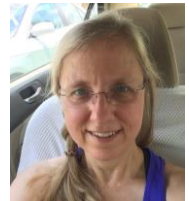![iDigBio - Integrated Digitized Biocollections]

# Biological collections data:
Best practices and trends for standards, digitization, and biodiversity informatics literacy for research use of collections data

**Deborah Paul**, Florida State University, iDigBio
Katja Seltmann, Cheadle Center for Biodiversity and Ecological Restoration (CCBER) #IslandBiology2016 University of the Azores at Angra do Heroísmo, Terceira Island, Azores, Portugal 19 July 2015
twitter @idbdeb @irene_moon

# An overview

- The need for high-quality data

- Digitizing collections and georeferencing

- Identifiers required

- Data sharing standards

- Researchers supplying, using data from aggregators

- Data gaps (Shelley)

- Biodiversity Informatics skills and literacy

# https://www.idigbio.org
# @iDigBio

# What Happens in an **Internet Minute?**

639,800 GB of global IP data transferred

**135** Botnet infections

**6** New Wikipedia articles published

**1,300** New mobile users

20 New victims of identity theft

**204 million** Emails sent

**47,000** App downloads

**$83,000** In sales

**100+** New Linkedin accounts

61,141 Hours of music

20 million Photo views

3,000 Photo uploads

**320+** New Twitter accounts

**100,000** New tweets

**277,000** Logins

**6 million** Facebook views

**2+ million** Search queries

**30** Hours of video uploaded

**1.3 million** Video views

## And **Future Growth** is **Staggering**

**Today**, the number of **networked devices** = the global population

By **2015**, the number of **networked devices** = **2x** the global population

In **2015**, it would take you **5 years** [IP] to view all video crossing IP networks each **second**
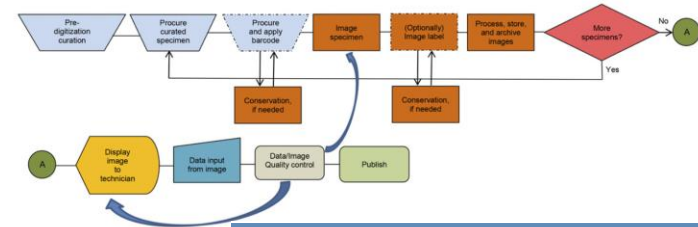
(intel)

# The Research Data Pipeline



Data > Knowledge > Application

# Digitization & Imaging

- Workflows, Protocols, ![ZooKeys - A peer-reviewed open-access journal. Launched to accelerate biodiversity research]
  - Best practices
- Prioritization trend
  - Research "digitization on demand"
- Curation
  - Physical and digital collections
- Working groups, webinars, publications, ...

