

Data at iDigBio

Feeding the Research Data Pipeline: long live the data



Deborah Paul, iDigBio Digitization and Training Specialist, [@idbdeb](#)
Missouri Botanical Garden, December 2 – 3, 2015

Content contributors: iDigBio staff - Alex Thompson, Joanna McCaffrey, Kevin Love, Matt Collins, Shelley James, Gil Nelson



Making data and images of millions of biological specimens available on the web

48,346,835

Specimen Records

13,055,353

Media Records

664

Recordsets

Search the Portal

Research



Why digitization matters

More about what we do and why



Digitization

Learn, share and develop best practices



Sharing Collections

Documentation on data ingestion



Working Groups

Join in, contribute, be part of the community



Proposals

New tool and workshop ideas



Citizen Scientists

How can you help biological collections?

<https://www.idigbio.org>
[@iDigBio](#)

Specimen Records Overview

Search Records Help Reset

search all fields

Must have image Must have map point

Filters | Mapping | Sorting | Download

Add a field Clear

Kingdom: x

Present Missing

Scientific Name: Add EOL Synonyms

Present Missing

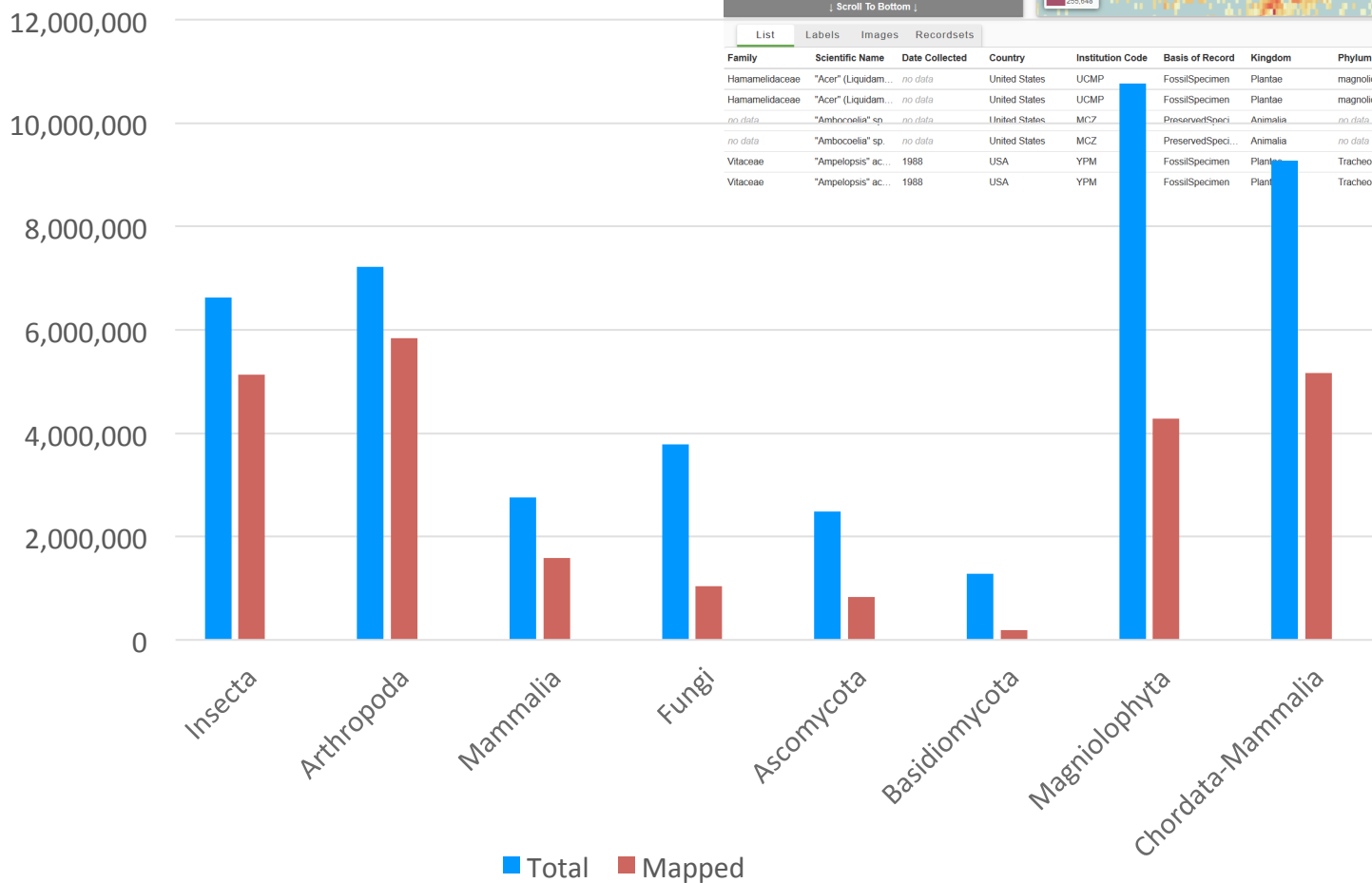
Date Collected: Start: End:

Present Missing

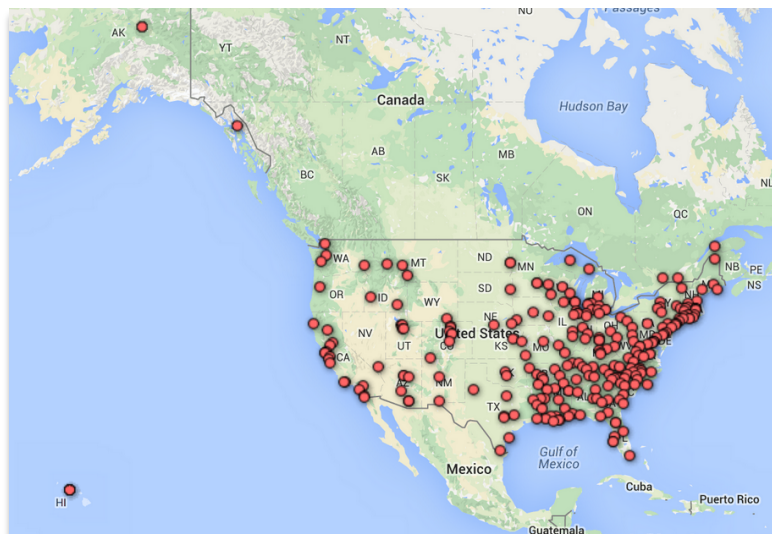
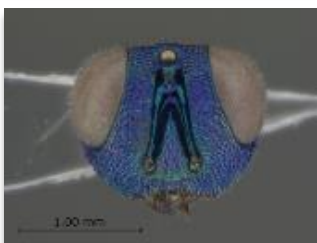
Scroll To Bottom

