

# Challenges Using Biodiversity Data: examples of what to look for

*presenter:* Arctic Data Center, DataONE, Environmental Data Initiative, ESIP, GBIF, iDigBio, NEON



#datahelpdesk

Ecological Society of America 2019 ESAUSSEE

Career Fair Center in the Exhibit Hall <https://esa.org/louisville/career-fair/>

Monday 12 August 230-300 PM

<http://bit.ly/datahelpesa2019>





# Advancing the Digitization of Biological Collections

## iDigBio Hub and Thematic (Museum) Collection Networks

total: 119,768,942

### Digitization

Workflows & Protocols  
Dissemination

### Research Use

Cyberinfrastructure  
Tool collaboration  
Portal development  
ENM workshop

### Research focus

Data quality  
APIs

### Training

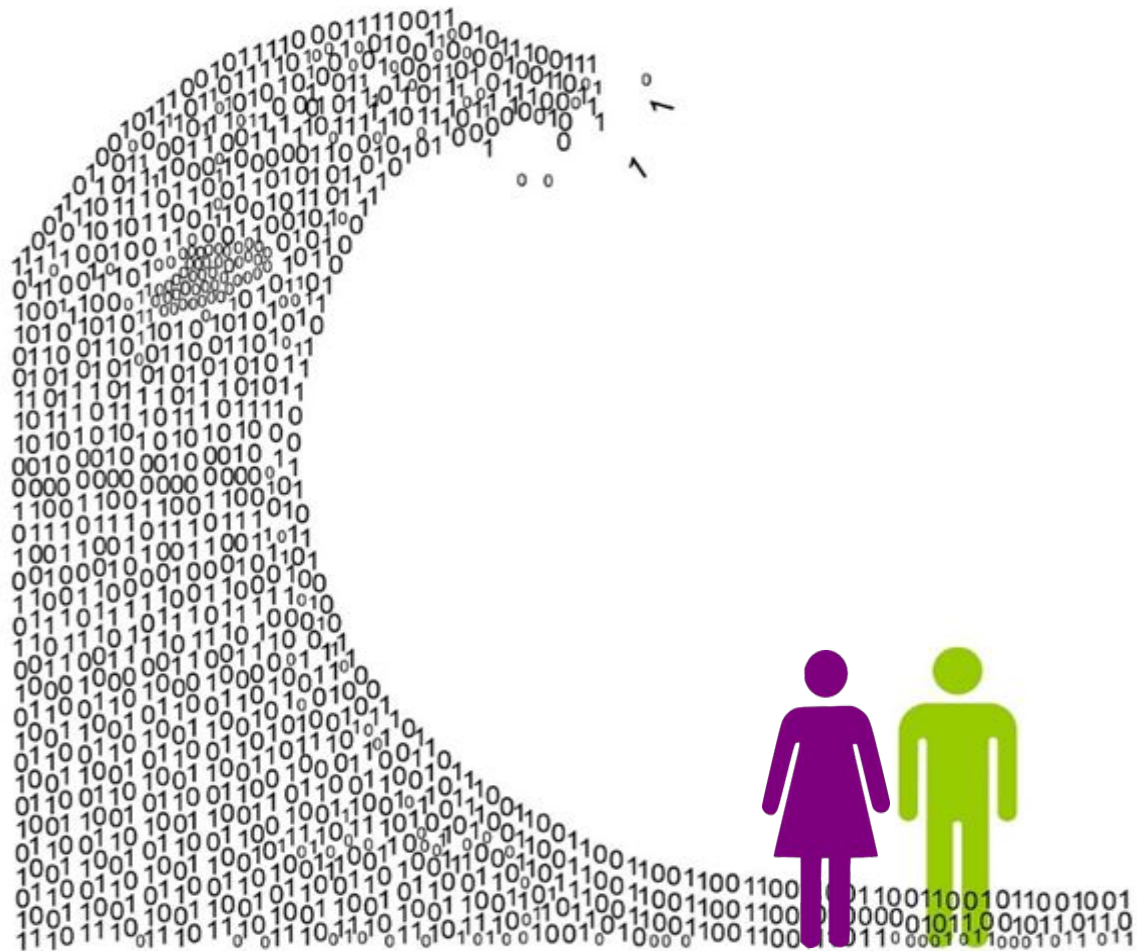
Biodiversity informatics  
Data skills and literacy  
Collections software  
Imaging  
Project Management

