

# **Data Management for Synthesis**

Ben Leinfelder  
Matthew B. Jones

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



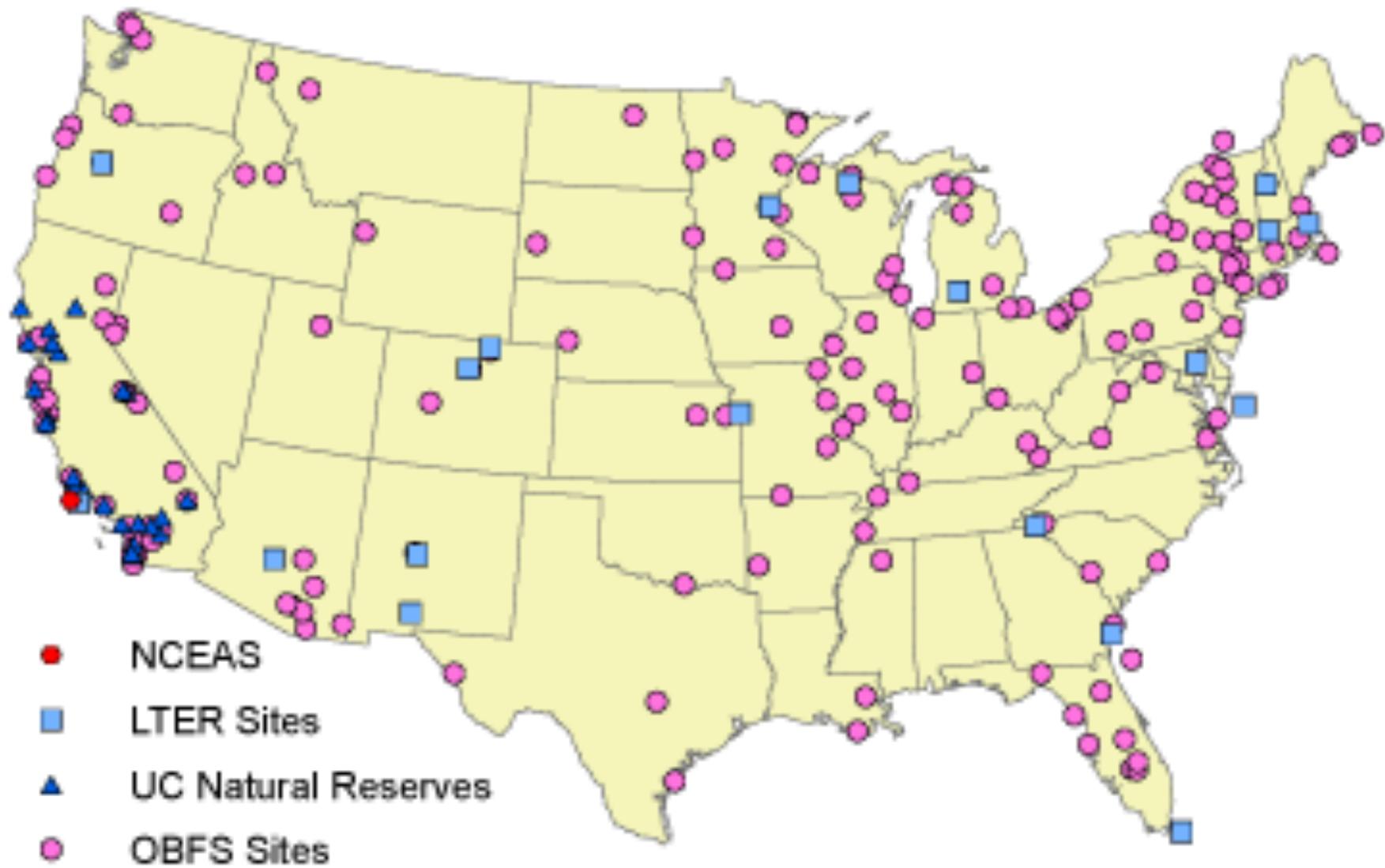
LNCC GBIF workshop  
August 25, 2014

# Barriers to Synthesis

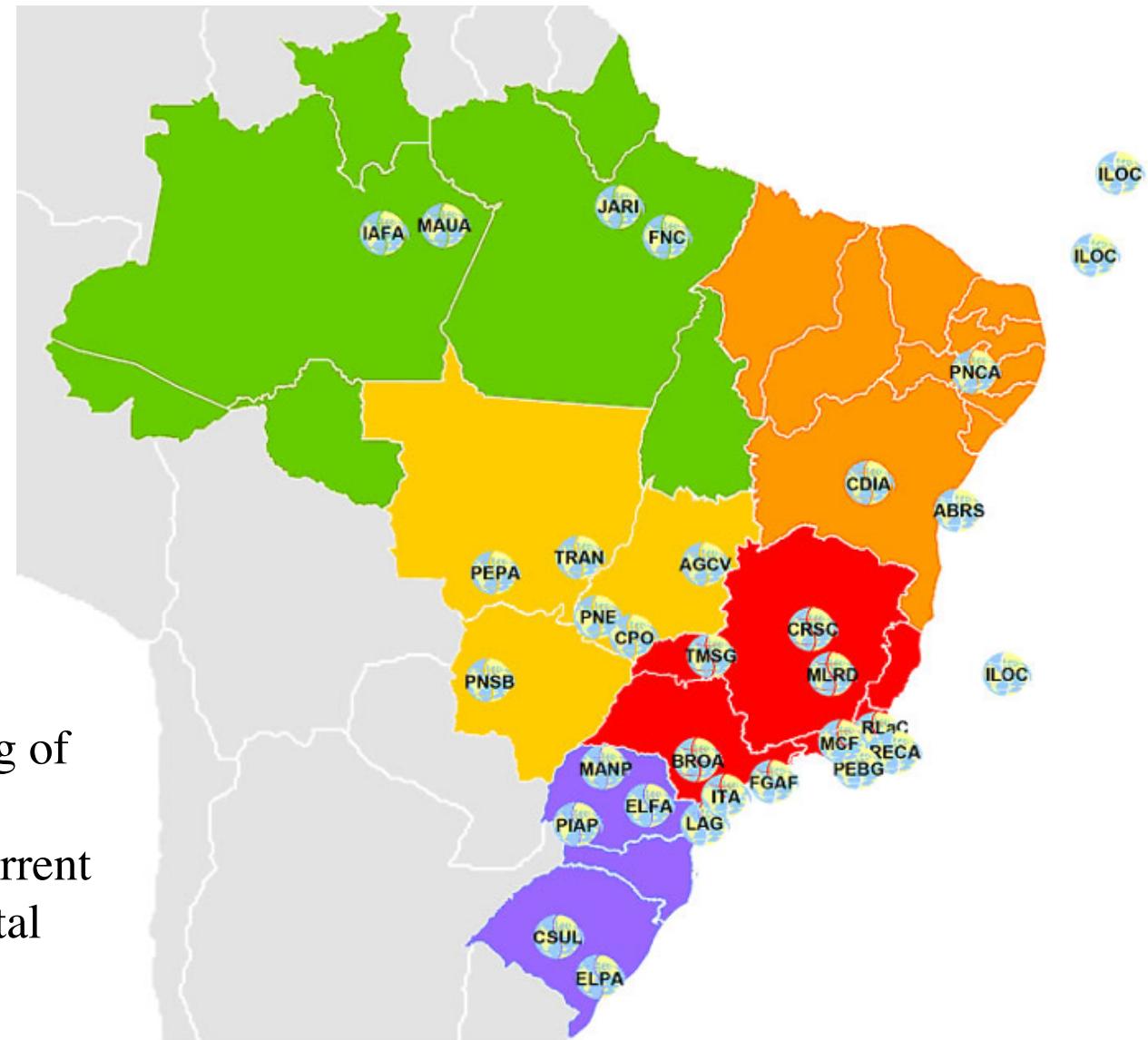
---