**Fossil Insect Collaborative**

**Fossil Marine Invertebrates (EPICC)**

**Great Lakes Invasives**

**InvertEBase**

**InvertNet**

**Lichens & Bryophytes**

**Macroalgal Consortium**

**Macrofungi Consortium**

**Microfungi Consortium**

**NEVP**

**PALEONICHES**

**SCAN**

**SERNEC**

**Tri-Trophic**

**Vouchered Animal Communication Signals**

# Standards?
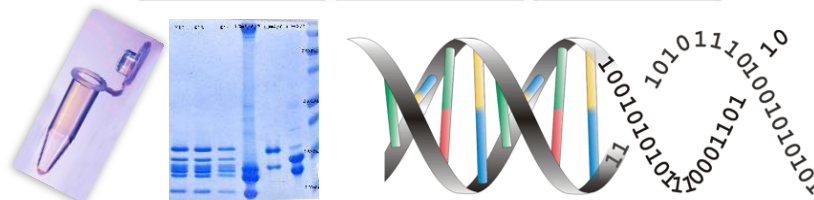
- All kinds of standards
- Data sharing standards

HTML

Biodiversity Information Standards
TDWG

# What are some examples of standards used for sharing biodiversity data? Where do they come from?

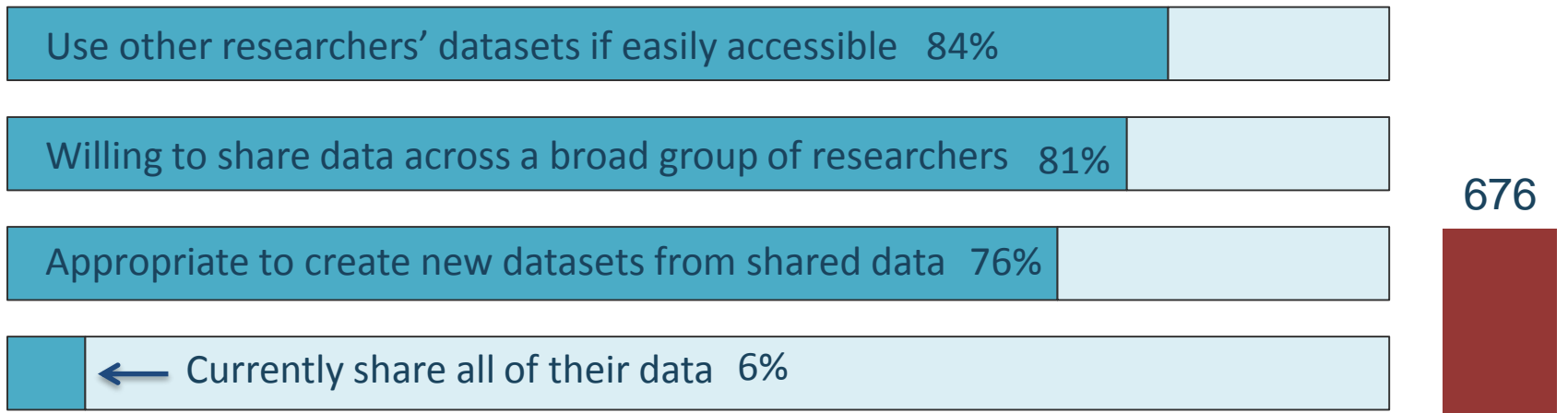| Data | Standards |
|---|---|
| specimens & observations | Darwin Core (DwC) |
| specimen & observation datasets | Ecological Metadata Language (EML) |
| media | Audubon Media Core |
| derivatives | Material Sample Core and GGBN Extensions, … |

WHO, WHAT, WHERE, WHEN

What's in the dataset?

# Scientists want to share data

Use other researchers' datasets if easily accessible   84%

Willing to share data across a broad group of researchers   81%

Appropriate to create new datasets from shared data   76%

← Currently share all of their data   6%

## Metadata standards

| DIF | DwC | DC | EML | FGDC | Open GIS | ISO | My Lab | none |
|-----|-----|-----|-----|------|----------|-----|--------|------|
| 12 | 21 | 26 | 95 | 95 | 96 | 97 | 266 | 676 |

# Standards – Why use them?

- Extend and expand data life

- Enhance sharing

- Facilitate re-use

- Increase likelihood of new uses

- Make linked-data initiatives possible
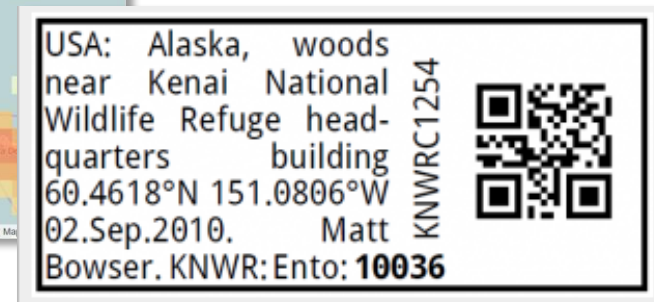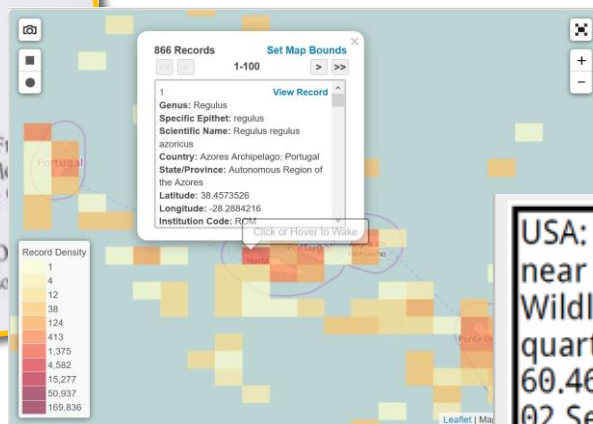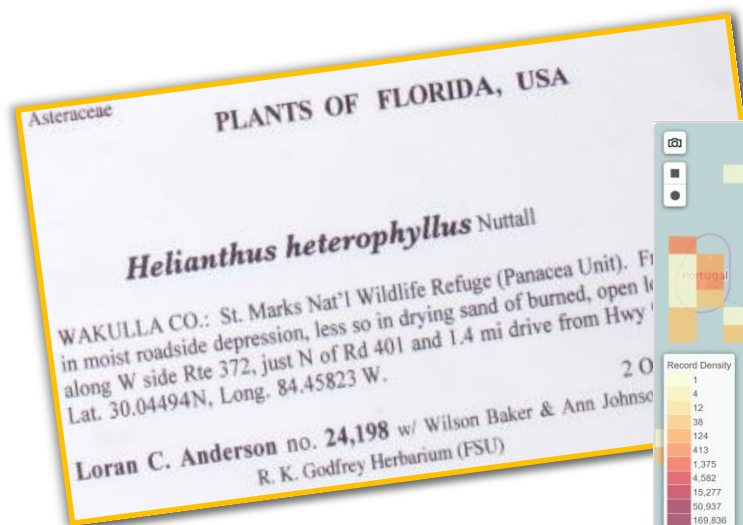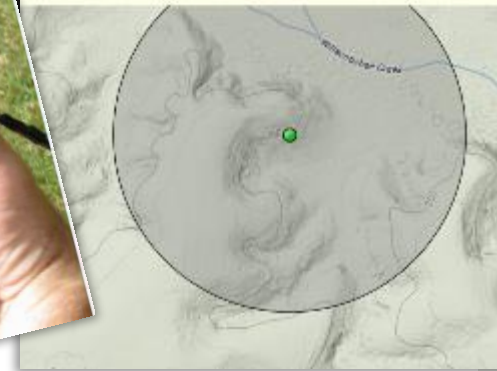
- Feedback, attribution