Total: 48,346,835

Family	Scientific Name	Date Collected	Country	Institution Code	Basis of Record	Kingdom	Phylum	Class	Order	Higher Taxon	Columns
Hamamelidaceae	"Acer" (Liquidam...	no data	United States	UCMP	FossilSpecimen	Plantae	magnoliophyta	Magnoliopsida	Hamamelidales	no data	view
Hamamelidaceae	"Acer" (Liquidam...	no data	United States	UCMP	FossilSpecimen	Plantae	magnoliophyta	Magnoliopsida	Hamamelidales	no data	view
	"Ambocoelia" sp.	no data	United States	MCZ	PreservedSpeci...	Animalia	no data	no data	no data	Animalia "Amboc...	view
	"Ambocoelia" sp.	no data	United States	MCZ	PreservedSpeci...	Animalia	no data	no data	no data	Animalia "Amboc...	view
Vitaceae	"Ampelopsis" ac...	1988	USA	YPM	FossilSpecimen	Plantae	Tracheophyta	magnoliopsida	vitales	Plantae, Tracheo...	view
Vitaceae	"Ampelopsis" ac...	1988	USA	YPM	FossilSpecimen	Plantae	Tracheophyta	magnoliopsida	vitales	Plantae, Tracheo...	view



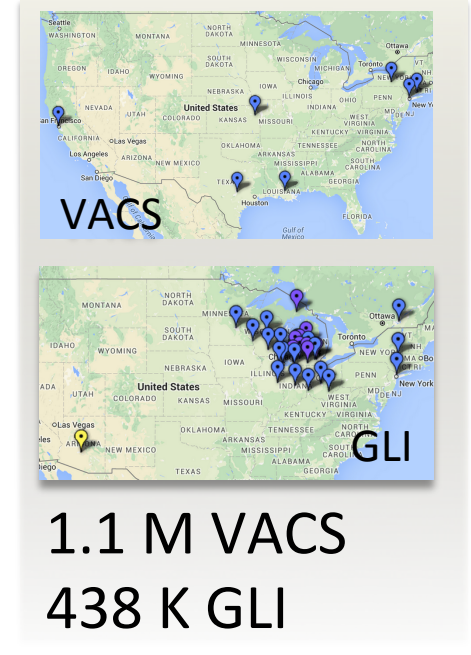
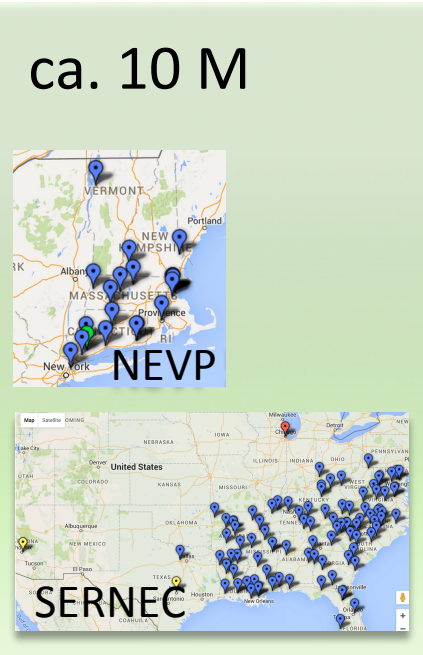
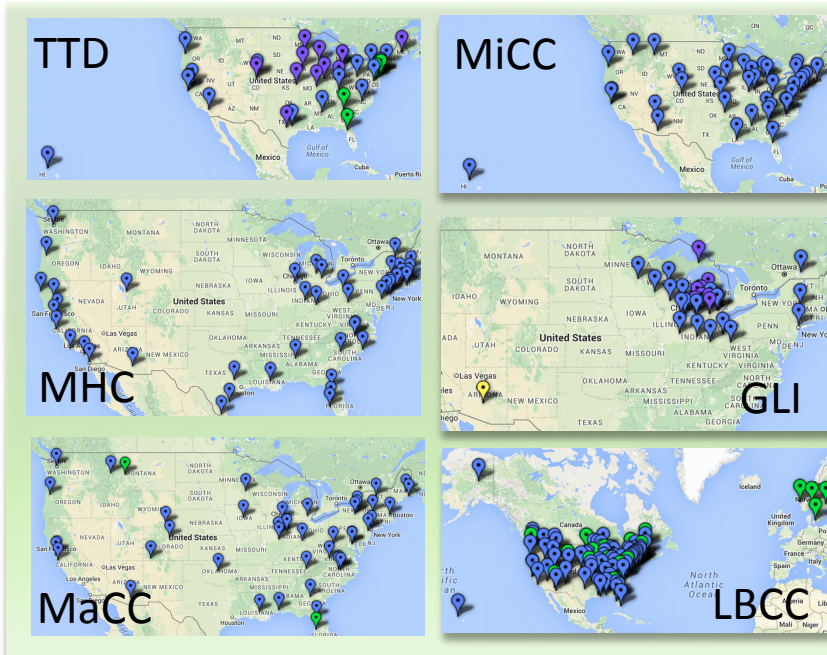
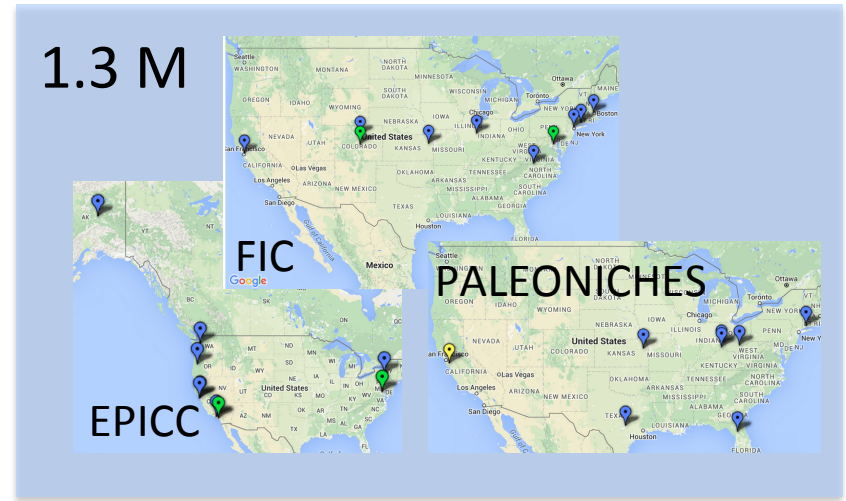
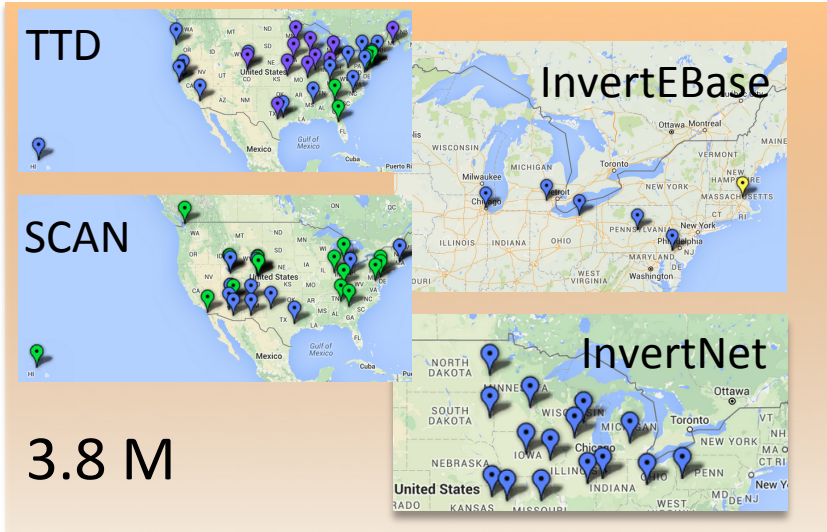
15 Thematic Collection Networks (TCNs)