- 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



# PELD Sites

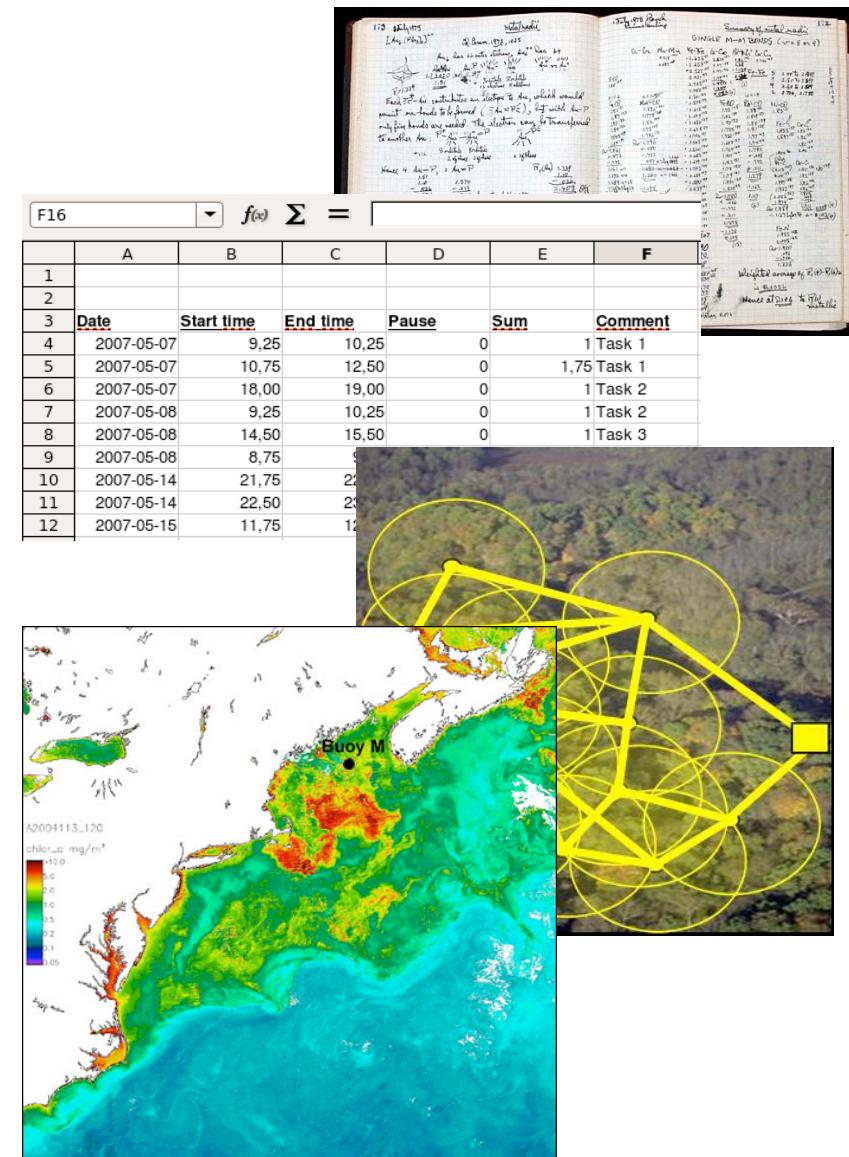


ILTER:

"increase understanding of  
global ecosystems and  
propose solutions to current  
and future environmental  
problems"

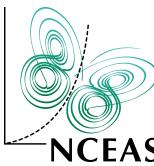
# What data are in scope?