# Georeferencing

- What's your georeferencing workflow?
  - Legacy or New data
  - Use best practices
- Darwin Core (dwc)

# Darwin Core terms
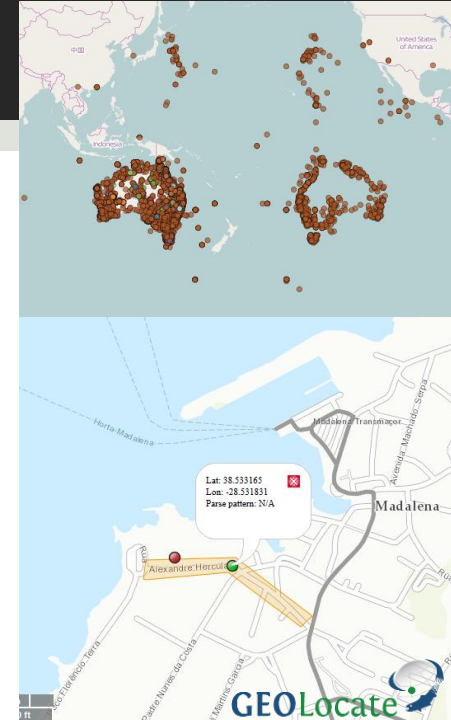## Location and Geological Context



locationID | higherGeographyID | higherGeography | continent | waterBody | islandGroup | island | country | countryCode | stateProvince | county | municipality | locality | verbatimLocality | minimumElevationInMeters | maximumElevationInMeters |minimumDepthInMeters | maximumDepthInMeters | verbatimDepth | minimumDistanceAboveSurfaceInMeters | maximumDistanceAboveSurfaceInMeters | locationAccordingTo | locationRemarks | decimalLatitude | decimalLongitude | geodeticDatum | coordinateUncertaintyInMeters | coordinatePrecision | georeferencedBy | georeferencedDate | georeferenceProtocol | georeferenceSources | georeferenceVerificationStatus | georeferenceRemarks

geologicalContextID | earliestEonOrLowestEonothem | latestEonOrHighestEonothem | earliestEraOrLowestErathem | latestEraOrHighestErathem | earliestPeriodOrLowestSystem | latestPeriodOrHighestSystem | earliestEpochOrLowestSeries | latestEpochOrHighestSeries | earliestAgeOrLowestStage | latestAgeOrHighestStage | lowestBiostratigraphicZone | highestBiostratigraphicZone | lithostratigraphicTerms | group | formation | member | bed

# Georeferencing Tools, Materials, Workflows



- http://georeferencing.org/
- GEOLocate
- GWG at iDigBio
  - Listserve, expertise
- Workshop materials, videos, powerpoints