Occurrence Data		Determination Hist.
Collector Info		
Catalog Number ?	Other Numbers ?	Collect
MSC-B-0000001		F.H. E
Associated Collectors		
J.E. Cantlon, A.L. Rebeck		
Latest Identification		
Scientific Name:		
Abietinella abietina		
ID Qualifier: ?		
Identified By:		
Locality		
Country	State/Province	
Locality:		
<input type="checkbox"/> Locality Security		
Latitude	Longitude	Uncertainty (met



- **InvertNet**: An Integrative Platform for Research on Environmental Change, Species Discovery and Identification (Illinois Natural History Survey, University of Illinois) <http://invertnet.org>
- **Plants, Herbivores, and Parasitoids**: A Model System for the Study of **Tri-Trophic Associations** (American Museum of Natural History) <http://tcn.amnh.org> (TTD)
- North American **Lichens and Bryophytes**: Sensitive Indicators of Environmental Quality and Change (University of Wisconsin – Madison) <http://symbiota.org/nalichens/index.php> <http://symbiota.org/bryophytes/index.php> (LBCC)
- Digitizing Fossils to Enable New Syntheses in Biogeography - Creating a **PALEONICHES**-TCN (University of Kansas)
- The Macrofungi Collection Consortium: Unlocking a Biodiversity Resource for Understanding Biotic Interactions, Nutrient Cycling and Human Affairs (New York Botanical Garden)(MaCC)
- Mobilizing New England Vascular Plant Specimen Data to Track Environmental Change (Yale University) (NEVP)
- **Southwest Collections of Arthropods Network (SCAN)**: A Model for Collections Digitization to Promote Taxonomic and Ecological Research (Northern Arizona University) <http://hasbrouck.asu.edu/symbiota/portal/index.php>
- **iDigPaleo**: Fossil Insect Collaborative: A Deep-Time Approach to Studying Diversification and Response to Environmental Change (FIC)
- Developing a Centralized Digital Archive of Vouchered Animal Communication Signals (VACS)
- The Macroalgal Herbarium Consortium: Accessing 150 Years of Specimen Data to Understand Changes in the Marine/Aquatic Environment (MHC)
- Documenting the Occurrence through Space and Time of Aquatic Non-indigenous Fish, Mollusks, Algae, and Plants Threatening North America's Great Lakes (GLI)
- **InvertEBase**: Reaching Back to See the Future: Species-rich Invertebrate Faunas Document Causes and Consequences of Biodiversity Shifts
- The Key to the Cabinets: Building and Sustaining a Research Database for a Global Biodiversity Hotspot (SERNEC)
- The Microfungi Collections Consortium: A Networked Approach to Digitizing Small Fungi with Large Impacts on the Function and Health of Ecosystems (MiCC)
- Documenting Fossil Marine Invertebrate Communities of the Eastern Pacific - Faunal Responses to Environmental Change over the last 66 million years (EPICC)



Kinds of data issues to think about

- 1100 unique strings in dwc:country field
- 1.7 M unique strings in dwc:scientificName
 - Using Apache Spark
- 6.7 M unique locality strings
 - The ones you need?
 - Georeferenced already?
 - With errors?
 - With metadata?

Flag
idigbio_isocountrycode_added ⓘ
dwc_continent_added ⓘ
dwc_kingdom_replaced ⓘ
dwc_class_added ⓘ
dwc_phylum_added ⓘ
dwc_kingdom_added ⓘ
dwc_order_added ⓘ
rev_geocode_eez ⓘ
dwc_phylum_replaced ⓘ

KINGDOMS

Taxa within GBIF backbone kingdoms.



- [Animalia](#) 1,913,825
- [Archaea](#) 463
- [Bacteria](#) 14,258
- [Chromista](#) 18,186
- [Fungi](#) 319,366
- [Plantae](#) 953,712
- [Protozoa](#) 41,539
- [Viruses](#) 5,724
- [Other](#) 6,248

RANKS

Number of accepted taxa by ranks.



- [kingdom](#) 9
- [phylum](#) 111
- [class](#) 272
- [order](#) 1,183
- [superfamily](#) 356
- [family](#) 20,760
- [genus](#) 390,790
- [species](#) 2,497,114
- [infraspecific name](#) 691
- [subspecies](#) 171,473
- [infrasubspecific name](#) 1
- [variety](#) 149,011
- [subvariety](#) 1,070
- [form](#) 40,473
- [subform](#) 7

Names

There are [1,143,026](#) synonyms in this dataset.

UNIQUE NAMES

There are 4,410,899 unique names in this dataset. On average 0.124% of the names are found in more than one taxon.