### Education Outreach

Citizen Science  
K-12 materials  
Undergraduate  
Fossil Clubs  
Mentor teachers

### Methods

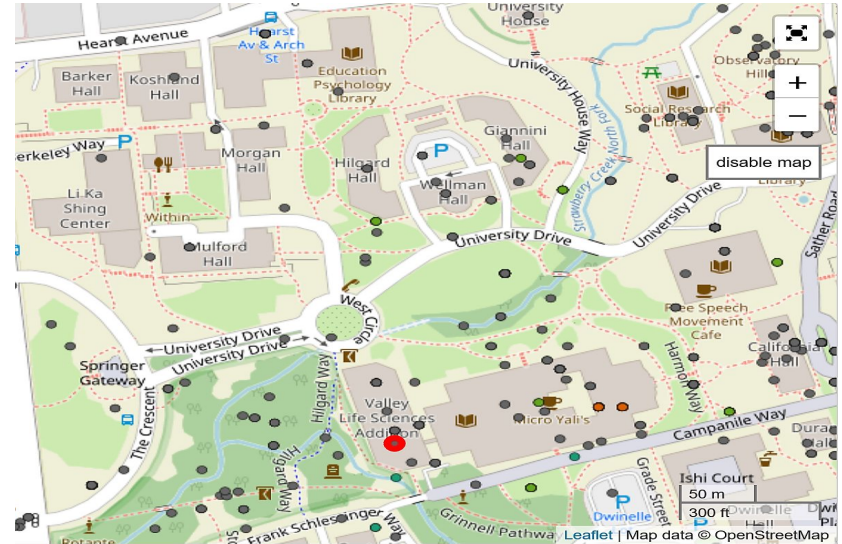
Workshops  
Webinars  
Symposia  
Conferences  
Working Groups  
Short Courses  
Adobe Connect  
Listservs  
Publications  
Social Media @idigbio



# Challenges researchers face

**exercise:** from the ecologist, collector, policy maker, or other downstream user of data (e.g. collections), what issues need to be addressed before applying your research methods?

1. List three major challenges
2. Actions, tools, and influence

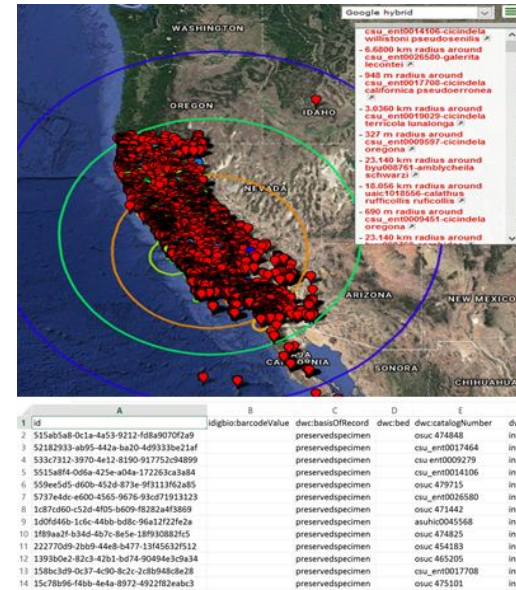




# Data lessons compiled - inspired by workshop

## *Georeferencing for Research Use of Museum Collections Data*

- Data mapped to standards
  - supports use and re-use (e.g. Darwin Core DwC, Ecological Metadata Language EML)
  - standards help with data validation and cleaning
- Data have issues
  - what are some you have experienced
  - need to be addressed before applying research methods
  - keep raw data raw
  - track your changes
- Issues can be grouped
  - Taxonomic / Nomenclatural
  - Locality / Place / Georeference
  - Time / Date
- Data visualization is key
  - QGIS lessons
  - Open Refine
  - R, etc.
- For future data - use data standards