- Good Localities Bad Localities: a Guide for your Field Notebook

iDigBio 2015 Field-To-Database Workshop

Field Number: US15-_____ Date:_____ - March-2015 Start Time:_____ End Time:_____

Country: UNITED STATES State: FLORIDA County: ALACHUA Lat:_____N Lon:_____W

Elev.:_____m GPS Error: +/-_____m Extent: _____m Datum: WGS84 Site Photo: yes / no
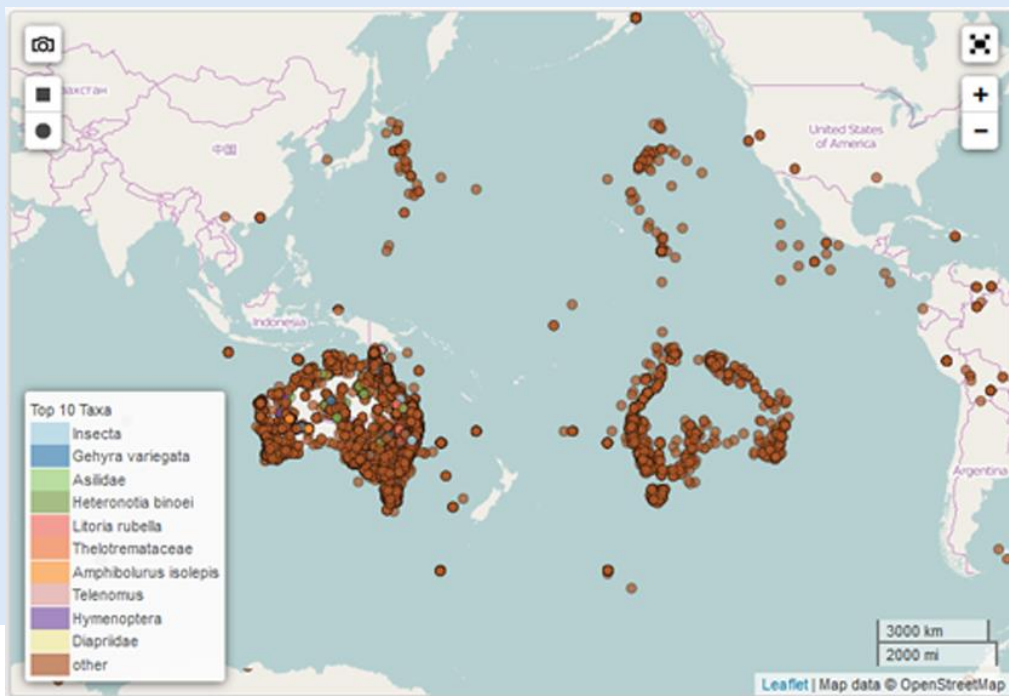
Specific Locality: Gainesville, University of Florida Campus:_____

General Site Description/Overview:

# iDigBio Data Quality (DQ) Flags enhance Digitization and Research Workflows

Example: spot and fix georeferencing issues.