2,497,114 | **4,416,347**

Species

Taxa

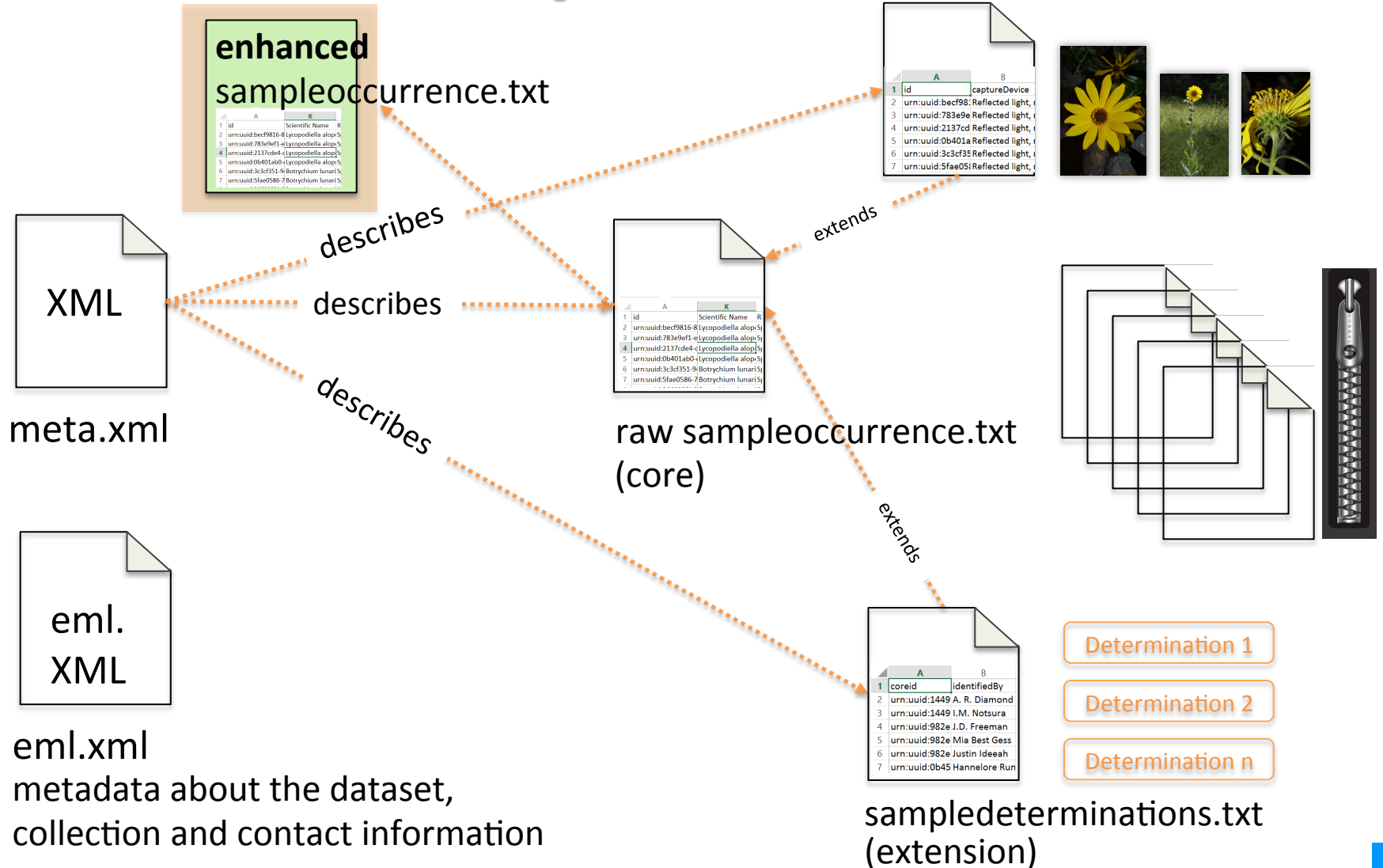
[View species](#)

Citation GBIF Secretariat: GBIF Backbone Taxonomy, 2013-07-01.

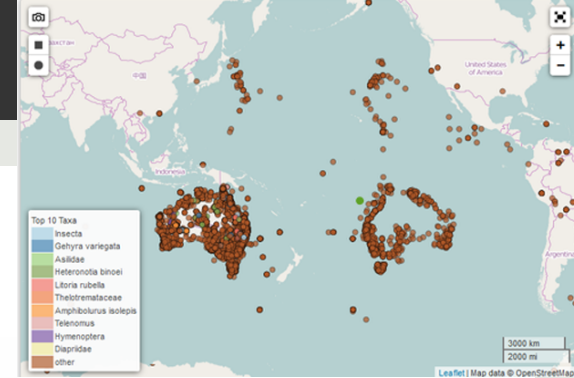
Accessed via <http://www.gbif.org/dataset/d7dddbf4-2cf0-4f39-9b2a-bb099caae36c> on 2015-12-01

Rights CC BY 3.0

What's Inside a DwC-A at iDigBio?



More data data quality flags



Data Corrected Data Use Raw

This table shows any data corrections that were performed on this recordset to improve the capabilities of iDigBio [Search](#). The first column represents the correction performed. The last two columns represent the number and percentage of records that were corrected. A complete list of the data quality flags and their descriptions can be found [here](#). Clicking on a data flag name will take you to a search for all records with this flag in this recordset.

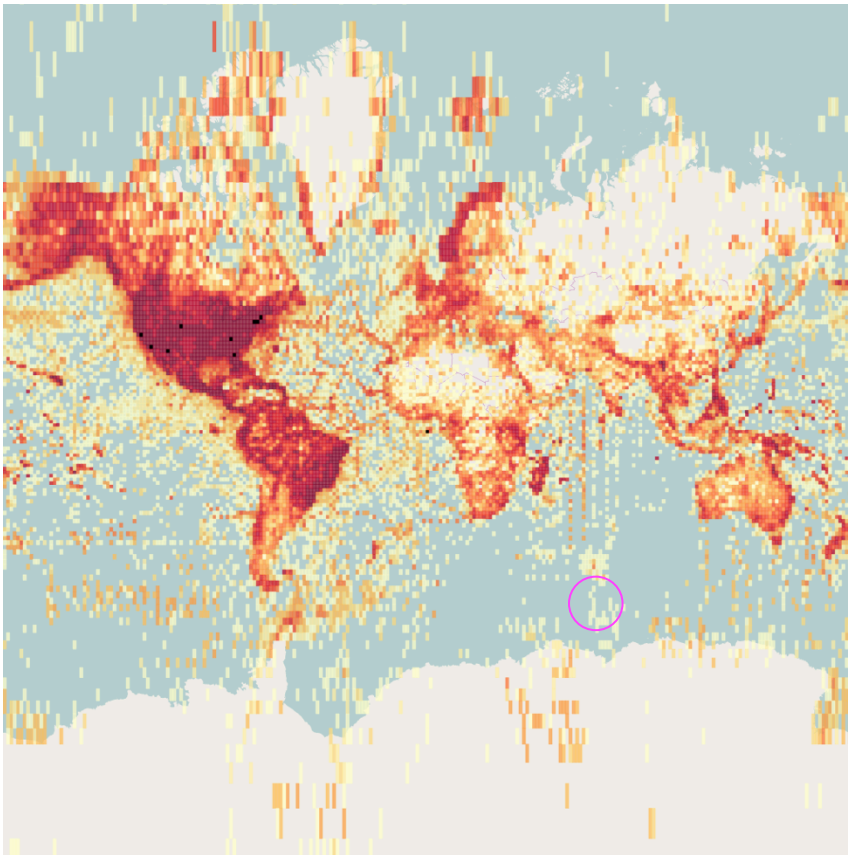
Flag	Records With This Flag	(%) Percent With This Flag
idigbio_isocountrycode_added ⓘ	2370	98.709
dwc_continent_added ⓘ	2369	98.667
dwc_kingdom_replaced ⓘ	784	32.653
dwc_class_added ⓘ	370	15.41
dwc_phylum_added ⓘ	370	15.41
dwc_kingdom_added ⓘ	365	15.202
dwc_order_added ⓘ	360	14.994
rev_geocode_eez ⓘ	339	14.119
dwc_phylum_replaced ⓘ	337	14.036

