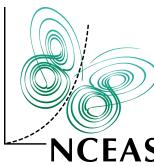
[Download all](#)

General

Title: Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013)

Identifier: jkbalch.15.2

Abstract: Intensifying land use practices and global climate change are leading to increased forest fragmentation and more severe droughts in the southeastern transitional forests of the Brazilian Amazon. The combination of land use and droughts has increased fire frequency, intensity, and extent in the region. Various plant traits can either limit or promote flammability, as demonstrated in many seasonally dry environments (e.g. Mediterranean shrublands), which are home to species that are adapted to frequent fire. However, little research has been done to associate flammability with the leaf traits of humid tropical forest species that are experiencing fire with greater frequency and severity than historically recorded. In this study, the leaf traits and consequent burning characteristics of seventeen species abundant in a transitional forest plot in Mato Grosso, Brazil were analyzed through controlled combustion experiments and trait measurements. Principal components analysis and linear regression were employed to compare leaf traits and their ability to predict flammability. The results reported herein show that these species demonstrate variability in flammability. Individual leaf surface area and leaf volume were found to be the most significant predictors of flammability among these species. The most flammable species were those with thin, lightweight leaves that become arranged in loosely packed, highly aerated fuelbeds and would be able to drive rapidly moving, high intensity fire across the forest floor, consuming a larger amount of fuel. The less flammable species were characterized by thick, heavy leaves that create dense fuelbeds, potentially inhibiting fire spread and the overall consumption of available fuels. Analyzing the dynamics between leaf traits and plant flammability with forest composition, successional trajectories and future fire behavior could lead to an examination of overall landscape effects in which certain species with fire-adaptive strategies could outcompete more fire-sensitive species, therefore influencing future fire regimes. Furthermore, given a scenario of increasing fire, this knowledge can lead to more effective conservation and land use policies in tropical regions.



Owner and Contact Metadata

People and Associated Parties

Data Set Creators

Individual **Amy Thissell**

Organization **The Pennsylvania State University**

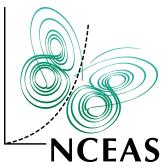
Position Primary Data Set Owner

Individual **Dr. Jennifer Balch**

Organization **University of Colorado**

Position Co-Owner

Address Guggenheim Building, 260 UBC, Room 110, University of Colorado,
Boulder, Colorado 80309 USA



Data file metadata

Data Table, Image, and Other Data Details

Technical Metadata

[Ecological Metadata Language \(EML\) File](#)

Data Table

Entity Name **alltemps_Thissell_2013.csv**

Object Name **alltemps_Thissell_2013.csv**

Online Distribution Info **jkbalch.14.1**

Size **15500059 byte**

Text Format

Number of Header Lines

1

Record Delimiter

#xA

Attribute Orientation

column

Simple Text

Field Delimiter

,

Number Of Records

214569

Column metadata

Attribute Information

Variables

scode

samp

sec

p10

p20

p40

t_avg

Name

scode

Label**Definition**

species code

Storage Type**Measurement Type**

nominal

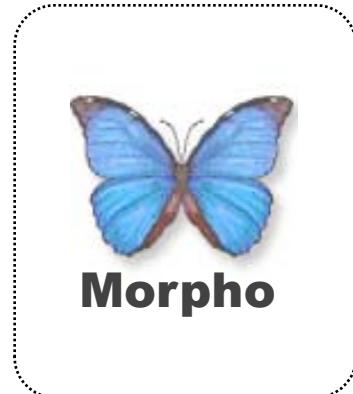
Measurement Domain

Definition First three letters of genus and species name

Missing Value Code

Morpho highlights

- Create metadata in EML format
- Manage data in EML packages
- Save, publish, and share data
- Search for data
- Multi-language
 - English, Spanish, Chinese, French, **Portuguese**, Japanese
- Export data and metadata
- Cross-platform, and open source



Wizards to create metadata

New Data Package Wizard

Welcome to the Data Package Wizard

This wizard guides you through the process of creating a new data package. If you have any questions or need help, please refer to the documentation. You can also contact the Data Package Manager for assistance.

If you have any questions or need help, please refer to the documentation. You can also contact the Data Package Manager for assistance.

Enter an abstract description of the data package. This will be used to describe the data itself.