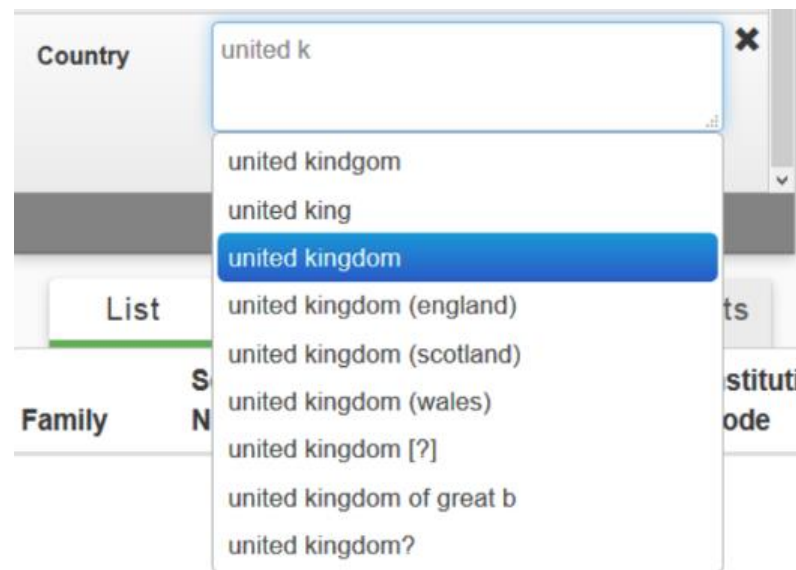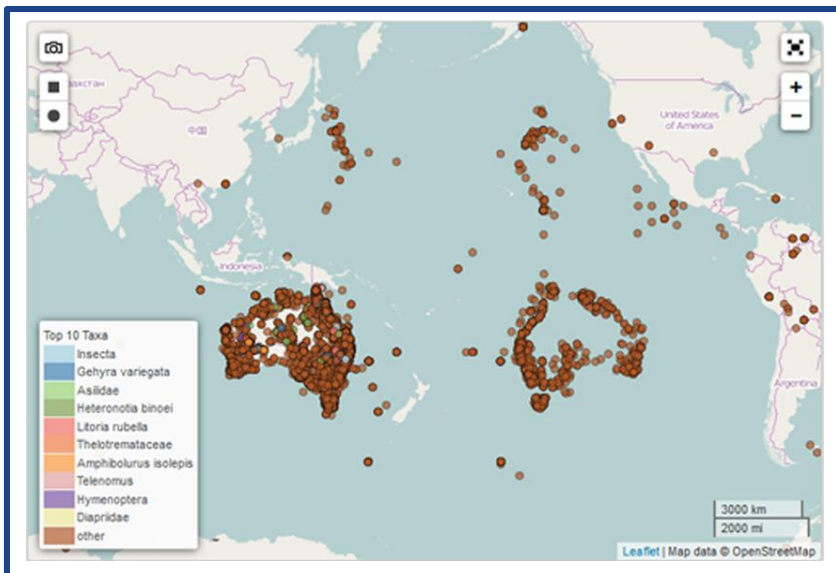| Flag |
|------|
| idigbio_isocountrycode_added ⓘ |
| dwc_continent_added ⓘ |
| dwc_country_replaced ⓘ |
| geopoint_datum_missing ⓘ |
| dwc_class_replaced ⓘ |
| dwc_phylum_replaced ⓘ |
| dwc_order_replaced ⓘ |
| geopoint_low_precision ⓘ |
| rev_geocode_eez ⓘ |
| dwc_stateprovince_replaced ⓘ |
| rev_geocode_mismatch ⓘ |
| dwc_order_added ⓘ |
| datecollected_bounds ⓘ |
| dwc_class_added ⓘ |
| dwc_kingdom_added ⓘ |
| dwc_phylum_added ⓘ |
| dwc_country_added ⓘ |
| rev_geocode_corrected ⓘ |
| rev_geocode_lon_sign ⓘ |

# Data quality: *an issue at many levels*



Hannah Frost
@feefifofannah
Following

From a @HydraInABox interview: "People will put anything and their dog in the date field. It's absolutely astonishing."



Top 10 Taxa
- Insecta
- Gehyra variegata
- Asilidae
- Heteronotia binoei
- Litoria rubella
- Thelotremataceae
- Amphibolurus isolepis
- Telenomus
- Hymenoptera
- Diapriidae
- other



Country: united k

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

196 Countries in the world, but 1100 distinct values in the country field

# Identifying specimens

- Determined by: Roman S. Wielgus
- Collected by: Roman S. Weilgus; Dalie Wielgus
- Date identified: 1968
- Sex: female
- Decimal Latitude: 34.7198786
- Georeferenced by: Jean-Batiste Quirino
- Date Collected: 1968-06-02
- Collector number: rsw1256
- Catalog number: ASUHIC0080642
- Institution code: ASU
- Collection code: ASUHIC
- Occurrence ID (GUID): afb73b66-ad30-40ea-bb86-8522448ad044
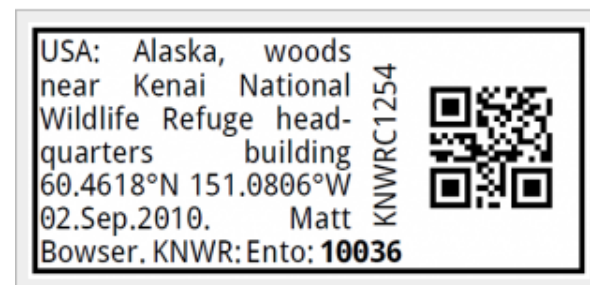
*Adelpha bredowii eulalia* (Geyer, 1837)

# Identifying specimens globally

- Specimens (and other bits)
  - need Globally Unique Identifiers (GUIDs)
- Plan ahead – Shelley's example and FIMS
- Journals using identifiers – Pensoft
- Store identifiers and **no re-use**

