

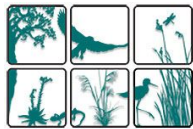


# Assuring the quality of your data: A natural history collection community perspective

Deborah Paul, Katja Seltmann, Laura Russell, David Bloom

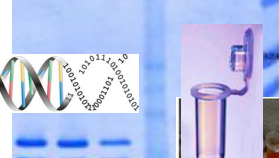
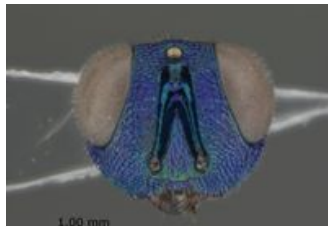
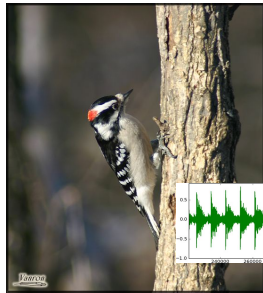


@idbdeb @iDigBio @VertNetOrg



Tuesday, January 12, 2016  
9 am Pacific / 10 am Mountain /  
11 am Central / 12 noon Eastern





Occurrence Data | Determination History | Images | Admin

Collector Info  
 Catalog Number MSC-B-000001 | Other Numbers | Collector F.H. Bormann | Number 1155 | Date 1953-07-20 | Dupe?  | Auto search

Associated Collectors  
 J.E. Cantlon, A.L. Reubek

Latest Identification  
 Scientific Name: Abietinella abietina | Author: (Hedw.) Fleisch.  
 D Qualifier: | Family: Thuidiaceae  
 Identified By: | Date Identified:

Locality  
 Country: | State/Province: | County: | Municipality: | Locality:

Locality Security

Latitude: | Longitude: | Uncertainty (meters): | Datum: | Elevation in Meters: | Verbatim Elevation: | Tools: | ft. |

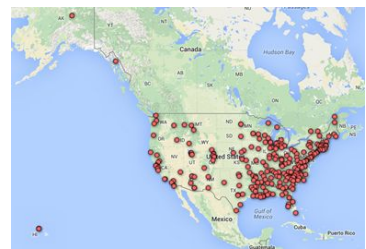
Misc  
 Habitat: | Substrate: | Associated Taxa: | Description: | Notes:

Label Processing

PLANTS OF ARCTIC ALASKA  
*Abietinella abietina* (Hedw.) C.H. det. by H. Crum  
 Colluvial willow on Red Mt. S-facing slope. 69° 25' N lat.  
 F.H. Bormann, J.E. Cantlon, A.L. Reubek NO. 1155 20 July 1953  
 Beat-Darlington Herbarium MICHIGAN STATE UNIVERSITY

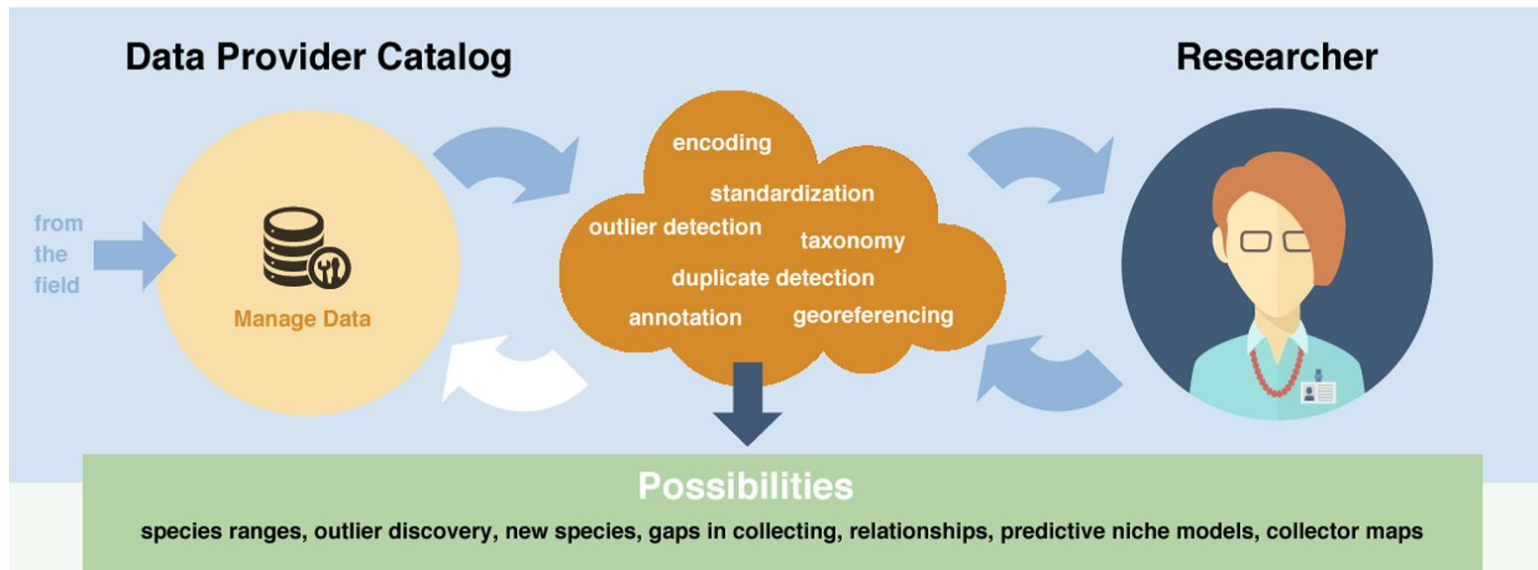
Options  
 OCR whole image  
 OCR w/analysis