- 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

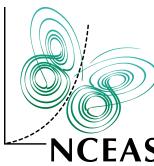
### (a) Georgia Coastal Ecosystems LTER (EML)

**(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

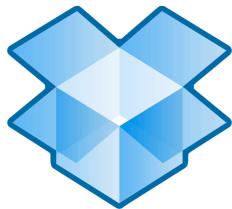
**(d) Macroecological data for fossil occurrences (Paleobiology Database)**



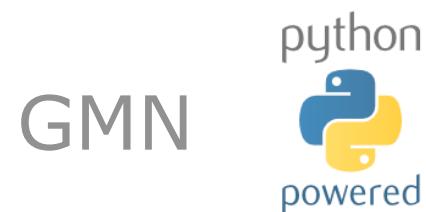
# Synthesis requires sharing data

- 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



Software diversity





find plots containing

 download 0 items

[advanced search](#) | [browse data](#)

[HOME](#)

[FAQ](#)

[SUBMIT DATA](#)

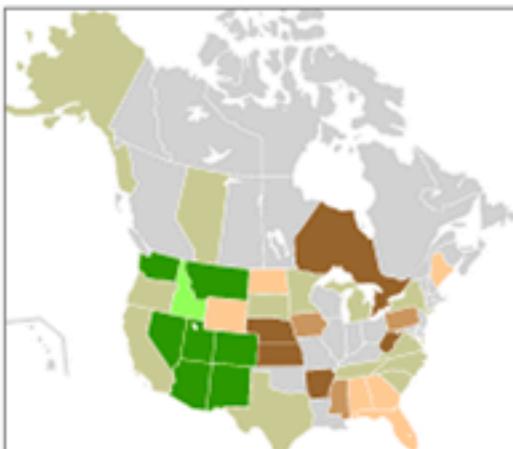
[ABOUT](#)

[MY ACCOUNT](#)

[SITE MAP](#)

## Find Plots

[Browse plots](#)  
[Simple search](#)  
[Search with a map](#)  
[Advanced plot search](#)



Map Key: plots [Larger Map](#)

1-49	50-99	100-249
250-999	1,000-3,000	> 3,000

## Recently Added Plots

Project ( <a href="#">view all</a> )	Added
EMNE - NatureServe	11-Sep-

## Plant Taxa

[What is a plant concept?](#)  
[Browse plants](#)  
[Search plants](#)  
[Submit plants](#)

## Plant Communities

[What is a community?](#)  
[Search communities](#)  
[Submit communities](#)

## Supplemental Data

[People](#)  
[Stratum methods](#)  
[Cover methods](#)  
[Projects](#)  
[References](#)  
[Search supplemental data](#)

## Data in VegBank

Plots	72,857
--Classified Plots	60,213
----to NVC communities	5,478

## News

» [Map plots: Example](#) | [Datacart](#) | [Multiple Datasets](#)  
(Requires Login)  
» [Save Your Datacart](#) | [Edit Datasets](#)  
» Create a [Constancy Table](#)

## My VegBank Account

[Edit profile information](#)  
[Manage datasets](#)

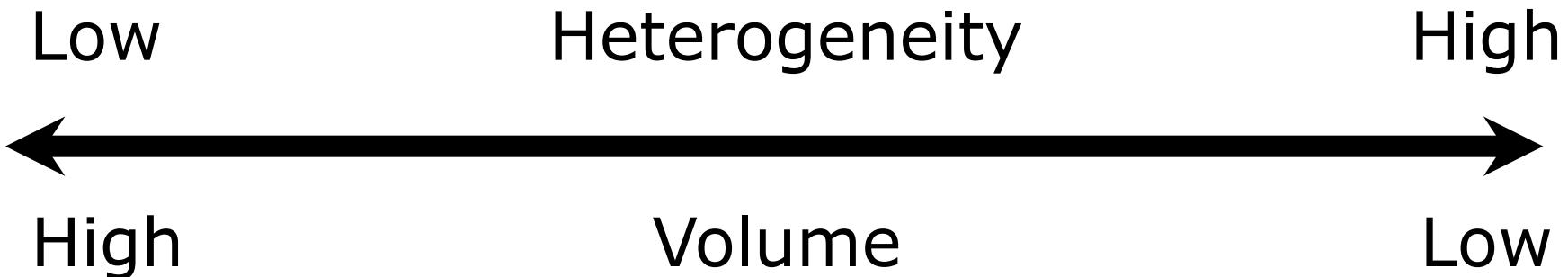
## Learn About VegBank

[What is VegBank?](#)  
[What is a plot?](#)  
[FAQ](#)  
[Tutorial](#)  
[Cite or link to VegBank](#)  
[Terms of use](#)  
[Site map](#)  
[Contact](#)

## Contribute Plot Data



# Data Heterogeneity



- Tight coupling
  - Simple subsetting
  - Explicit semantics

- Loose coupling
  - Hard subsetting
  - Limited semantics





# Solutions

- Preserve data
- Adopt standards

<EML/>



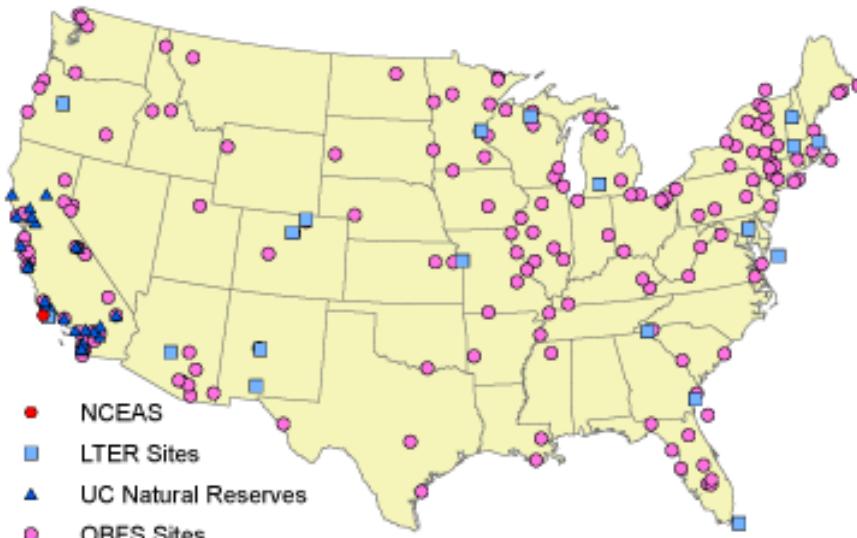
- Create networks
- Create interoperable software