Downstream Use Cases

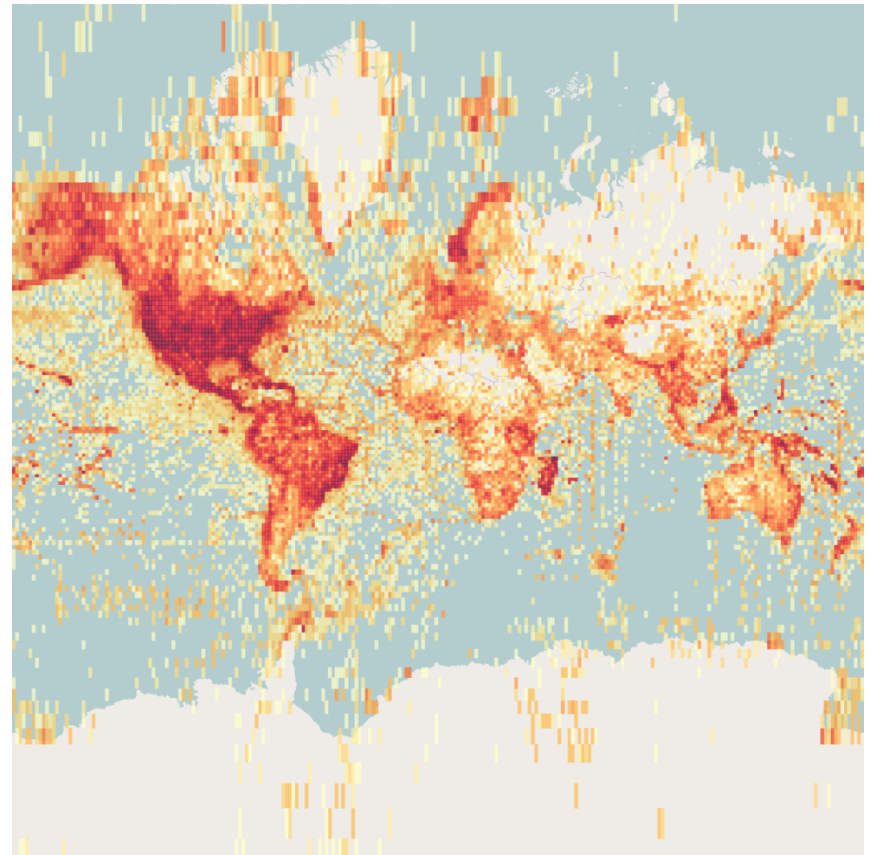
- Environmental Ontology Checking - Pier Buttigieg & Grant Godden <https://www.idigbio.org/content/webinar-shaping-semantic-layer-mining-digitised-data-encounter-between-idigbios-plant>
- Data Cleaning - Heather Appleby & Katja Seltmann <https://www.idigbio.org/content/summer-learning-r-clean-data-idigbio-portal-recordset-correction-feature>
- Lifemapper - KU & ACIS Collaboration <http://lifemapper.acis.ufl.edu>
- PhyloJive - Joe Miller & ACIS Collaboration <http://phylojive.acis.ufl.edu>
- Proximity and Correlation: Two new computer programs for mining phytosociological information held in herbarium databases using central Arizona as a test case. Daryl Lafferty & Les Landrum, Arizona State University <http://dx.doi.org/10.12705/645.9> *Taxon*, Volume 64, Number 5, 28 October 2015, pp. 998-1016(19) https://www.idigbio.org/wiki/index.php/Specimens_Full_Circle_SPNHC_2015