Title: **New Data Package Wizard**

People or Organizations Associated With This Data Package

Owners

Enter info about the organization or person associated with this data package.

Description: Enter a description of the geographic coverage. Enter a general description of the geographic area in which the data were collected. This can be a bounding box or a polygon.

One or more organizations can be associated with this data package.

Abstract:

Step 2 of 15

Coverage

Method

Access

Step 5 of 15

Note: Required information is highlighted in yellow. It is highly recommended that you provide this information.

Named Regions:

Step 1 of 15

Define Temporal Coverage:

Choose date type:

Bounding Box:

Enter start date:

Enter end date:

Enter start time:

Enter end time:

Enter start date:

Enter end date:

Enter start time:

Enter end time:

Define Access:

Select a user or group from the list below:

Name	Email / Description / Distinguished Name
Access Tree	
SDSC	
OBFS	
UCNRS	
A Tester	atester@ucnrs.org, rnotrott@nceas.ucsb.edu
Alexander Glazer	aglazer@ucnrs.org, alexander.glazer@ucop.edu
Alicia Flammia	aflammia@ucnrs.org, alicia_flammia@hotmail.com
Allan Muth	amuth@ucnrs.org, deepcanyon@mindspring.com
Andrew Brooks	abrooks@ucnrs.org, brooks@lifesci.ucsb.edu
Arnulfo Lozoya	alozoaya@ucnrs.org, lozoya@gte.net

Refresh the user list...

Allow selected user(s) Read access

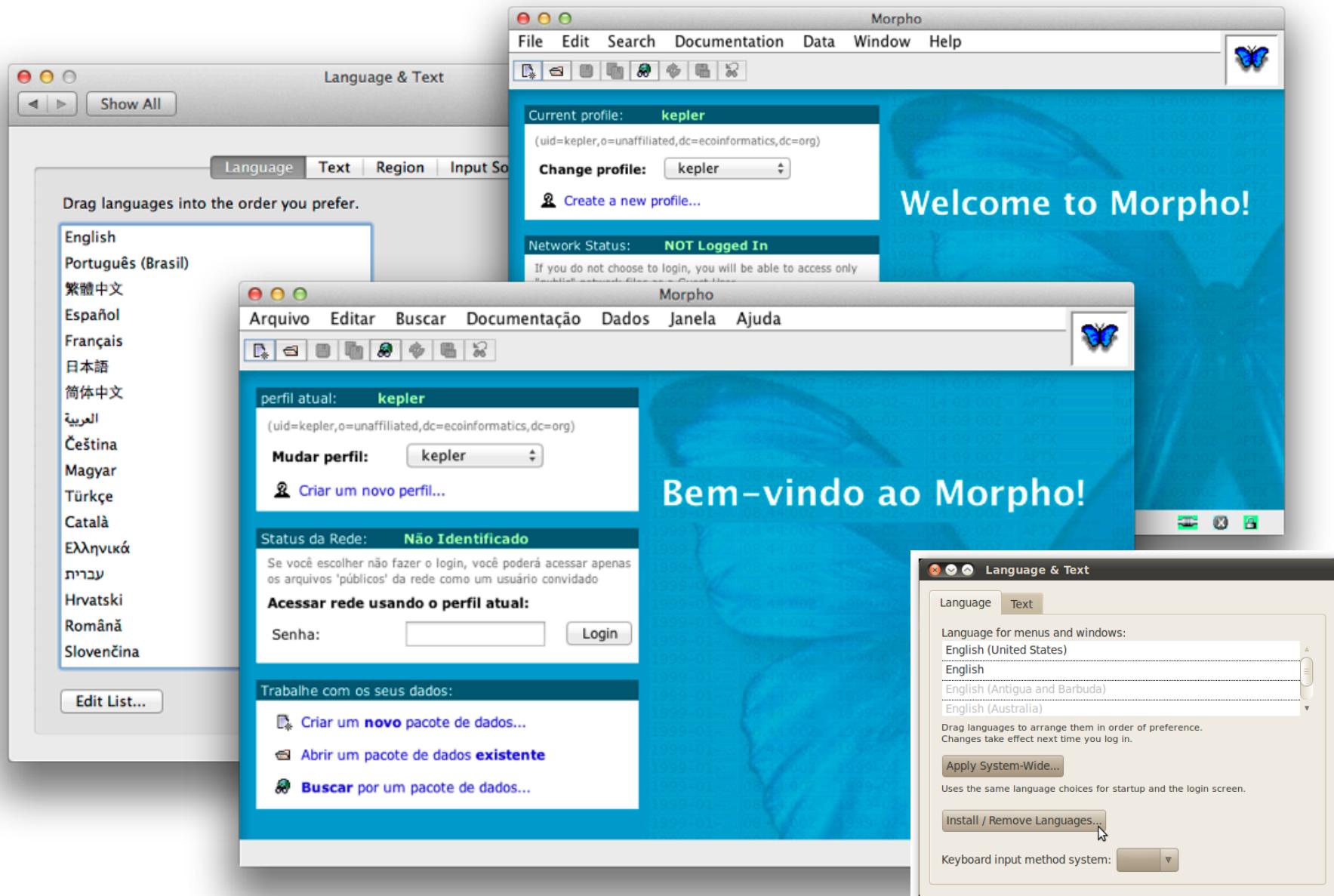
Description of access levels:

- Read: Able to view data package.
- Read & Write: Able to view and modify data package.
- Read, Write & Change Permissions: Able to view and modify datapackage, and modify access permissions.
- All: Able to do everything (this is the same as Read, Write & Change Permissions)

OK Cancel



Multilingual support



The image displays three screenshots of the Morpho software interface, illustrating its multilingual support across different platforms and components.

- Top Screenshot:** A Mac OS X window titled "Morpho". The menu bar includes "File", "Edit", "Search", "Documentation", "Data", "Window", and "Help". The title bar shows "Morpho". The main area displays a "Welcome to Morpho!" message with a blue butterfly icon. A sidebar on the left lists languages: English, Português (Brasil), 繁體中文, Español, Français, 日本語, 简体中文, العربية, Českina, Magyar, Türkçe, Català, Ελληνικά, עברית, Hrvatski, Română, and Slovenčina. A "Change profile" dropdown is set to "kepler".
- Middle Screenshot:** A Windows-style window titled "Morpho". The menu bar includes "Arquivo", "Editar", "Buscar", "Documentação", "Dados", "Janela", and "Ajuda". The title bar shows "Morpho". The main area displays a "Bem-vindo ao Morpho!" message with a blue butterfly icon. A sidebar on the left lists languages: perfil atual: kepler, Mudar perfil: kepler, Criar um novo perfil..., Status da Rede: Não Identificado, Acessar rede usando o perfil atual:, Senha: [input field], and Login. A note states: "Se você escolher não fazer o login, você poderá acessar apenas os arquivos 'públicos' da rede como um usuário convidado".
- Bottom Screenshot:** A Mac OS X window titled "Language & Text". The menu bar includes "File", "Edit", "Search", "Documentation", "Data", "Window", and "Help". The title bar shows "Language & Text". The main area displays language selection options: Language for menus and windows: English (United States) and English (Antigua and Barbuda). A note says: "Drag languages to arrange them in order of preference. Changes take effect next time you log in." Buttons include "Apply System-Wide...", "Install / Remove Languages...", and "Keyboard input method system: [dropdown menu]".



CREATE NETWORKS

Global Metacat deployments



TFRI Metacat Data Catalog

metacat.tfri.gov.tw/tfri/

Taiwan Forestry Research Institute Data Catalog

林業試驗所研究資料目錄

Home Analysis Tools EML Parser

Metacat ver 2.4.1 powered by KNB, localized by TFRI

You are NOT logged in ([Login](#)). You may see some features that are not available to you, but will have access only to "public" data (see below).

Enter a search phrase (e.g. biodiversity) to see what data is available. You can also browse by category using the links below.