- Why use them?
  - Support linked data
  - Data reusability
  - Data discovery

Examples

- e20d9229-4dfc-41c8-be05-afa7b40245a3
- urn:catalog:CAS:IZ:25668
- urn:uuid:6fdb535e-d43d-40c7-a387-e60c07cec5d1
- urn:lsid:biosci.ohio-state.edu:osuc_occurrences:OSUM_Birds_B_13589
- http://ucjeps.berkeley.edu/cgi-bin/new_detail.pl?RSA460997
- occurrenceID: afb73b66-ad30-40ea-bb86-8522448ad044

*Adelpha bredowii eulalia* (Geyer, 1837)



USA: Alaska, woods near Kenai National Wildlife Refuge head-quarters building 60.4618°N 151.0806°W 02.Sep.2010. Matt Bowser. KNWR: Ento: **10036**

KNWRC1254

SEMC0993403
KUNHM-ENT

AM_ENT
AMNH_PBI 00388325

Online import of occurrence records directly into a manuscript!

GBIF | BOLD SYSTEMS | IDigBio | PlutoF

Occurrence 1
Occurrence 2
...
Occurrence n

ARPHA WRITING TOOL

**Edit Materials**

You may place multiple ID's separated by " | " here

e20d9229-4dfc-41c8-be05-afa7b40245a3 | 6fdb535e-   Add

BOLD record ID   (example: ACRJP618-11 | ACRJP619-11)

BOLD BIN   (example: BOLD:AAA5125 | BOLD:AAA5126)

GBIF via Occurrence ID   (example: urn:catalog:HYO:ENT:B1367540 | 4b7b4bb4-0db7-4592-b3f9-1b15b6235360)

GBIF ID   (example: 1061574007 | 240843113)

iDigBio UUID   (example: 1db58713-1c7f-4838-802d-be784e444c4a | d957ac64-ce51-4d40-801e-670b345aa7b6)

PlutoF record ID   (example: FM178343 | EU343855)

PlutoF SH ID   (example: 10.15156/CH487435.07FU | SH487425.07FU)

Save   Close

Taxonomic manuscript

submission

Biodiversity Data Journal   http://bdj.pensoft.net
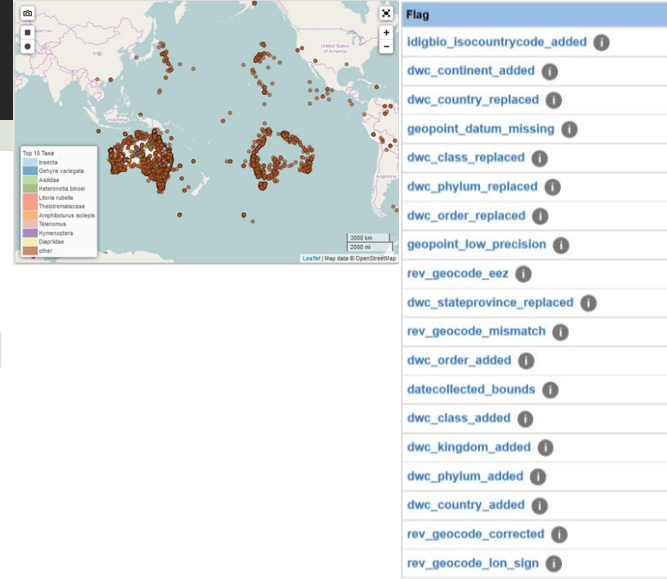
REPOSITORIES

ARPHA WRITING TOOL

MANUSCRIPT

PUBLISHED ARTICLE

# Natural History Museum
# Specimen Data Aggregators and You

- iDigBio, GBIF, VertNet, ALA, ...
  - providing data quality information
- iDigBio dataset downloads
  - **original and enhanced data!**
  - **dataset citation** (Matt)
- when using museum specimen data
  - use / cite the guids (dwc:occurrenceID) provided
- share your data widely
- ask collections you work with if they are sharing their data
  - easy to do
    http://www.idigbio.org/wiki/index.php/Data_Ingestion_Guidance

Once digitized, why share and publish your data? Can you share it in more than one place? Yes!