Carabidae (beetles) of California



# Synthesis of issues

*evaluating the research fitness-for-use of these data*

*creating a list of data quality checks*

- Timey-wimey stuff
  - date issues like formats
- Geography
  - place name issues
  - out of expected bounds
  - missing metadata
- Taxonomy
  - taxon name issues\*
  - concepts
  - authority files
  - parsing

doi: 10.3897/rio.4.e32449

**GRU Workshop Conversation on Data Quality Considerations and Checks.**

**Data quality (dq) considerations and checks – an annotated bibliography**

Workshop participants discuss data quality checks

**Georeferencing for Research Use (GRU): An integrated geospatial training paradigm for biocollections researchers and data providers**

▼ Katja C. Seltmann, Sara Lafia, Deborah L. Paul, Shelley A. James, David Bloom, Nelson Rios, Shari Ellis, Una Farrell, Jessica Utrup, Michael Yost, Edward Davis, Rob Emery, Gary Motz, Julien Kimmig, Vaughn Shirey, Emily Sandall, Daniel Park, Christopher Tyrrell, R. Sean Thackerdeen, Matthew Collins, Vincent O'Leary, Heather Prestridge, Christopher Evelyn, Ben Nyberg

**Abstract**

Georeferencing is the process of aligning a text description of a geographic location with a spatial location based on a geographic coordinate system. Training aids are commonly created around the georeferencing process to disseminate community standards and ideas, guide accurate georeferencing, inform users about new tools, and help users evaluate existing geospatial data. The *Georeferencing for Research Use (GRU)* workshop was implemented as a training aid that focused on the creation and research use of geospatial

**Authors**

- **Katja C. Seltmann** - Corresponding author  
Credentia Center for Biodiversity and Ecological Restoration, University of California - Santa Barbara, Santa Barbara, CA, United States of America  
Articles by this author in: Crossref | PubMed | Google Scholar
- **Sara Lafia**  
Department of Geography, Center for Spatial Studies, University of California - Santa Barbara, Santa Barbara, CA, United States of America  
Articles by this author in: Crossref | PubMed | Google Scholar
- **Deborah L. Paul** - Corresponding author  
iDigBio, Gainesville, United States of America  
Florida State University, Tallahassee, United States of America  
Articles by this author in: Crossref | PubMed | Google Scholar

**states (or other placeholder values researchers use to represent field blank when no date (or other information) is available, rather than a placeholder.**



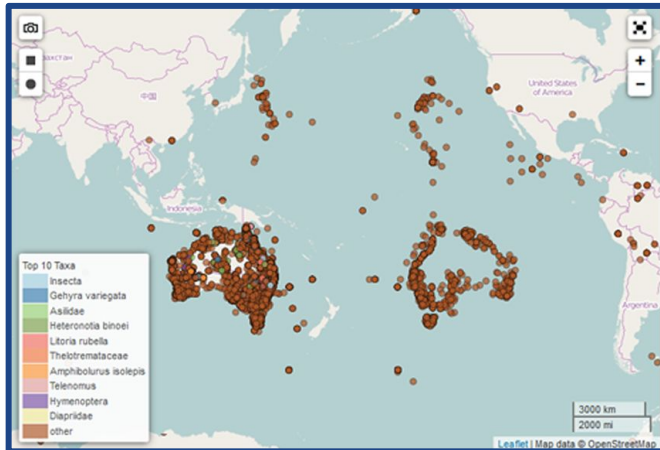
# Grouping data issues

- time (date)
- locality
- place

**Hannah Frost**  
@feeffofannah

⚙️
Following

From a [@HydralnABox](#) interview: "People will put anything and their dog in the date field. It's absolutely astonishing."



**Country**

✕

- united kindgom
- united king
- united kingdom
- united kingdom (england)
- united kingdom (scotland)
- united kingdom (wales)
- united kingdom [?]
- united kingdom of great b
- united kingdom?

**List**

---

**Family**



# Date and time issues



Following

From a [@HydralnABox](#) interview: "People will put anything and their dog in the date field. It's absolutely astonishing."