Search Title, Abstract, Keywords, Person
 Search all fields (Slower)

中文

植物, 森林, 生態, 育林, 人工林, 保護, 經營, 經濟, 氣候, 環境, 土壤, 地理, 文化, 化學, 木材纖維, 利用, 竹材, 碳吸存

search for data

Repositório de Dados do PPBio

ppbio.inpa.gov.br/knb/style/skins/ppbio/

PPBio Programa de Pesquisa em Biodiversidade

Página do PPBio Repositório de Dados

Repositório de Dados do PPBio

Bem vindo ao Repositório de Dados do Programa de Pesquisa em Biodiversidade - PPBio. Esse repositório contém dados de levantamentos realizados no âmbito do Programa de Pesquisa em Biodiversidade e projetos parceiros. Os dados armazenados nesse repositório em breve estarão conectados à rede KNB ("Knowledge Network for Biocomplexity"), um repositório internacional de dados. Algumas das séries de dados contidas aqui foram geradas por pesquisadores em esforços individuais, enquanto outras são resultado de esforços conjuntos de grupos de contribuintes. As descrições de cada conjunto de dados contêm mais informações sobre pessoas e instituições envolvidas. Os dados seguem a política de dados do PPBio, disponível em <http://ppbio.inpa.gov.br/politicadados>.

Para perguntas, comentários e sugestões, por favor contate: ppbio@inpa.gov.br.

Busca por Dados

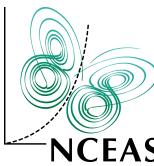
Buscar

Buscar somente nos campos "Título", "Resumo", "Palavras-Chave", "Pessoas Envolvidas"
 Buscar todos os campos

Esta ferramenta permite a busca por conjuntos de dados de interesse. Ao inserir um texto no quadrado e clicar no botão "Buscar", a busca será conduzida apenas nos campos "Título", "Resumo", "Palavras-Chave" e "Pessoas Envolvidas". Ao optar pela opção "Buscar todos os campos", a busca ocorrerá em todos os campos (isso fará com que a busca leve mais tempo).

O caractere "%" pode ser usado como um "coringa" (ou "wildcard") nas buscas (por exemplo, "%biodiversidade%" localizará qualquer frase que contenha a palavra biodiversidade).