- power of aggregation
- accessibility (researchers, funders, collaborators)
- discoverability
- enhancement (scripts)
- tools and methods
- data quality checks
- off site copy
- visibility
- new uses for data

Researchers sentiments about data:
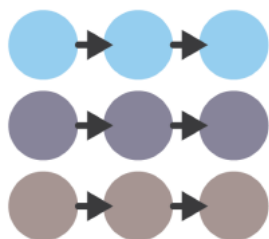
(BEACON, SESYNC, NESCent, iPlant, iDigBio)

- I usually manage data in Excel and it's terrible.
- I'm organizing GIS data and it's becoming a nightmare.
- I'm having a hard time analyzing microarray, SNP or multivariate data with Excel and Access.
- I want to use public data.
- I work with faculty at undergrad institutions and want to teach data practices, but I need to learn it myself first.
- I'm interested in going in to industry and companies are asking for data analysis experience.
- I'm trying to reboot my lab's workflow to manage data and analysis in a more sustainable way.
- I'm re-entering data over and over again by hand; there must be a better way.
- I have overwhelming amounts of data.
- I'm tired of feeling out of my depth on computation and want to increase my confidence.

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

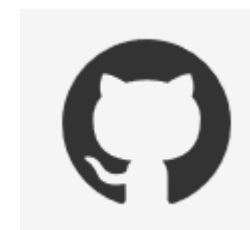# People seeking skills and efforts to scale up to meet needs



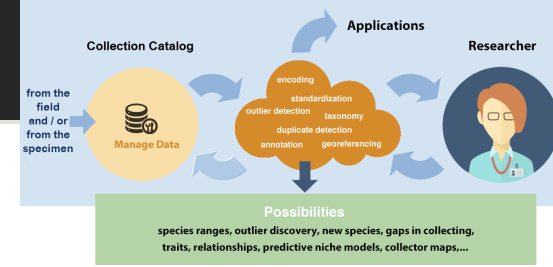Reproducible Science Curriculum

## Library Carpentry

Software Skills Training for Librarians

http://librarycarpentry.github.io/city-november-2015/

@GitHub

# Some Biodiversity Informatics short courses at iDigBio

# Ten Simple Rules for the Care and Feeding of Scientific Data

Alyssa Goodman, Alberto Pepe ✉, Alexander W. Blocker, Christine L. Borgman, Kyle Cranmer, Merce Crosas, Rosanne Di Stefano, Yolanda Gil, Paul Groth, Margaret Hedstrom, David W. Hogg, Vinay Kashyap, Ashish Mahabal, Aneta Siemiginowska, Aleksandra Slavkovic

- Love your data, and help others love it, too
- Share your data online, with a permanent identifier
- Conduct science with a particular level of reuse in mind
- Publish workflow as context
- Link your data to your publications as often as possible
- Publish your code (even the small bits)
- State how you want to get credit
- Foster and use data repositories
- Reward colleagues who share their data properly
- Be a booster for data science

# And now, more about collections data fit-for-research use,

facebook.com/iDigBio

twitter.com/iDigBio

vimeo.com/idigbio

idigbio.org/rss-feed.xml

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

**www.idigbio.org**

# TASK GROUPS ON FITNESS FOR USE:
## INVASIVE ALIEN SPECIES AND DNA EVIDENCE

**2016 group** on invasive alien species:
**Melodie McGeoch** + 5 experts
- meetings, teleconferences and survey
- report in November 2016

\*  \*  \*

Exploring data use and data publishing needs of molecular biodiversity research: **DNA** survey, GGBN conference and TDWG

Interested? Contact Dmitry Schigel
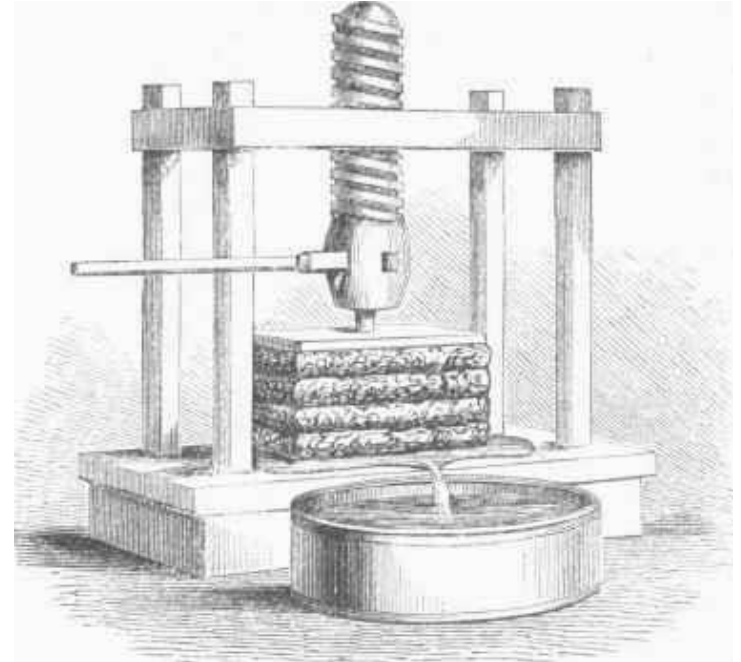dschigel@gbif.org