PLANTS OF ARCTIC ALASKA  
*Abietinella abietina* (Hedw.) det. by H. Crum  
 Colluvial willow on Red Mt. S-facing slope. 69° 25' N lat.  
 F.H. Bormann, J.E. Cantlon, A.L. Reubek NO. 1155 20 July 1953  
 Beat-Darlington Herbarium MICHIGAN STATE UNIVERSITY



# Our data producers & data users

**Data Quality** starts here, before collection of specimens and field data



Naturalists, Field biologists,  
Nature explorers, Research institutions, Citizen  
scientists, Curators

Ecologists, biogeographers; Analysts, modellers;  
Conservation planners;  
Nature managers; Policy managers;  
Funding agencies; Industry; ?, ...

# Data, data types, data standards

Darwin Core (DwC)

Audubon Media (AC)

Ecological Metadata Language (EML)

Global Genome Biodiversity Network (GGBN)

, . . .

## DwC Categories

- Record Level (19)
- Occurrence (19)
- Organism (7)
- Material Sample (1)
- Event (15)
- Location (44)
- Geological Context (18)
- Identification (8)
- Taxon (33)

## Record Level

**dcterms:type** | **dcterms:modified** |  
**dcterms:language** | **dcterms:license** |  
**dcterms:rightsHolder** | **dcterms:**  
**accessRights** | **dcterms:**  
**bibliographicCitation** | **dcterms:**  
**references** | **institutionID** | **collectionID**  
| **datasetID** | **institutionCode** |  
**collectionCode** | **datasetName** |  
**ownerInstitutionCode** | **basisOfRecord** |  
**informationWithheld** |  
**dataGeneralizations** | **dynamicProperties**

# Data Publication Fallacies (and Truths)

## **The Fallacy of Perfection:**

Data must be perfect before publication.

## The Truth:

Data will never be perfect in every aspect.

# Data Publication Fallacies (and Truths)

## **The Fallacy of Petrification:**

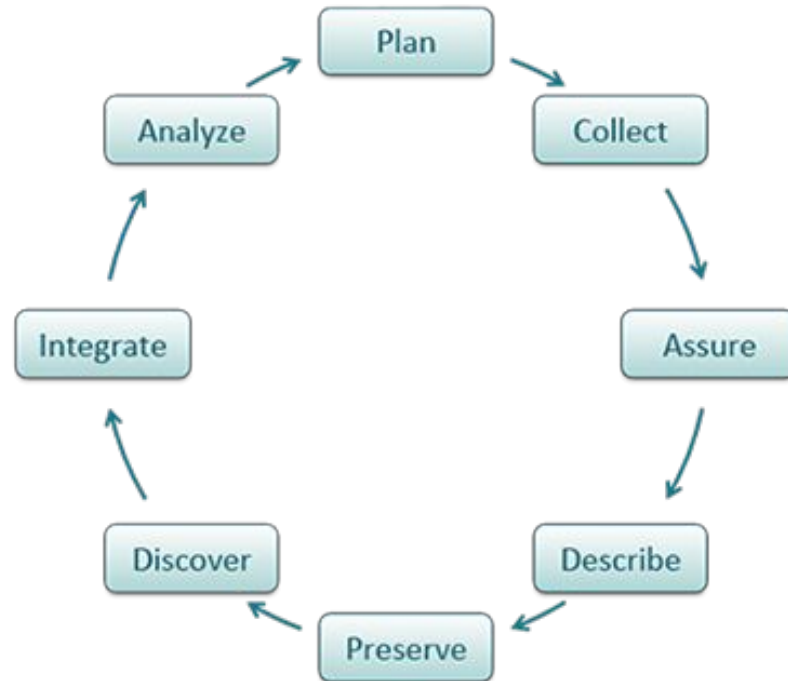
Data do not change (or don't need to change) once they are in a ledger or database.