Buscar todos os metadados cadastrados no Repositório do PPBio

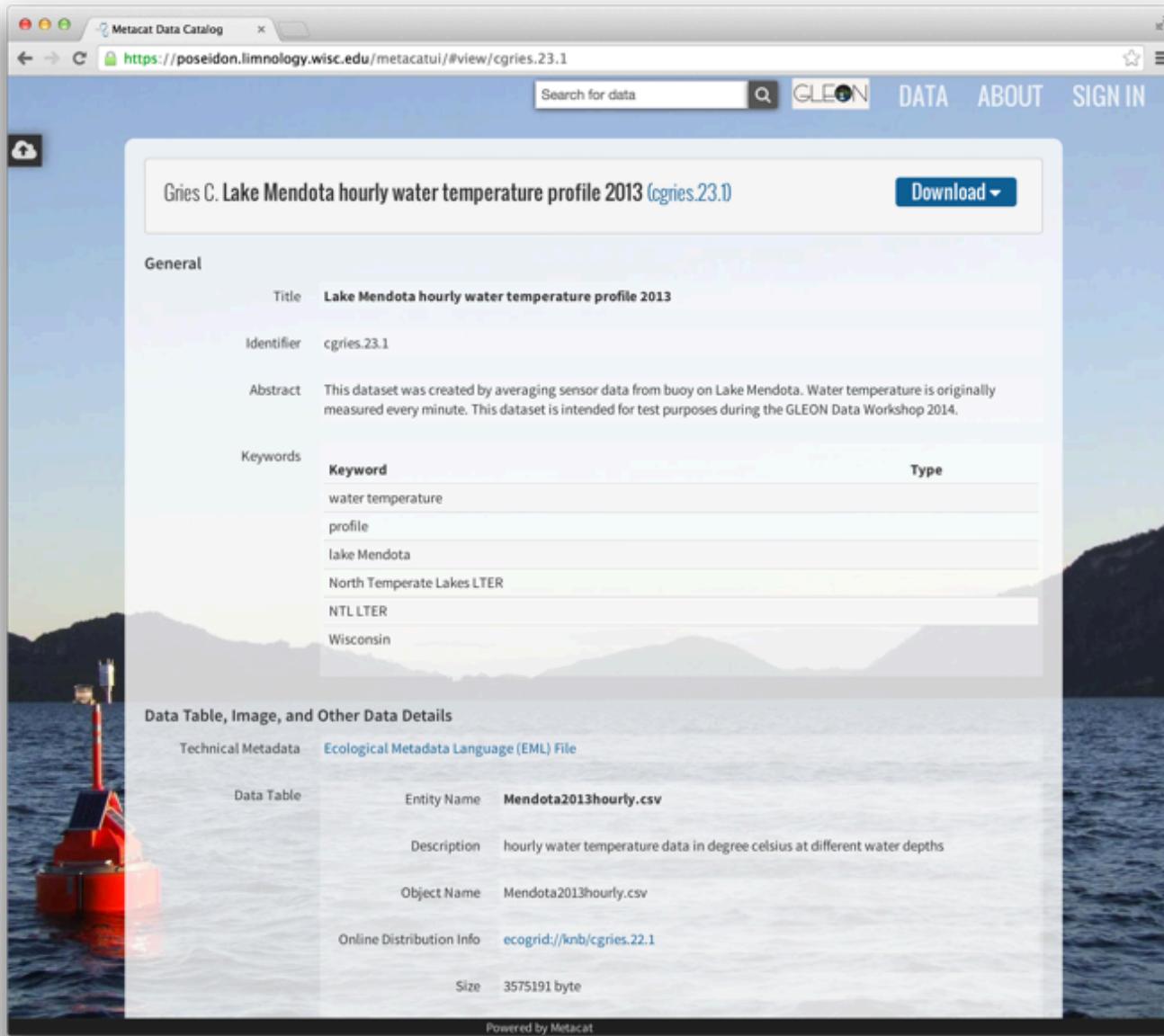


UI Themes

The collage illustrates four distinct UI themes for scientific data repositories:

- Top Left:** KNB | Knowledge Network. A dark-themed interface with a red header bar containing the KNB logo, followed by "ABOUT", "DATA", "SHARE", and "TUTORIALS". Below is a large image of a red flower.
- Top Middle:** Metacat Data Catalog. A light blue-themed interface with the NCEAS logo and "NCEAS News" button. It features a search bar and a sidebar with links like "DATA", "ABOUT", and "SIGN IN".
- Bottom Left:** Gulf of Alaska Data Portal. A dark-themed interface with a "Gulf of Alaska Data Portal" header. It includes a "Share Your Data" button and various filtering options such as "FILTERS", "DATA ATTRIBUTES", and "ONLY RESULTS WITH DATA".
- Bottom Middle:** Metacat Data Catalog. A light blue-themed interface showing a map of the Americas with data density overlays. A "Filter" sidebar is open, displaying search fields for "Anything", "Data attribute (density)", "Creator", "Year", "Taxon", and "Location". To the right is a list of datasets, each with a thumbnail, title, views, and download link.

GLEON Data Catalog



The screenshot shows a web browser displaying the Metacat Data Catalog at <https://poseidon.limnology.wisc.edu/metacatui/#view/cgries.23.1>. The page title is "Gries C. Lake Mendota hourly water temperature profile 2013 (cgries.23.1)". The main content area is titled "General" and contains the following information:

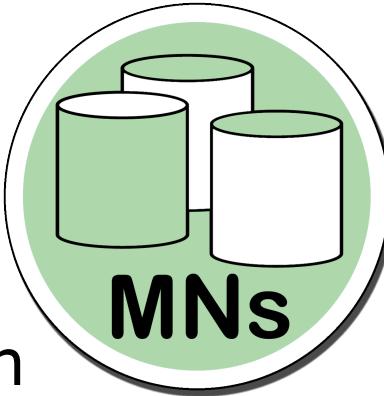
	Title	Identifier	Abstract	Keywords	Type
	Lake Mendota hourly water temperature profile 2013	cgries.23.1	This dataset was created by averaging sensor data from buoy on Lake Mendota. Water temperature is originally measured every minute. This dataset is intended for test purposes during the GLEON Data Workshop 2014.	water temperature profile lake Mendota North Temperate Lakes LTER NTL LTER Wisconsin	

Below this section is a "Data Table, Image, and Other Data Details" section. It includes a thumbnail image of a red buoy in the water. The table details are as follows:

Technical Metadata	Ecological Metadata Language (EML) File
Data Table	Entity Name: Mendota2013hourly.csv
	Description: hourly water temperature data in degree celsius at different water depths
	Object Name: Mendota2013hourly.csv
	Online Distribution Info: ecogrid://knb/cgries.22.1
	Size: 3575191 byte

At the bottom of the page, it says "Powered by Metacat".