# The date field (and others)

1. Order matters for dates
  - use yyyy-mm-dd
2. What to do when you have no value to share?
  - leave it blank
  - avoid placeholder values
3. 0 has meaning

dwc:eventDate	dwc:eventDate
1900-01-01	1900-01-01
1900-04-01	1900-04-01
03-04-03	1903-04-03
1901-08-17	1901-08-17
1901-08-17	1901-08-17
0	
999	
1903-05-02	1903-05-02
1903-05-02	1903-05-02
1903-05-02	1903-05-02



# Season of observation / collection

1. Does the date fit the organism?
2. Is it an outlier?
3. Or an error?

CM	CN
dwc:eventDate	dwc:scientificName
1900-01-01	calosoma semilaeve
1900-04-01	platynus brunneomarginatus
1903-04-03	calathus ruficollis ruficollis
1901-08-17	cicindela trifasciata sigmaidea
1901-08-17	cicindela trifasciata sigmaidea
	cicindela senilis
	cicindela senilis
1903-05-02	omus tularensis
1903-05-02	omus californicus
1903-05-02	omus tularensis



# From gene names to dates?

Ziemann *et al.* *Genome Biology* (2016) 17:177  
DOI 10.1186/s13059-016-1044-7

Genome Biology

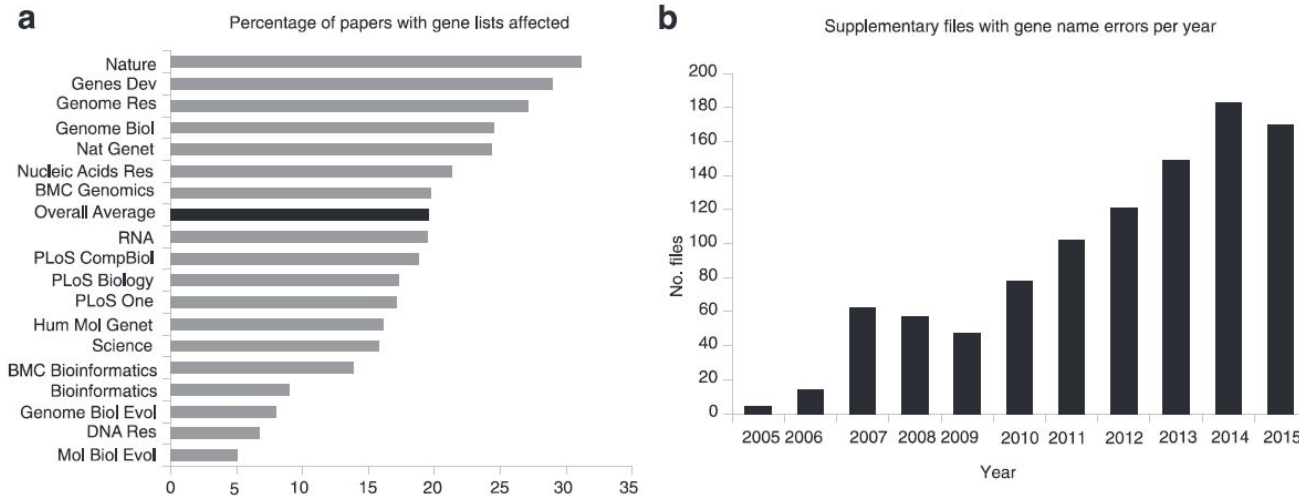
COMMENT

Open Access

## Gene name errors are widespread in the scientific literature



Mark Ziemann<sup>1</sup>, Yotam Eren<sup>1,2</sup> and Assam El-Osta<sup>1,3\*</sup>