Advanced Mapping API Capabilities

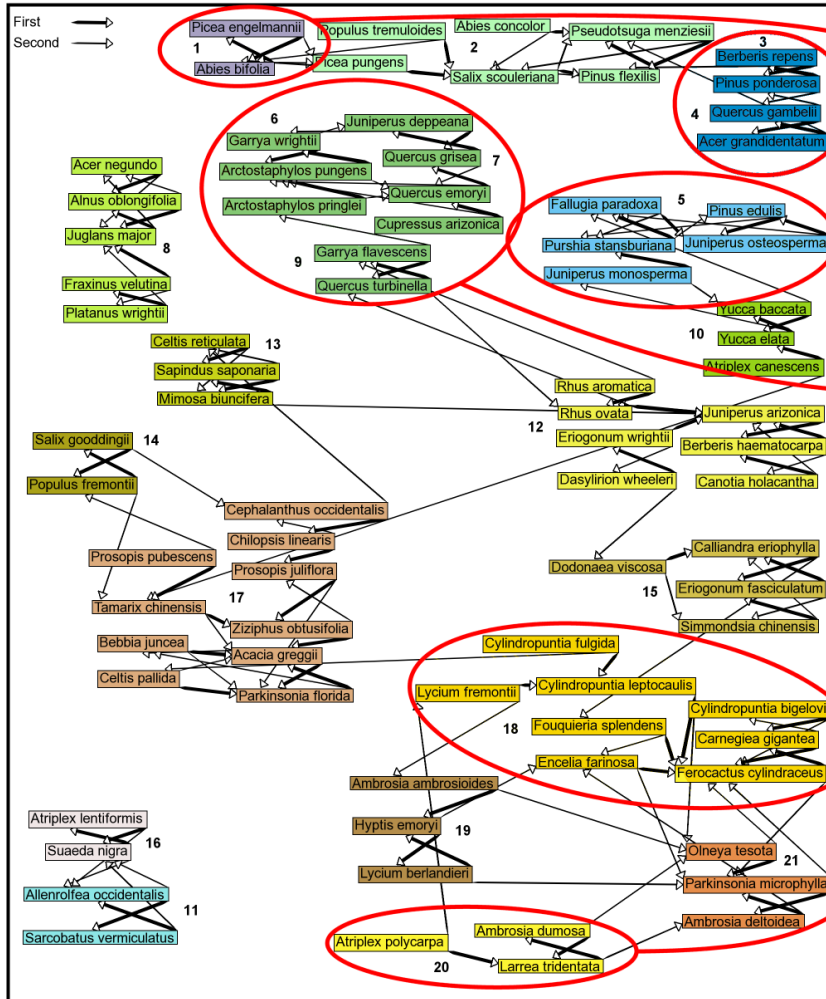
“Species Density” Map



Specimen Density Map

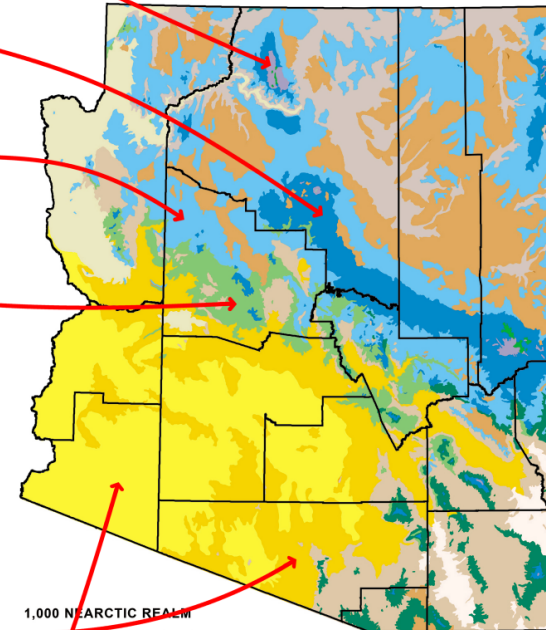


Mining and Using Associated Taxa Data



BIOTIC COMMUNITIES OF ARIZONA

David F. Brown, Thomas C. Brennan, and Andrew T. Holycross



1,000 NEARCTIC REALM

- 1111 ARCTIC AND ALPINE TUNDRAS
 - 111.0 Rocky Mountain and Great Basin Alpine Tundra
- 121 BOREAL AND SUBALPINE FORESTS AND WOODLANDS
 - 121.3 Rocky Mountain and Great Basin Subalpine Conifer Forest
- 122 COLD TEMPERATE FORESTS AND WOODLANDS
 - 122.6 Rocky Mountain Monocot Conifer Forest
 - 122.7 Great Basin Conifer Woodland
- 123 WARM TEMPERATE FORESTS AND WOODLANDS
 - 123.3 Madroan Evergreen Forest and Woodland
- 133 WARM TEMPERATE SCRUBLANDS
 - 133.3 Southwestern (Arizona) Interior Chaparral
- 142 COLD TEMPERATE GRASSLANDS
 - 142.12 Plains and 142.13 Intermountain Grassland

- 143 WARM TEMPERATE GRASSLANDS
 - 143.1 Semidesert Grassland
- 152 COLD TEMPERATE DESERTLANDS
 - 152.1 Great Basin Desertscrub
- 153 WARM TEMPERATE DESERTLANDS
 - 153.1 Mohave Desertscrub
 - 153.2 Chihuahuan Desertscrub
- 3,000 NEOTROPICAL REALM
- 154 TROPICAL-SUBTROPICAL DESERTLANDS
 - 154.1 Sonoran Desertscrub
 - 154.11 Lower Colorado River Valley Subdivision
 - 154.12 Arizona Uplands Subdivision