# PRESERVE DATA

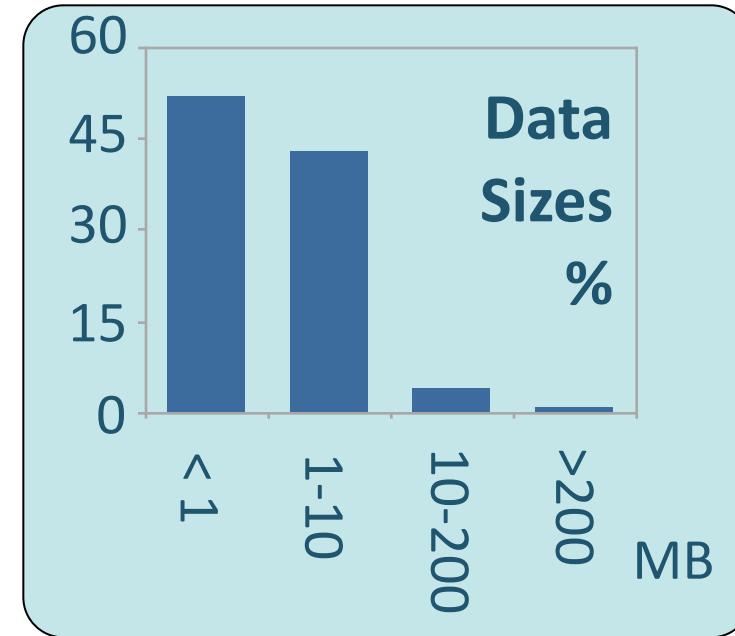
# Preserve data in the KNB



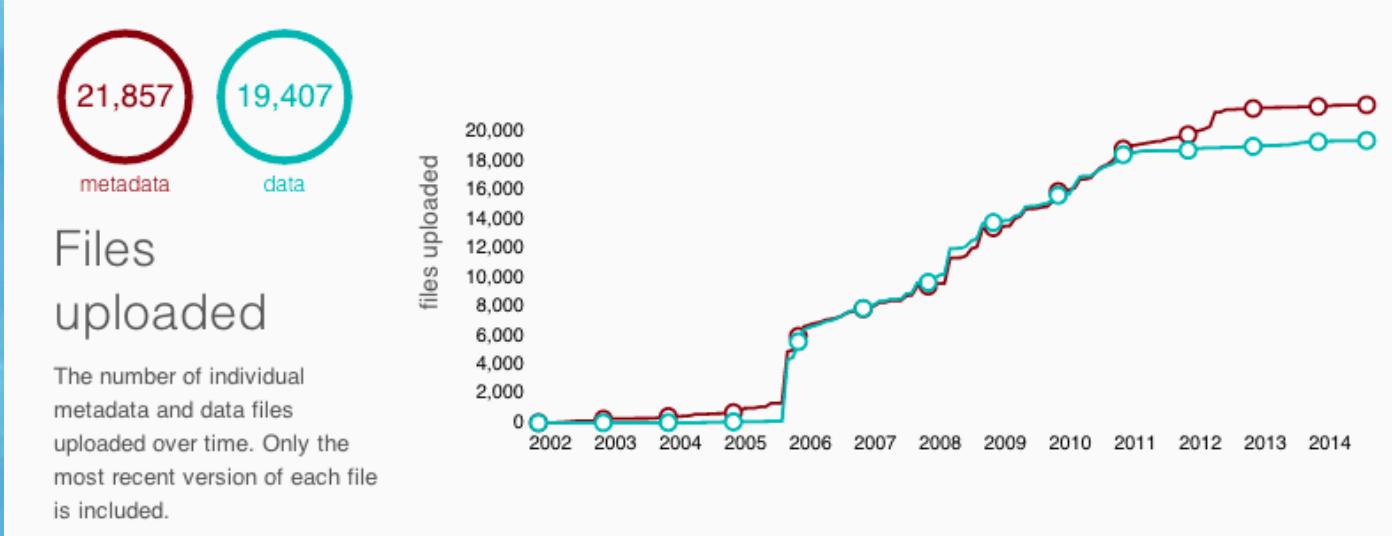
- Diverse Contributors
- Individual investigators
- Field stations and networks
- Government agencies
- Non-profit partnerships
- Scientific Societies
- Synthesis centers

## Data Types

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



# Knowledge Network for Biocomplexity 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 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 "Only results with all spatial coverage inside the map" (with a checkbox and a help icon). To the right of the map is a list titled "Mapping 1 to 112 of 112 data packages". The list is sorted by "Most recent" and includes the following entries:

- 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

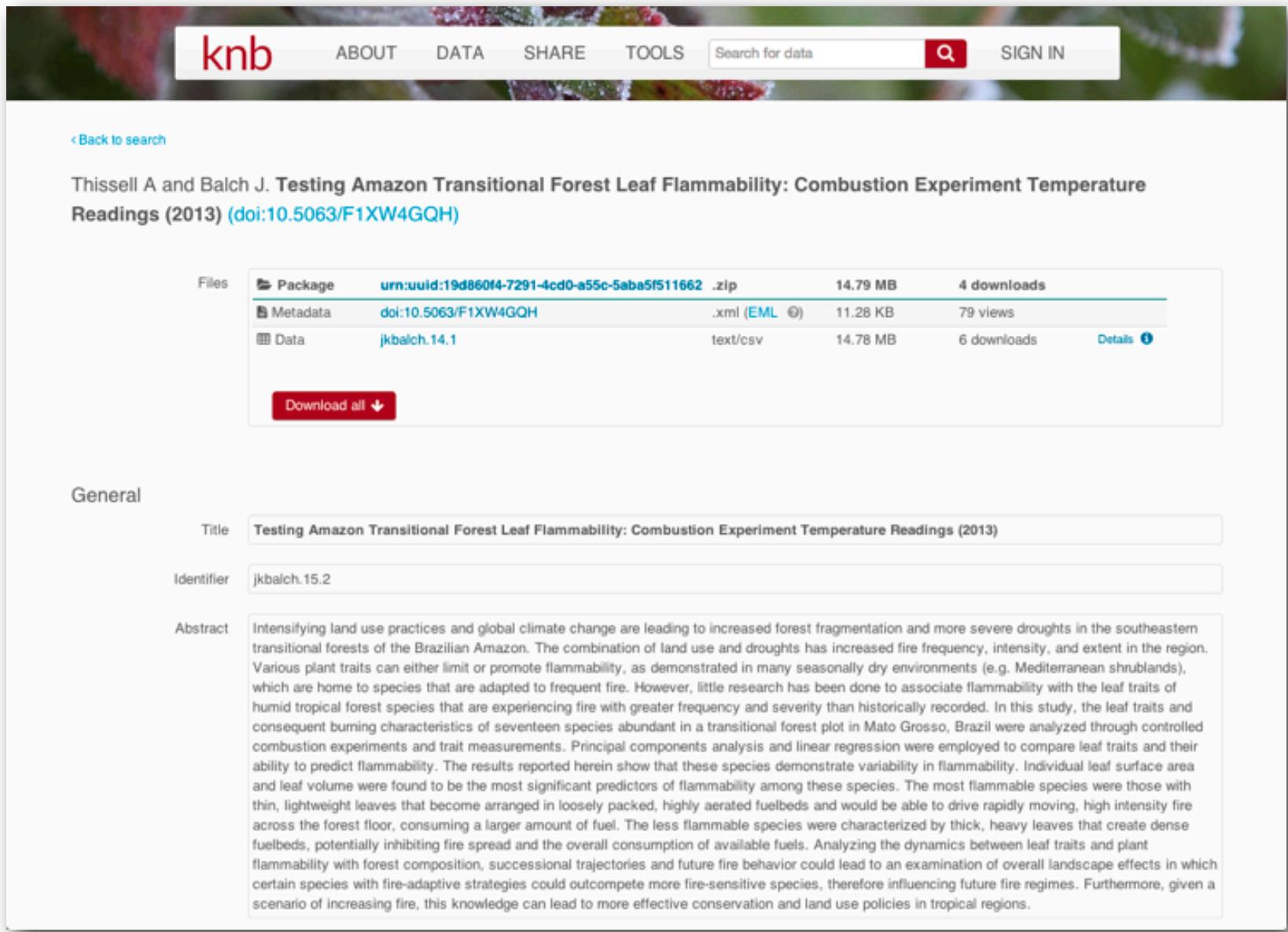


# ADOPT STANDARDS

# 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

# Metadata



The screenshot shows the KNB (Knowledge Network for Biocomplexity) metadata interface. At the top, there is a navigation bar with links for ABOUT, DATA, SHARE, TOOLS, a search bar labeled "Search for data" with a magnifying glass icon, and a "SIGN IN" button. Below the navigation bar, the title of the study is displayed: "Thissell A and Balch J. Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013) (doi:10.5063/F1XW4GQH)".

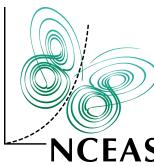
The main content area displays the study's files. There are three items listed:

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

A red "Download all" button is located at the bottom of this section.

Below this, under the "General" heading, there are three input fields:

- Title: Testing Amazon Transitional Forest Leaf Flammability: Combustion Experiment Temperature Readings (2013)
- Identifier: jkbalch.15.2
- Abstract: A detailed paragraph describing the study's purpose, methods, and findings related to forest leaf flammability in the Brazilian Amazon.



# Owner and Contact Metadata

## People and Associated Parties

### Data Set Creators

Individual **Amy Thissell**

Organization **The Pennsylvania State University**

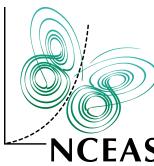
Position Primary Data Set Owner

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**

# 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:

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

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:

- 6
- 13
- 20
- 27

People or Organizations Associated With This Data Package

Owners

Enter info about the people or organizations associated with this data package.

Description: 

One or more owners

Define Temporal Coverage:

Choose date type:

Enter start date: February 2013

Bounding Box:

Set the geographic bounding box containing the fractional degrees values.

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	alozoya@ucnrs.org, lozoya@gte.net

Refresh the user list...

Allow selected user(s) Read access

Description of access levels:

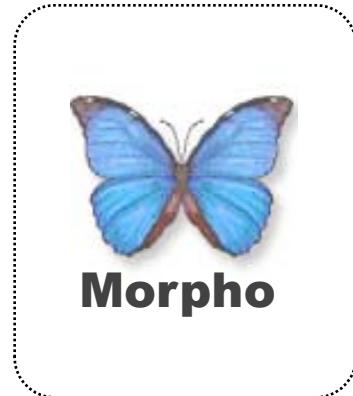
- Read: Able to view data package.
- Read & Write: Able to view and modify data package.
- Read, Write & Change Permissions: Able to view and modify datapackage, and modify access permissions.
- All: Able to do everything (this is the same as Read, Write & Change Permissions)

OK Cancel



# 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





# 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>)

- NCEAS can issue DOI identifiers for publicly archived data sets:

doi:10.xxxx/AA/gulfwatch.9.15

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



# CREATE NETWORKS

# Global Metacat deployments





## The US Long Term Ecological Research Network

[LTER Home](#)  
[Login](#)  
[Search](#)  
[Browse](#)

### Welcome to the LTER Data Catalog

Data are one of the most valuable products of the LTER program. The LTER Network seeks to inform the LTER and broader scientific community by creating well designed and well documented databases and to provide fast, effective, and open access to LTER data via a network-wide information system designed to facilitate data exchange and integration. Currently, the LTER Data Catalog contains entries for over 6000 ecological datasets from 26 LTER Network research sites, and thousands of additional datasets from numerous other ecological field stations and research institutions.