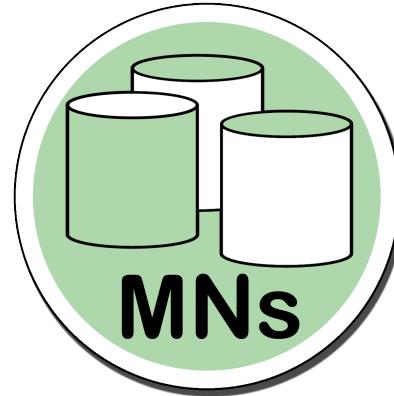
Member Nodes



- **Authoritative** members of the Federation
- **Curate** data holdings
 - *Provide unique identifiers for each object*
 - *Ensure availability, quality, and reliability*
- **Replicate** holdings for other MNs
- Provide access and **access control**
- **Log** and report accesses to objects
- Engage with DataONE community
- Deploy DataONE-compatible software systems



Member Nodes



Avian
Knowledge
Network

...and many more!



A DataONE Search Tool for Scientific Data

Search For:

Hint: boolean operators and phrases are allowed. ex: precipitation or (rain and "moisture content")

Results/Page

10

[SEARCH](#)[Show/Hide Advanced Options](#)[Help](#)**Fielded Search**

FullText	OR
FullText	OR
FullText	

[Help](#) | [clear](#)**Date Search**

- Collection Date
 Publication Date
 Either

during

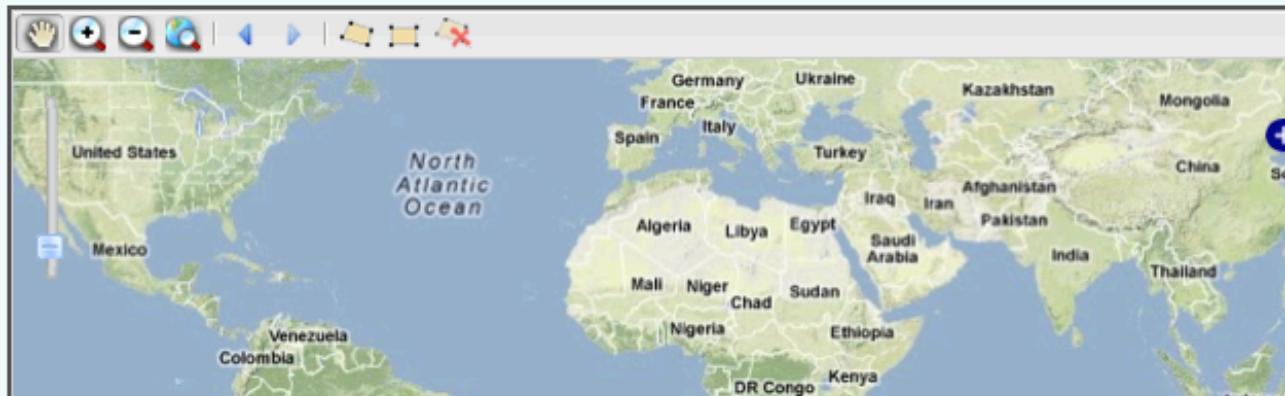


thru



mm/dd/yyyy