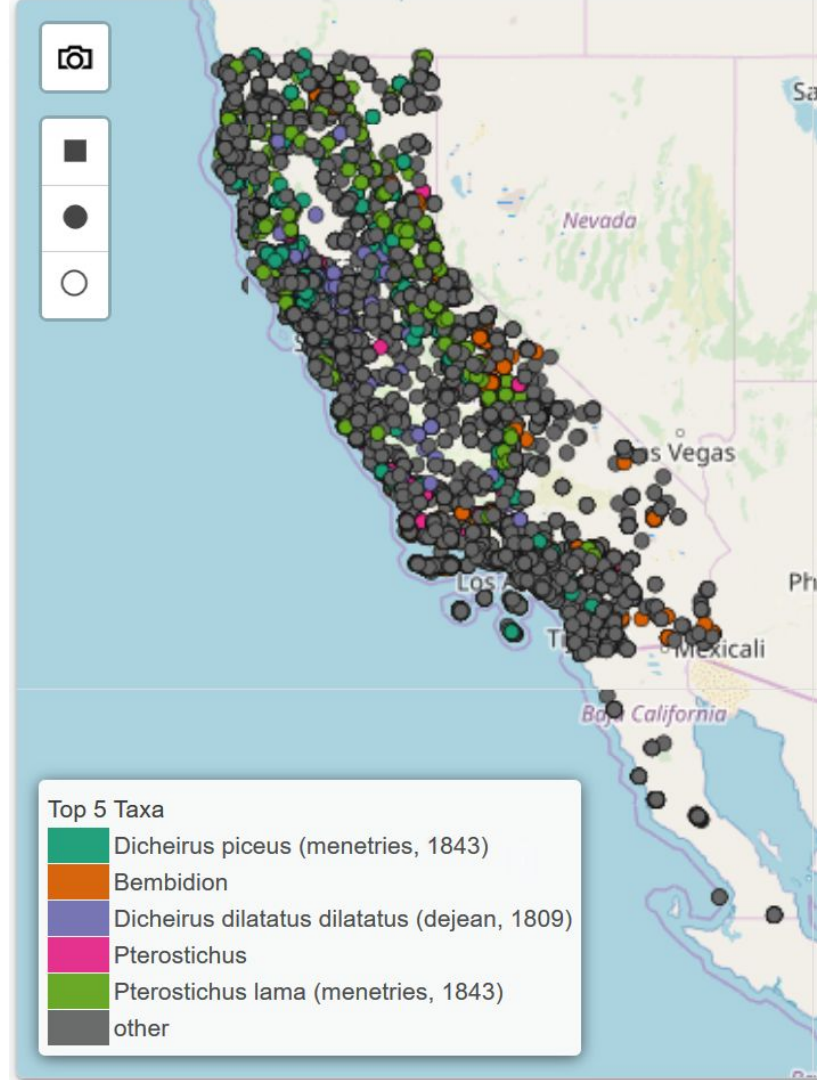


**Fig. 1** Prevalence of gene name errors in supplementary Excel files. **a** Percentage of published papers with supplementary gene lists in Excel files affected by gene name errors. **b** Increase in gene name errors by year



## Taxonomic names

1. Check spatial bounds for taxa
  - a. expected or not?
2. Check downloaded data
  - a. raw and tweaked data comparison needed





# What level of taxonomic determination is required for your research question?

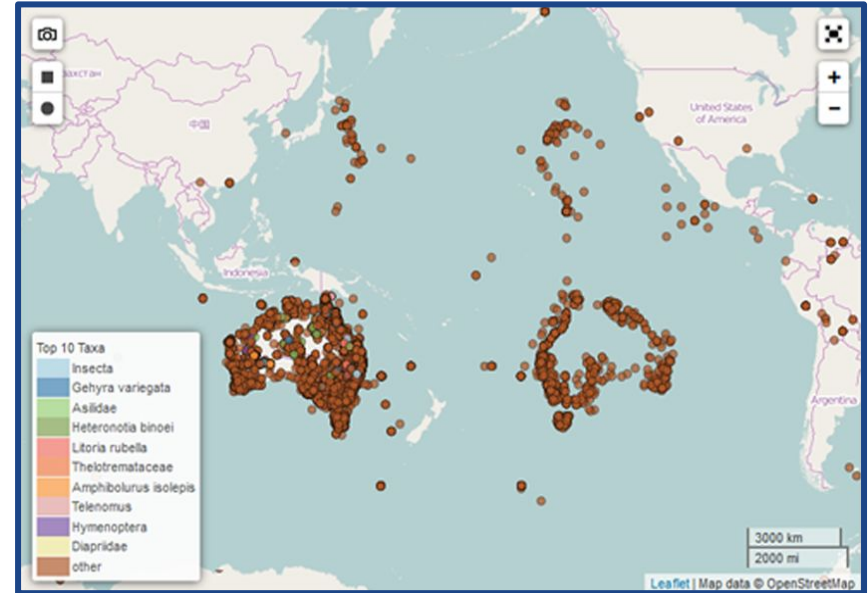
1. Remove records with undetermined, indet, empty
2. Use **dwc:taxonRank** to remove data outside ranks needed
3. Check raw and automated taxonomic names
4. Check endings - are they really the same?
5. Check synonymies, decide which names to use