## The Truth:

Data are dynamic and require regular curation.

# Data Publication Fallacies (and Truths)

## DataONE Data Life Cycle



# Data Publication Fallacies (and Truths)

## **The Fallacy of Fitness for Use:**

The fitness for use of data depends upon how and why the data were collected.

## The Truth:

Fitness depends upon the questions being asked (value is in the eye of the beholder).



# Challenges for data publication and quality

## **Individual:**

Education, Experience, Funding, Presence of support

## **Institutional:**

Inter/Intra-collection communication, Mixed technology, Legal limitations

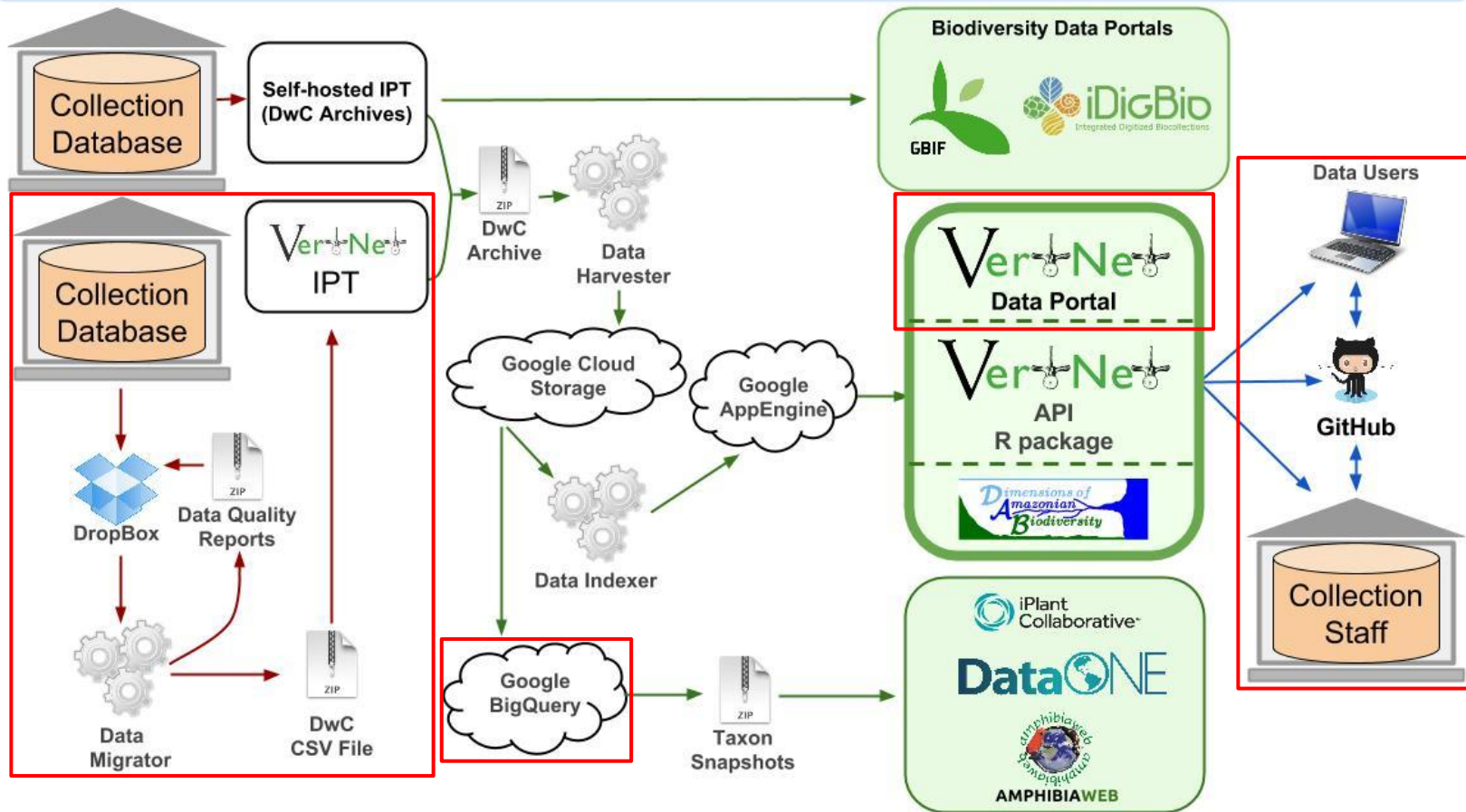
## **Structural:**

Interoperability, Related/Derivative sources

# Data assurance practices - VertNet

- Data publishing
  - VertNet Toolkit
    - Data migrator
    - Data quality reports
- Data acquisition and processing
  - Google BigQuery
- Data access - VertNet portal
  - Spatial quality
  - Flag data issues
- Feedback
  - Issue tracking and feedback via GitHub