mm/dd/yyyy

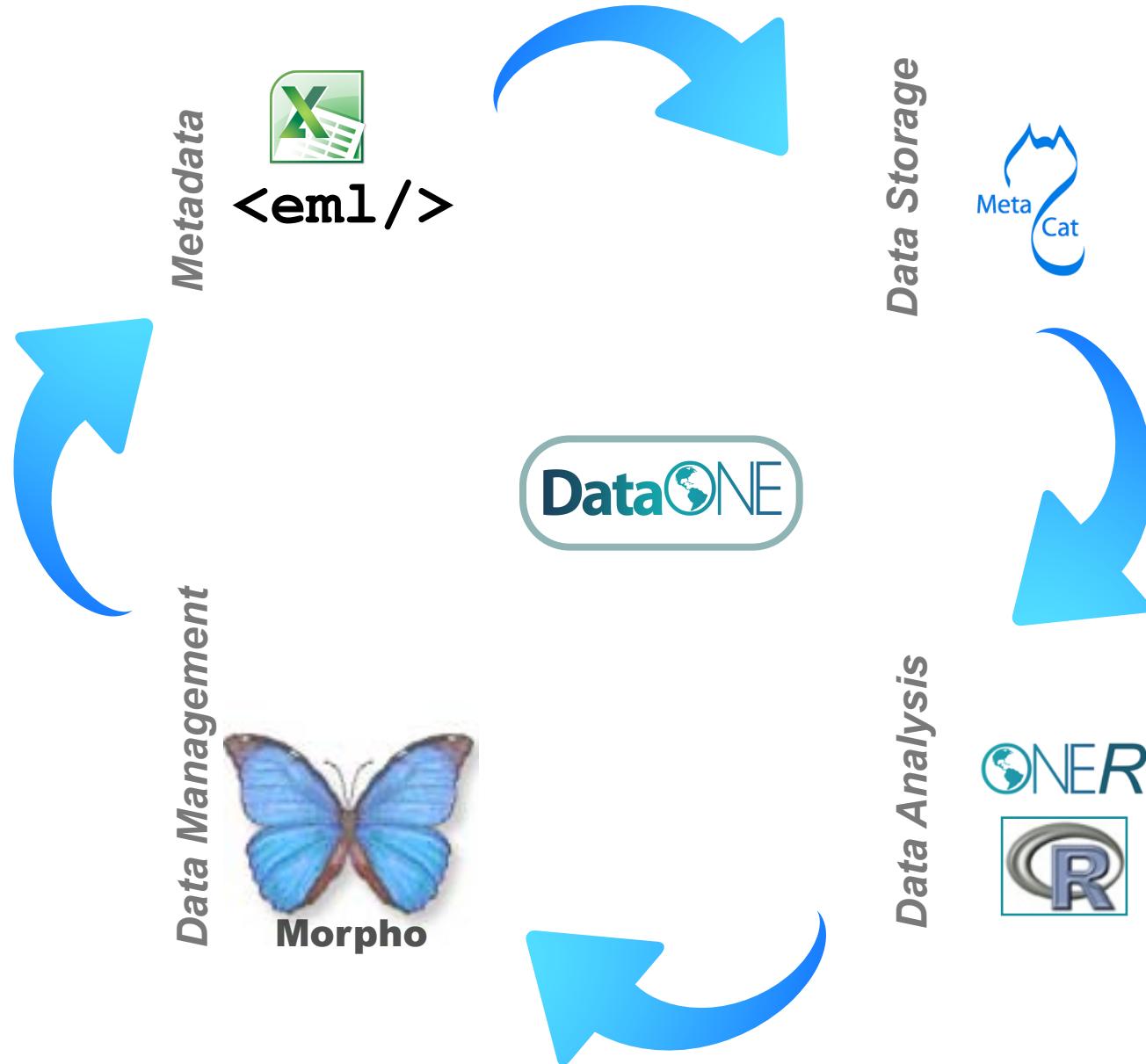
[Help](#) | [clear](#)**Geographic Search****List Areas in:** USA WORLD[Select from list](#)**Search Area:** overlaps encloses

North



CREATE INTEROPERABLE SOFTWARE

Data life-cycle



KNB System Components

Document

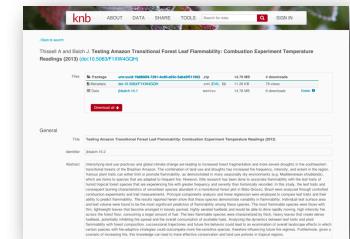
Share

Analyze

Communicate



W
E
B

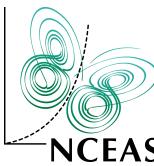


Metadata

Data

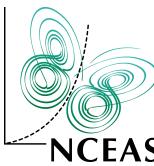
Workflows

Results



Data Registration Activity

- Create an account
 - <https://poseidon.limnology.wisc.edu/metacatui/#signup>
 - or
 - <https://identity.nceas.ucsb.edu/>
- Register a dataset!
 - <https://poseidon.limnology.wisc.edu/metacatui/#share>



Questions?

- Contact:
 - Ben Leinfelder <leinfelder@nceas.ucsb.edu>
 - Matt Jones <jones@nceas.ucsb.edu>
- Links
 - <http://www.nceas.ucsb.edu/ecoinfo/>
 - <https://knb.ecoinformatics.org/>
 - <http://www.dataone.org>