#### Data Catalog

The LTER Data Catalog includes content from both LTER and non-LTER data sources including PISCO, KNB, etc. By default, search results display only LTER data sources. You may include non-LTER data sources in your search by selecting the check box below.

**NEW!** When you begin typing in the search box form below, an auto-completion dialog will suggest ranked terms that originate from key-words and titles within the Data Catalog.

Search Term:

[Advanced Search](#)

Include non-LTER data

[Search](#)

[Reset](#)

#### LTER Data Policies

The [LTER data policy](#) includes three specific sections designed to express shared network policies regarding the release of LTER data products, user registration for accessing data, and the licensing agreements specifying the conditions for data use.

#### Other Databases

Additional information is available through these value-added data products:

- [LTER/USFS Climate / Hydrology Data](#)
- [Annual Net Primary Productivity Data](#)



**PPBio**

Research Program in Biodiversity  
 Programa de Pesquisa em Biodiversidade  
 Programa de Investigación en Biodiversidad

[Home](#) [Repository](#) [Register](#) [Logout](#)

### Data Set Citation

Drucker D of BR-LTER. Abundância e Distribuição de Ervas Terrestres em Parcelas Ripárias na Reserva Ducke: Variação Lateral.

drucker.3.1 (<http://aranha.inpa.gov.br/knb/metacat/drucker.3.1/ppbio>).

*Metadata download:*

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

### Data Set Owner(s):

*Individual:* MSc. Debora Drucker

*Organization:* BR-LTER

*Email Address:* deboradrucker@gmail.com

### Abstract:

Os dados aqui disponibilizados são produto do trabalho realizado por Debora Drucker durante seu mestrado. O objetivo central foi investigar a abundância e distribuição espacial de ervas terrestres (apenas as espécies que germinam e passam todo o seu ciclo de vida no solo) em 20 Parcelas ripárias paralelas aos igarapés na Reserva Florestal Adolpho Ducke.

### Keywords:

- Ervas
- Parcelas Ripárias
- Reserva Ducke
- Floresta de Terra Firme

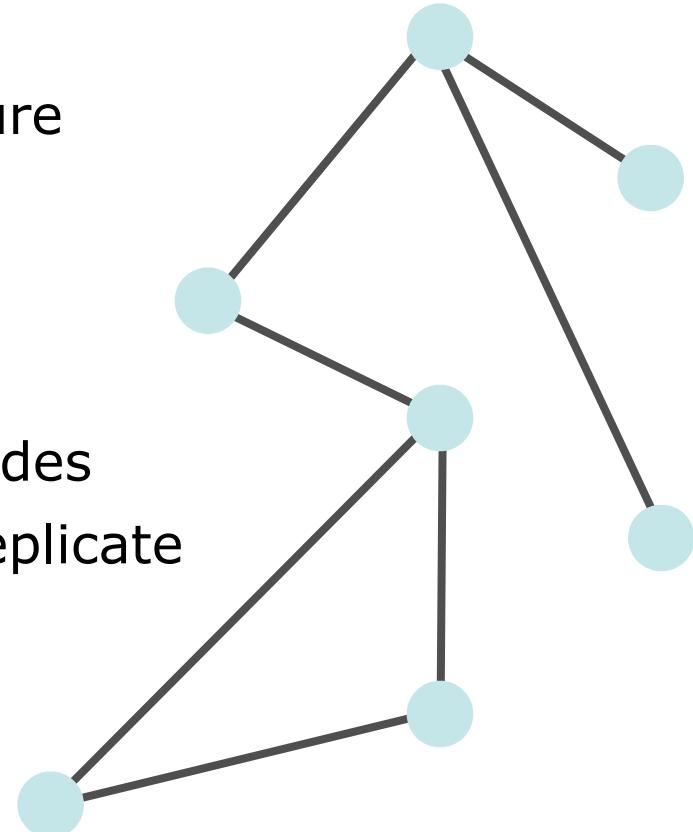
### License and Usage Rights:

- **Diverse Federation == Resilience**

- Failover for temporary outages
- Insurance against institutional failure
- Avoid data loss in the long-term

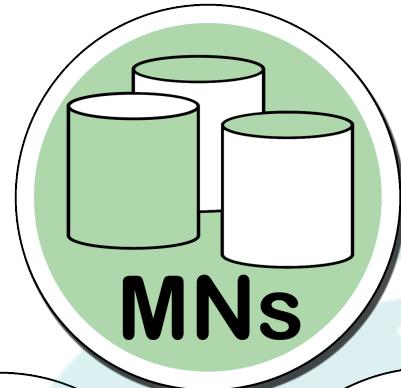
- **Diverse Federation == Scalability**

- Storage increases with Member Nodes
- Incremental costs to each MN to replicate
- Distributes sustainability costs