Original Data Sources  Data Acquisition & Processing  Access  Feedback & Stats



General toolkit for working with VertNet data. We call these data "migrators." Once customized to an original data source, it converts that original data into Darwin Core ready for upload to an Integrated Publishing Toolkit (IPT) resource. — Edit

91 commits

1 branch

1 release

2 contributors

Branch: develop

New pull request

New file

Find file

HTTPS

https://github.com/VertNet/



Download

tucotuco Amended table name LegacyCoordinates. Added template AudubonCoreExtens... Latest commit 72840f6 27 days ago

ForIPT	Amended table name LegacyCoordinates. Added template AudubonCoreExtens...	27 days ago
bkp	Cleaning up base branch to latest revision.	2 years ago
reports	Added processing for duplicate occurrenceIDs.	5 months ago
source	Major upgrade to the migrator. No longer supports DiGIR provider output.	a year ago
templates	Amended table name LegacyCoordinates. Added template AudubonCoreExtens...	27 days ago
workspace	Adding in new migrator scripts for non-verts and added executable pat...	7 months ago
1 - FullDataPreparation.bat	Major upgrade to the migrator. No longer supports DiGIR provider output.	a year ago
1a - RunMigrators.bat	Added processing for inverts. Updated coordinate processing to accomm...	4 months ago
1b - CleanMigratedTables.bat	Added templates and processing for Entomology and Fungi data sets bas...	8 months ago
1c - RunAggregators.bat	Adding in new migrator scripts for non-verts and added executable pat...	7 months ago
BlankLineIssues.awk	Syncing to latest migrator.	3 years ago

## DATA QUALITY

We believe in publishing the highest quality and most complete data possible. After we talk with you about your data, we'll review your data for possible data improvements, such as duplicate catalog numbers, indeterminate and non-standard geography, inconsistent taxonomy, and terms not compliant with **Darwin Core**. Then we'll provide you with a detailed **report** so that you can update your database locally. Only publishers who host their data on the **VertNet IPT** benefit from these services automatically. If you host your own data set, you can request this service to make your data more complete.

```
1 SELECT class FROM [dumps.vertnet_latest] where class is not null group by class
```

**RUN QUERY**

Save Query

Save View

Format Query

Show Options

Query complete (4.1s)

Results	Explanation
Row	class
155	Insecta
156	LEPOSPONDYL
157	LISSAMPHIBIAN
158	Leptocardii
159	Liliopsida
160	Lycopodiopsida
161	MAMM INSECT TRACE
162	MAMM INVERT
163	MAMM PLANT
164	MAMM TRACE
165	MAMMALIA

Warning: Some validations could not be performed. Check below.

## Data completeness

Are coordinates present?	✓	Yes
Is the country value present?	✓	Yes
Are both coordinates 0 (zero)?	✓	No
Do coordinates have three or more decimal figures?	✓	Yes
Do coordinates have datum?	✓	Yes

## Data inconsistencies

Are coordinates within specified country? <sup>1</sup>	✓	Yes
Distance outside of specified country (in degrees) <sup>1</sup>	✓	0
Distance outside of species range map (in degrees) <sup>1</sup>	⊘	Could not be assessed

## Data Errors

Is latitude between 90 and -90?	✓	Yes
Is longitude between 180 and -180?	✓	Yes
Are coordinates transposed? <sup>1</sup>	✓	No
Is latitude hemisphere correct? <sup>1</sup>	✓	Yes
Is longitude hemisphere correct? <sup>1</sup>	✓	Yes

<sup>1</sup>Assessed with [Map Of Life](#) validation tools

## Submit data issue

Help data publishers track data issues using [GitHub!](#)

Eleutherodactylus gryllus - Specimen in Aguada?

This specimen is listed as in [Aguada](#), way outside its known range. Perhaps a wrong location or misidentification of the species?

Submit data issue

Easy for users to submit. Users do need a free GitHub account.

Easy for data publishers to receive and manage feedback about their published data.