CN	CO
dwc:scientificName	dwc:taxonRank
calosoma semilaeve	species
platynus brunneomarginatus	species
calathus ruficollis ruficollis	subspecies
carabidae	family
cicindela trifasciata sigmoidea	subspecies
cicindela senilis	species
indet.	
omus tularensis	species
omus californicus	species
omus	genus
omus tularensis	species
omus tularensis	species
pterostichus lama (menetries, 1843)	species



# Got geopoints?

1. Map them
2. use layers to check reasonableness
  - a. bounding boxes
  - b. habitat, terrain
3. look for outliers and transpositions
4. Place could suggest cultivated, port, or zoo record





# Presence data or precise data?

1. How fine-grained of a georeference do you need for your research question/model?
2. How many decimal places should make you suspect a data point? (fake precision)

WHAT THE NUMBER OF DIGITS IN YOUR COORDINATES MEANS

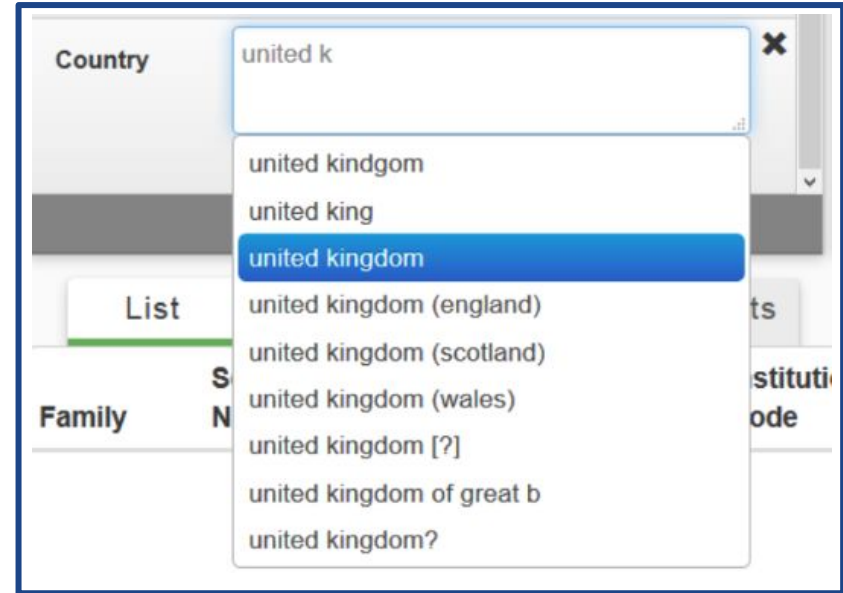
LAT/LON PRECISION	MEANING
28°N, 80°W	YOU'RE PROBABLY DOING SOMETHING SPACE-RELATED
28.5°N, 80.6°W	YOU'RE POINTING OUT A SPECIFIC CITY
28.52°N, 80.68°W	YOU'RE POINTING OUT A NEIGHBORHOOD
28.523°N, 80.683°W	YOU'RE POINTING OUT A SPECIFIC SUBURBAN CUL-DE-SAC
28.5234°N, 80.6830°W	YOU'RE POINTING TO A PARTICULAR CORNER OF A HOUSE
28.52345°N, 80.68309°W	YOU'RE POINTING TO A SPECIFIC PERSON IN A ROOM, BUT SINCE YOU DIDN'T INCLUDE DATUM INFORMATION, WE CAN'T TELL WHO
28.5234571°N, 80.6830941°W	YOU'RE POINTING TO WALDO ON A PAGE
28.523457182°N, 80.683094159°W	"HEY, CHECK OUT THIS SPECIFIC SAND GRAIN!"
28.523457182818284°N, 80.683094159265358°W	EITHER YOU'RE HANDING OUT RAW FLOATING POINT VARIABLES, OR YOU'VE BUILT A DATABASE TO TRACK INDIVIDUAL ATOMS. IN EITHER CASE, PLEASE STOP.