Plantae records pulled via portal and API

owltools TSV export from ENVO's OWL representation

EnvO

R

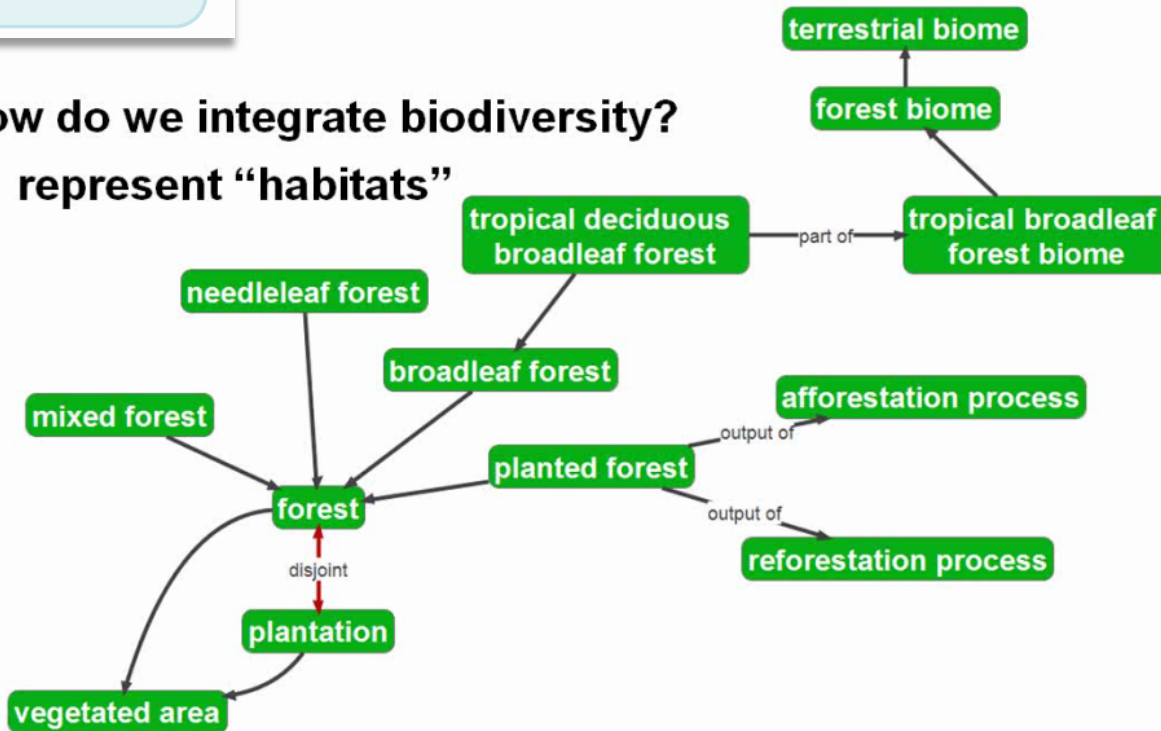
Text mining iDigBio records

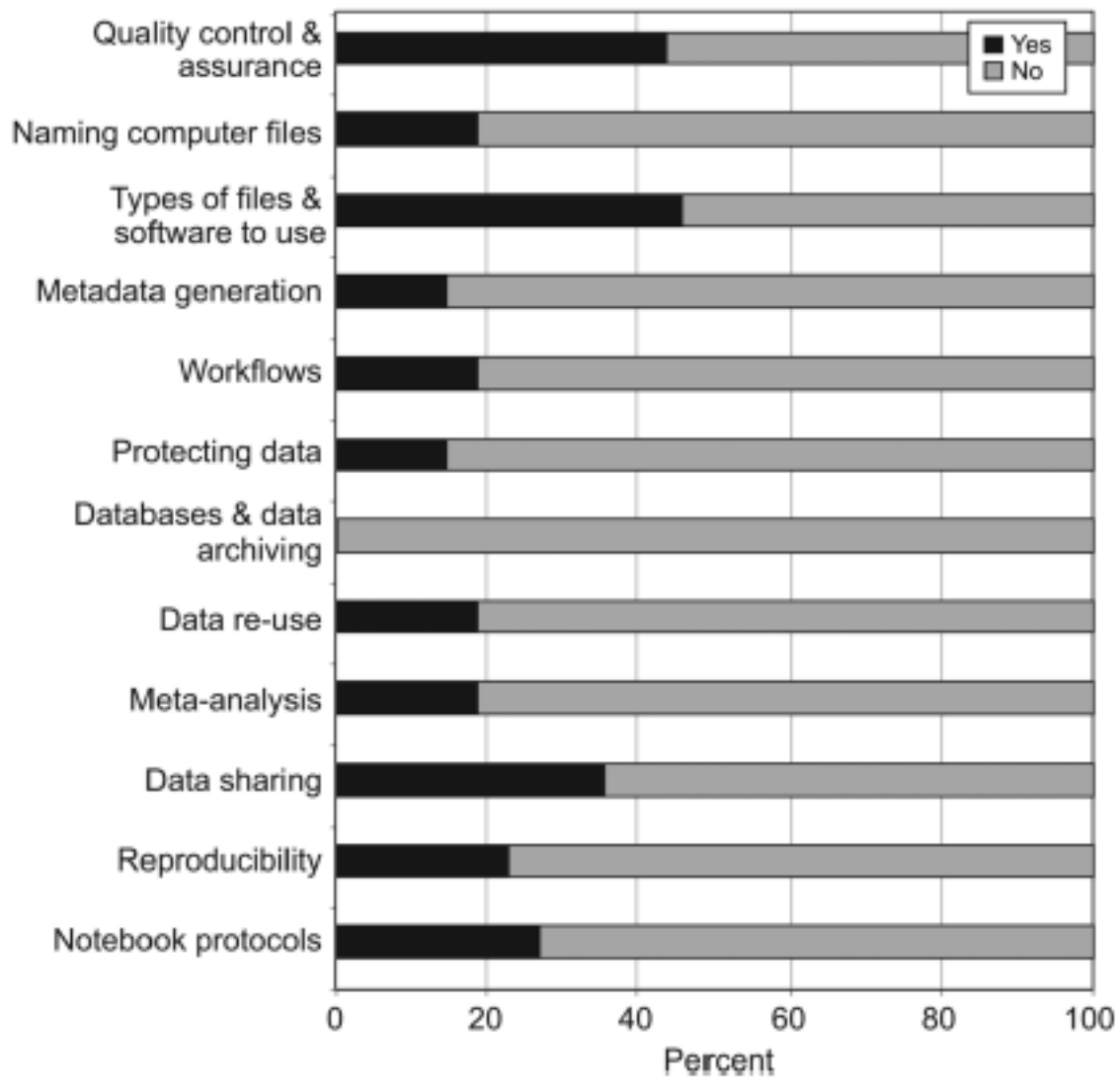
- R packages openNLP, tm, RWeka
- defined domain- and attribute- specific
 - stopwords (s, n, sw, se ...)
 - delimiters (among, between, ...)
- scanned for common n-grams
- explored associations of common terms
- began to align ENVO's content to our observations

Mining Aggregated Data to Enhance Ontologies

In support of future research

How do we integrate biodiversity? represent "habitats"





% Ecology courses that address / teach the listed data management principles



The path ahead

- refine our corpora and TM approaches to deliver robust and more directly actionable output
 - account for synonymy and, if feasible, multilingual issues
- engage experts on specific taxa to evaluate the representation of environments relevant to their work
 - create taxon-specific mining routines and link to efforts in the semantic representation of habitats
- Perform gap analyses using ENVO as a reference to find undersampled environment types
- Pre-compose classes in ENVO and related ontologies which align to the most frequent (or important) environment descriptions in the available corpora

3-gram	Frequency
montane rain forest	41
evergreen broadleaf forest	38
subtropical evergreen broadleaf	36
substrate loam disturbance	34
baserock granite substrate	22
granite substrate loam	20
disturbance agriculture slope	15
loam disturbance agriculture	15

Providing semantic feedback to iDigBio?



engaging botanists (and other biologists), text-miners, ontologists, database engineers is the key to creating semantically-aware, approachable, and sustainable infrastructure

digitiser training is

Looking forward to more data and enhanced data being used in research! Up next, the TCNs...



www.idigbio.org



facebook.com/iDigBio



twitter.com/iDigBio



vimeo.com/idigbio



idigbio.org/rss-feed.xml



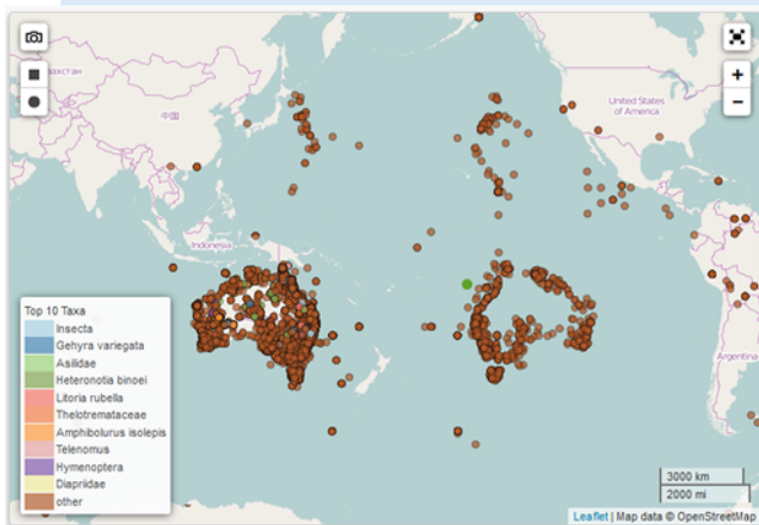
<webcal://www.idigbio.org/events-calendar/export.ics>

<http://tinyurl.com/GWGmikeyostGPSapps>

iDigBio Data Quality (DQ) Flags enhance Georeferencing Workflows!

iDigBio DQ flags making it easier to spot and fix georeferencing issues.
Great for researchers, great for data providers

See <https://www.idigbio.org/portal/publishers>



Search Records Help

search all fields

Must have image Must have map point

Filters Mapping Sorting Download

Add a field Clear

Data Flags rev_geocode_mismatch

Present Missing

Flag
dwc_class_added i
dwc_kingdom_added i
dwc_order_added i
dwc_phylum_added i
idigbio_isocountrycode_added i
dwc_continent_added i
geopoint_datum_missing i
rev_geocode_eez i
dwc_stateprovince_replaced i
dwc_country_replaced i
rev_geocode_mismatch i
rev_geocode_lon_sign i
rev_geocode_lat_sign i
geopoint_similar_coord i
rev_geocode_flip_lat_sign i
datecollected_bounds i

A Locality Service?