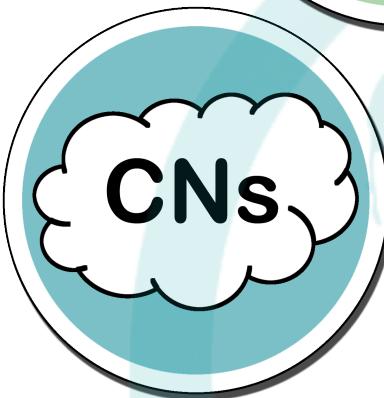
- **Member Nodes (MNs)**

- Heart of the federation
- Harness the power of local curation



- **Coordinating Nodes (CNs)**

- Services to link Member Nodes



- **Investigator Toolkit (ITK)**

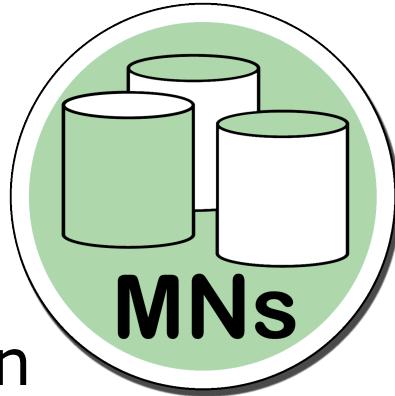
- Tools for the whole data lifecycle



Interoperability

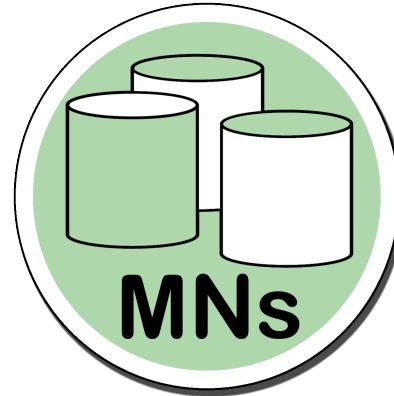
## Member Nodes

- **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



**knb**



**ONEShare**

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	

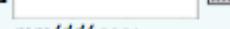
[Help](#) | [clear](#)**Date Search**

- Collection Date  
 Publication Date  
 Either

during

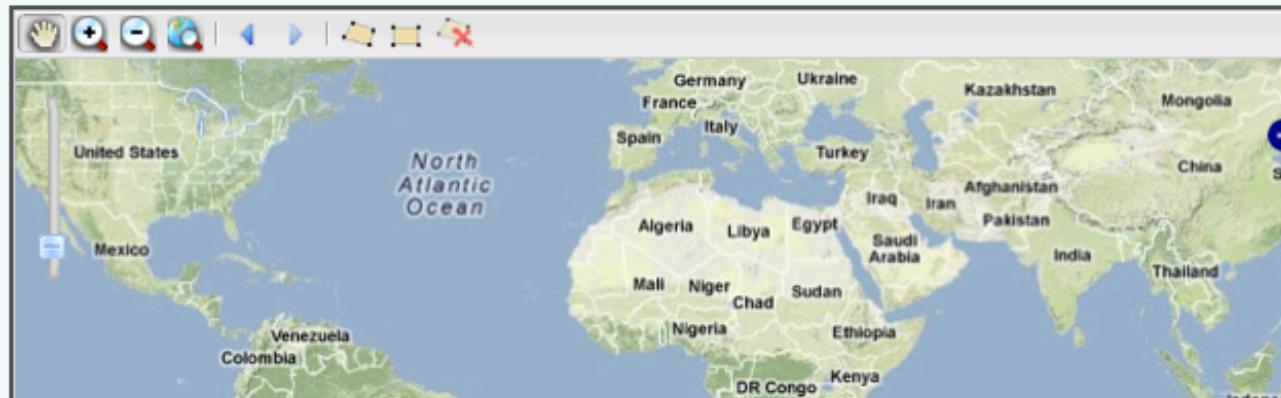


thru



mm/dd/yyyy

mm/dd/yyyy

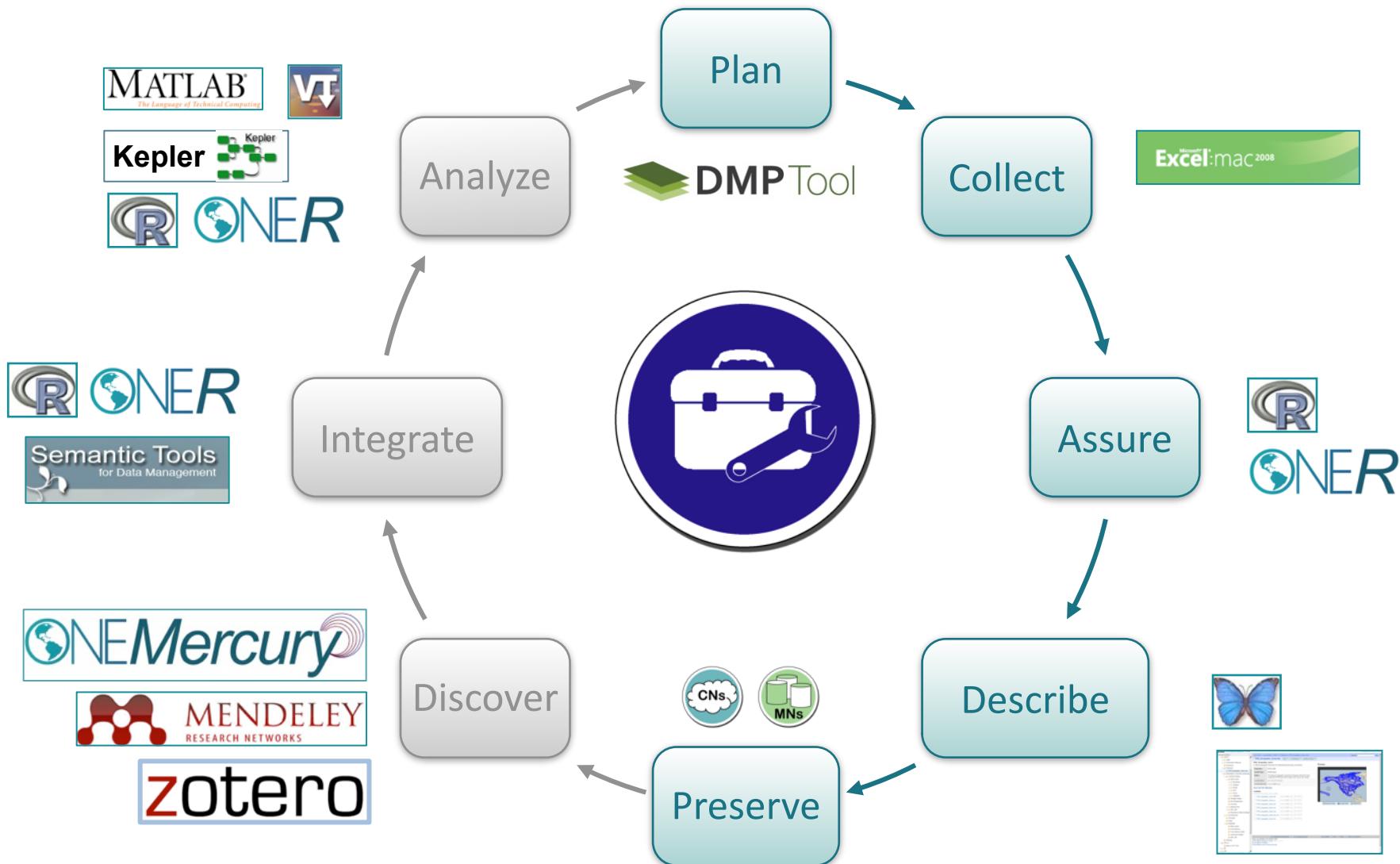
[Help](#) | [clear](#)**Geographic Search****List Areas in:** USA  WORLD[Select from list](#)**Search Area:** overlaps  encloses