Coordinate Precision from xkcd <https://xkcd.com/2170/>



## Names for people, and in this case places

- Country names as example
  - check multiple values
  - check for abbreviations, etc.



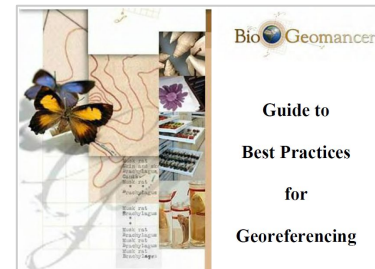




# Time travel to love your past self

*use data standards*

- Darwin Core (dwc)
- Ecological Metadata Language (EML)
- Best practices guide to georeferencing
- QGIS lesson for visualizing data to before analysis
- ...ask about more tools and resources at the #datahelpdesk in the exhibit hall



# Collections Biodiversity Data - *expected and emerging uses*

## Important Human Uses

- Evolutionary medicine,
- Disease discovery, tracking, and treatment
- Food security,
- Biodiversity conservation and sustainability,
- Computation,
- Design,
- Evolution and justice,
- Development of new types of biodiversity theories that accommodate newly emerging data streams.

## Emerging Research Angles

- Supplementing existing datasets with digital layers to enhance niche and species distribution modeling,
- Use of 3D/CT data for generating and testing new hypotheses,
- Implementation of convolutional neural networks (CNN) and deep learning in the analysis of image,
- Data for taxonomic determination and specimen curation,
- Delineation of traits in specimen images,
- Determination and identification to genus or species from, sediment-deposited pollen grains.

Nelson G, Ellis S. 2018. The history and impact of digitization and digital data mobilization on biodiversity research. *Philosophical Transactions of the Royal Society B* Volume 374 Issue 1763 <https://doi.org/10.1098/rstb.2017.0391>

# ESASUSSEE Data Help Desk

## *Who we are and how to find us (in the exhibit hall)*

Amber Budden, @aebudden, @DataONEorg, aebudden@nceas.ucsb.edu

Christina Alba, @botanic, christina.alba@botanicgardens.org

Christine Laney, @cmlaney, @NEON\_sci, claney@battelleecology.org

Deborah Paul, @idbdeb, @idigbio, dpaul@fsu.edu

Dmitry Schigel, @dschigel, @GBIF, dschigel@gbif.org

Jeanette Clark, @sjeanetteclark, @ArticDataCtr, jclark@nceas.ucsb.edu

Karl Benedict, @kbene, @ESIPfed, kbene@unm.edu, president@esipfed.org

Kelsey Yule, @kmyule, biorepo.neonscience.org, kmyule@asu.edu

Kristin Vanderbilt, @vanderbik, @EDIgotdata, krvander@fiu.edu

Kyle Copas, @kylecopas, @GBIF, kcopas@gbif.org

Laura Brenskelle, @lbrensk, @idigbio, lbrensk@ufl.edu

Margaret O'Brien @obrienmobb, @EDIgotdata, margaret.obrien@ucsb.edu

Megan Jones, @MeganAHJones, @NEON\_sci, mjones01@battelleecology.org

Nico Franz, @taxonbytes, biorepo.neonscience.org, nico.franz@asu.edu

Rebekah Wallace, www.eddmaps.org, bekahwal@uga.edu

Stevan Earl, @StevanEarl, stevan.earl@asu.edu

William Michener, @DataONEorg, william.michener@gmail.com



Amber



Deborah



Karl



Kyle



Megan



Stevan

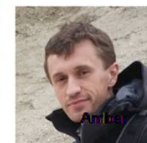


William

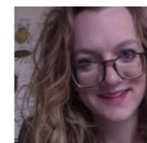
## ESASUSSEE 2019 Data Help Desk



Christina



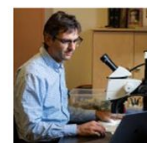
Dmitry



Kelsey



Laura



Nico



Christine



Jeanette



Kristin



Margaret



Rebekah