GBIF

# DATA QUALITY, ANALYSIS AND USE

Promoting data quality culture:
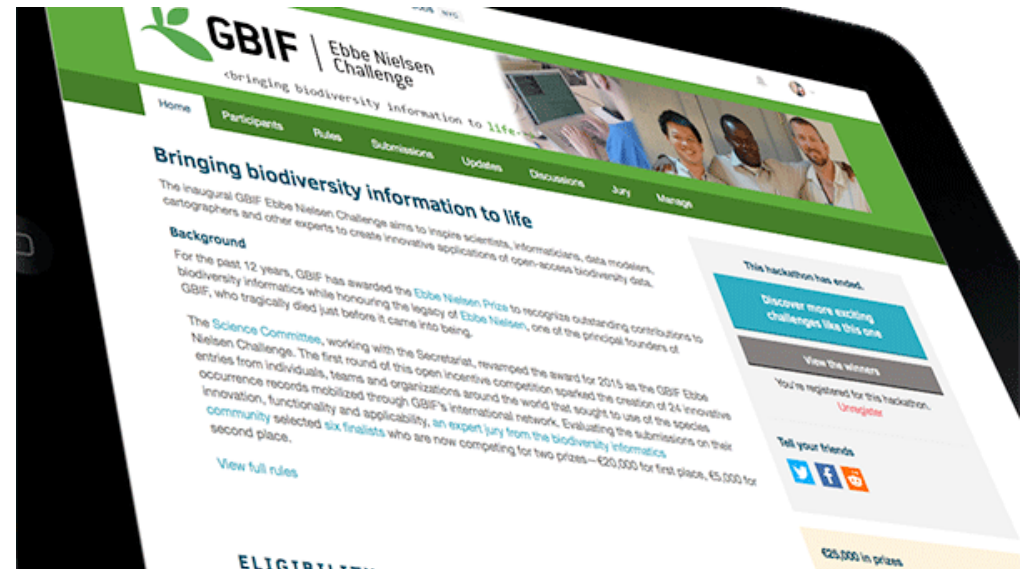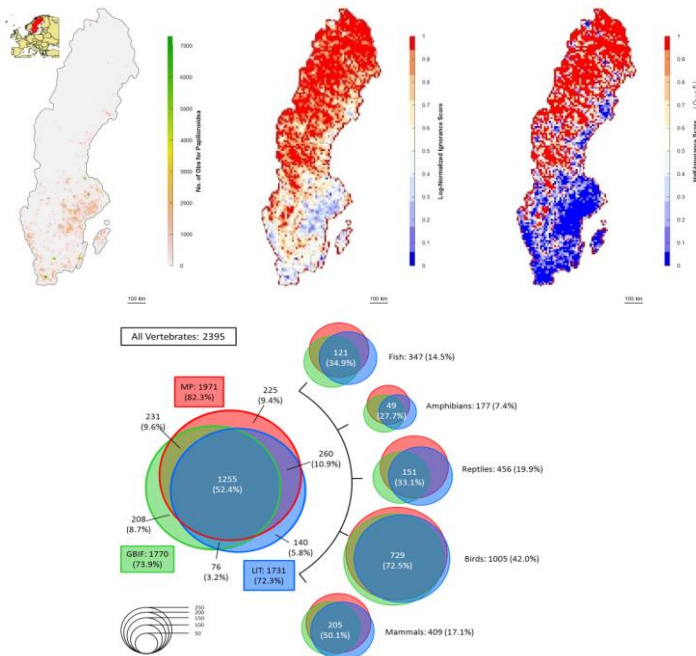talks, teaching, [publications]


TDWG / GBIF
data quality interest group, chairs:
**Antonio Saraiva** & **Arthur Chapman**


Join the work at the **GBIF Community Site!**

Biodiversity
Information
Standards
TDWG

GBIF

# GAPS AND EBBE NIELSEN CHALLENGE 2016



Challenge opens on 29 July 2016
First prize **€20,000**, second prize **€5,000**
Winners announced on 26 October 2016

**Dealing with gaps is the theme of the Ebbe Nielsen Challenge 2016**

http://www.gbif.org/about/awards

# Public Participation in Digitization

**Transcription Blitz with FL Native Plant Society**

**Herbarium Imaging Blitz**