North



# **CREATE INTEROPERABLE SOFTWARE**

# Software Interoperability



# KNB System Components

Document

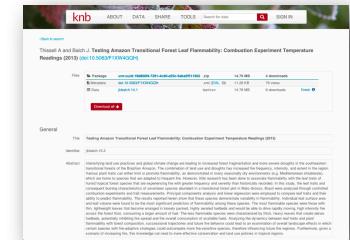
Share

Analyze

Communicate



W  
E  
B



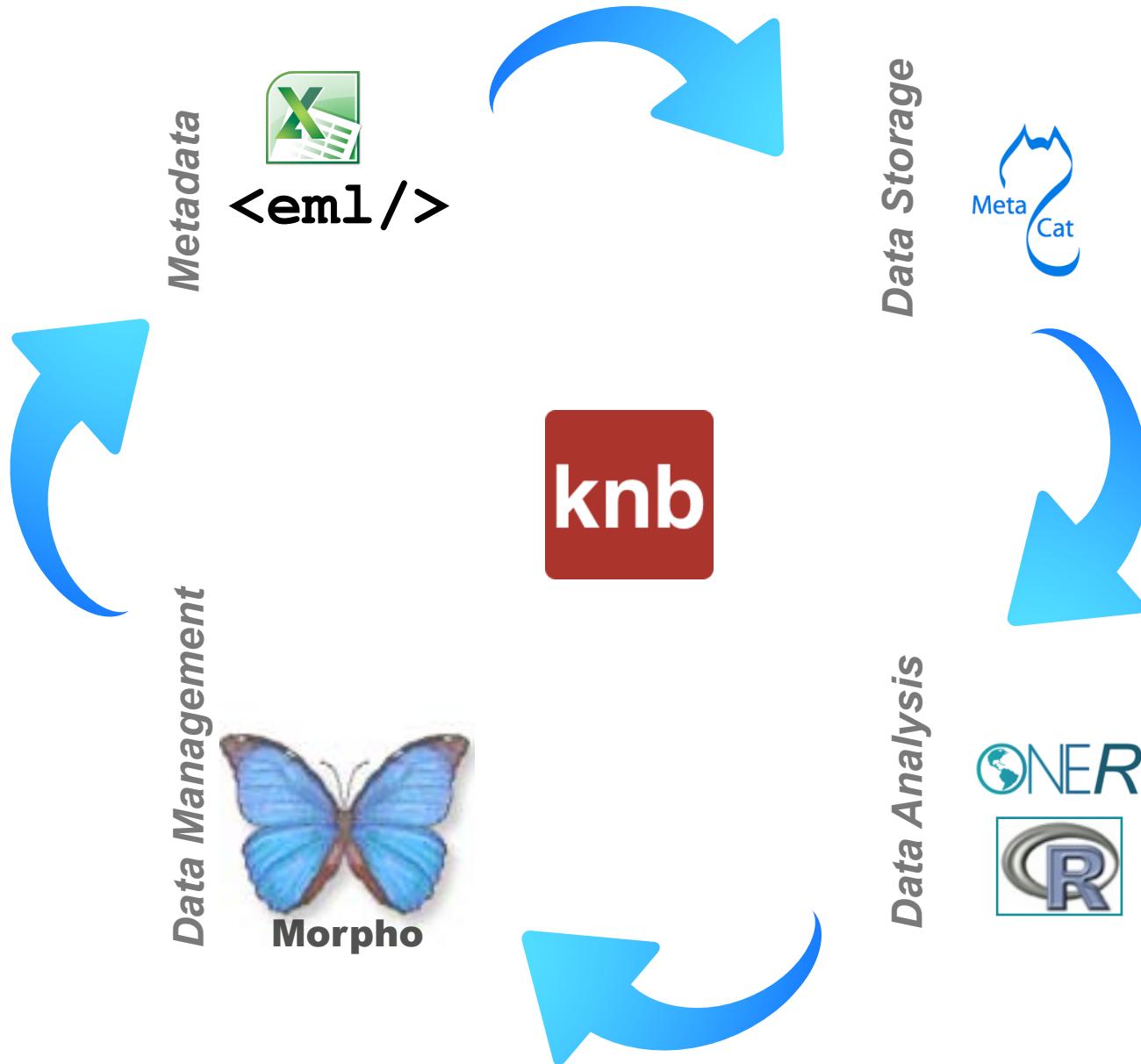
Metadata

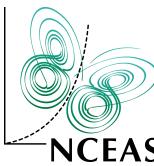
Data

Workflows

Results

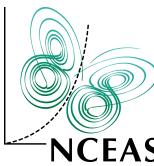
# Data life-cycle





# How do we harness the long tail?

- Efficient data federation
  - Focus on individual contributors
- Late binding in informatics systems
  - Loose coupling
  - Schema-less storage
- Central search for discovery
- Interoperable software



# Data Registration Activity

- <https://identity.nceas.ucsb.edu/>
- Register sample dataset
  - <https://dev.nceas.ucsb.edu/#share>



# Questions?

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