What is it? It can speed up georeferencing!

Ask us about it. Let's reduce re-georeferencing.

- Data by organismal group
 - Unique data
 - Measurement or fact
 - **Host associations**
 - Associated source materials, like recordings
 - **Locality data** for mining
- Mining aggregated data
 - Enriching ontologies to support research
 - ENVO
 - Other ontologies that could do this
- The data sets.
- Ecologists want to know
- Using data to test models
 - Historic (back-casting)
- Clarifying fit-for-useness
 - How to indicate what type of research a dataset is good for
 - Georef to centroid?
 - Fine-grain georef?
 - Tdww
- Aggregators, Data Generators, Researchers – friends with benefits (#valueOfAggregation)
- Unique data
 - Habitat
 - Localities

What we hear

- *I've borrowed my colleague's computer.*
- Scaling Up?: I'm running analysis on three different computers.
- Excel is a database, isn't it?
- Should I learn R (Python,...)? Is it worth my time?
- How do I visualize this data?
- What's the best way to share data with colleagues?
- How do I work with a txt (csv, hdf5,...) file?
- It's difficult to replicate my (others) research.
- What is an API?
- How do I use APIs to enhance my research, ...

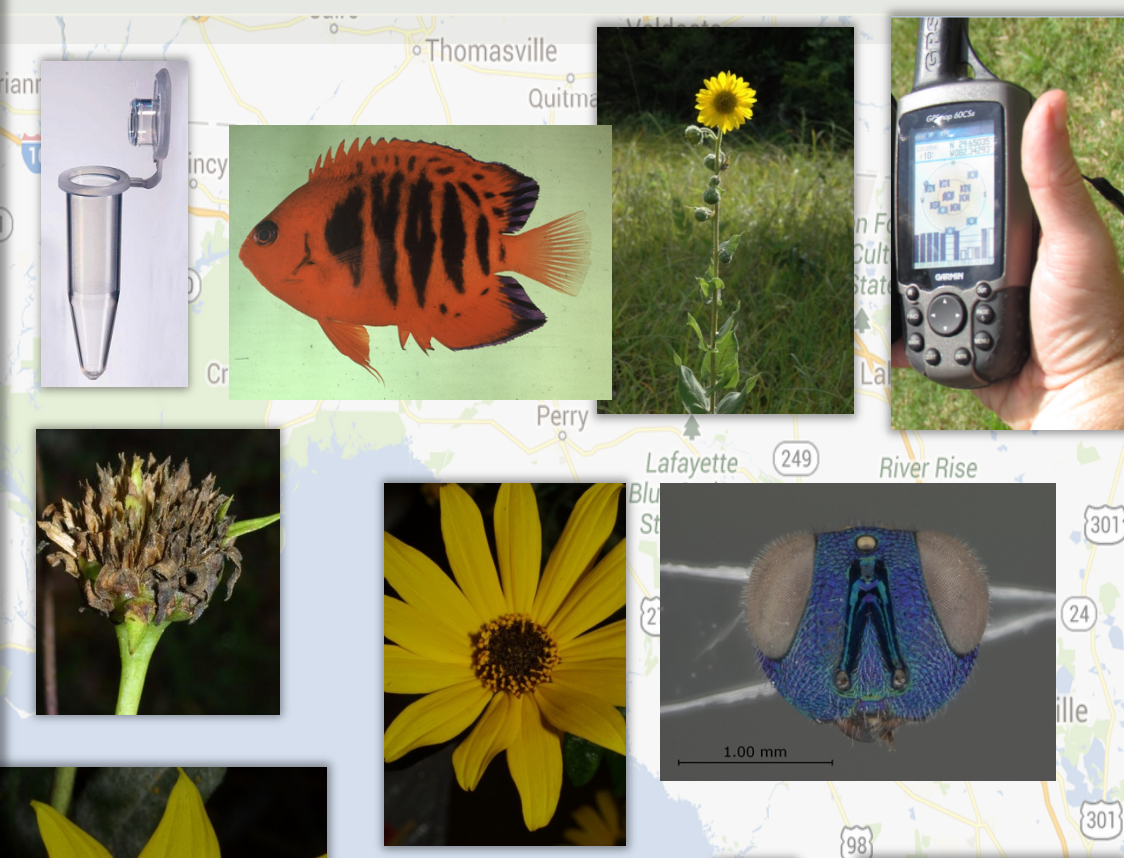
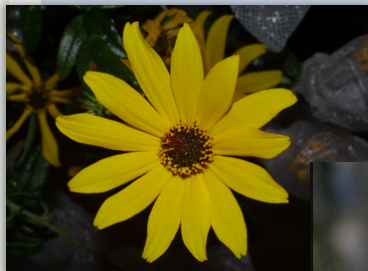
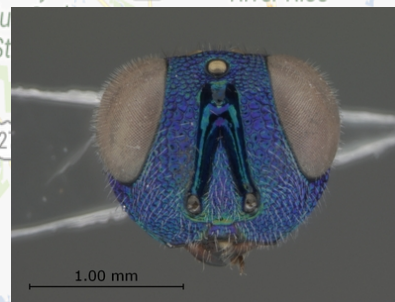
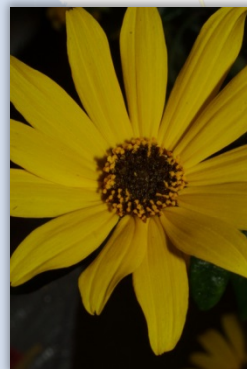
What do we hear after the workshops?

- “the course is well suited for beginners; I received the courage to venture and learn more about R”
- “I feel I can start using R again”
- “I appreciate all the hands-on work”
- “the class is well organized and easy to follow”
- “really great working through all the examples/challenges, lots of examples = lots of practice = understanding”
- wish I’d known this before graduate school!
- we need to hire a data manager.
- I can now give cleaner data to aggregators like iDigBio and GBIF

What are some of these biodiversity informatics skills and who needs them?

- For the researcher
- Collections Manager
- Data Manager
- Student
- Standards
- File management, organization, and literate programming
- Better Spreadsheet Skills!
- Workflows
 - Digitization
 - Research
- Data Analysis
- Data Visualization
- Data Mobilization





Data and Metadata.

It's about **discovery** and **data re/use**.

It's about **feedback** and **accountability**.

It's about **credit** and **attribution**.

Make sure your data's not under a rock.