## [CM Herps 36065] Eleutherodactylus gryllus - Specimen in Aguada? #19

[Open](#) | [ljvillanueva](#) opened this issue 8 days ago · 0 comments



[ljvillanueva](#) commented 8 days ago

This specimen is listed as in [Aguada](#), way outside its known range. Perhaps a wrong location or misidentification of the species?

You can [view the original detail page](#) on VertNet. Here are the original record contents:

Term	Value
Modified	2015-10-11
Language	en
AccessRights	<a href="http://vertnet.org/resources/norms.html">http://vertnet.org/resources/norms.html</a>
References	<a href="http://portal.vertnet.org/o/cm/herps?id=36065">http://portal.vertnet.org/o/cm/herps?id=36065</a>

# Data assurance practices



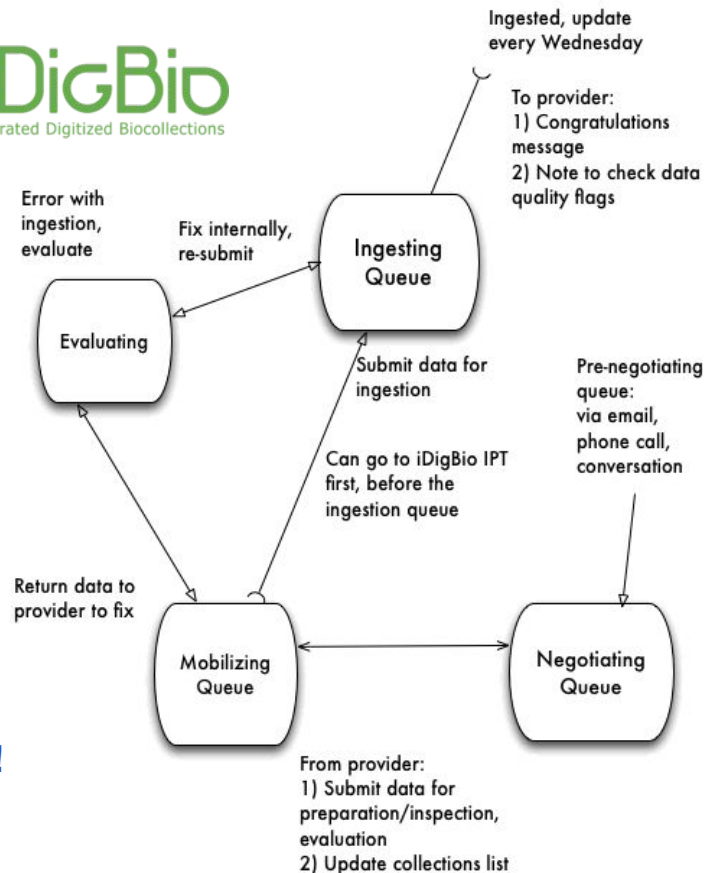
Data Publication - an overview of the process data goes through before ingestion at iDigBio

1. Negotiating
  - a. Evaluation (human, scripts, reports)
  - b. *iDigBio Data Quality (DQ) Flags* ↩
2. Mobilizing (human, scripts)
3. Ingestion

**a Perk of data sharing!**



[https://www.idigbio.org/wiki/index.php/Data\\_Ingestion\\_Guidance](https://www.idigbio.org/wiki/index.php/Data_Ingestion_Guidance)



These queue name roughly line up with the ingestion report here:  
[https://www.idigbio.org/wiki/index.php/Data\\_Ingestion\\_Report](https://www.idigbio.org/wiki/index.php/Data_Ingestion_Report)

Graphic by Joanna McCaffrey, iDigBio Biodiversity Informatics Manager



# Search and download using the iDigBio DQ Flags

## DQ Flags

- enhance and improve data
- enhance discoverability
- visualization aid
- transparency builds community trust
- facilitate potential development of automated updates by provider

The screenshot displays the iDigBio search portal. At the top, a green navigation bar includes links for 'iDigBio Home', 'Portal Home', 'Search Records', 'Tutorial', 'Data', 'Research Tools', and 'Feedback', along with a user profile 'dpaul'. The main search area features a 'Search Records' section with a search box and filters for 'Must have image' and 'Must have map point'. Below this are tabs for 'Filters', 'Mapping', 'Sorting', and 'Download', with an arrow pointing to the 'Download' button. A 'Data Flags' section is highlighted with an arrow, containing a search box and checkboxes for 'Present' and 'Missing'. To the right, a map shows 'Record Density' with a color scale from 1 (yellow) to 255,991 (dark red). Below the map, a table lists search results with columns for Family, Scientific Name, Date Collected, Country, Institution Code, Basis of Record, Kingdom, Phylum, Class, Order, Higher Taxon, and Common Name. The total number of records is 48,308,211.

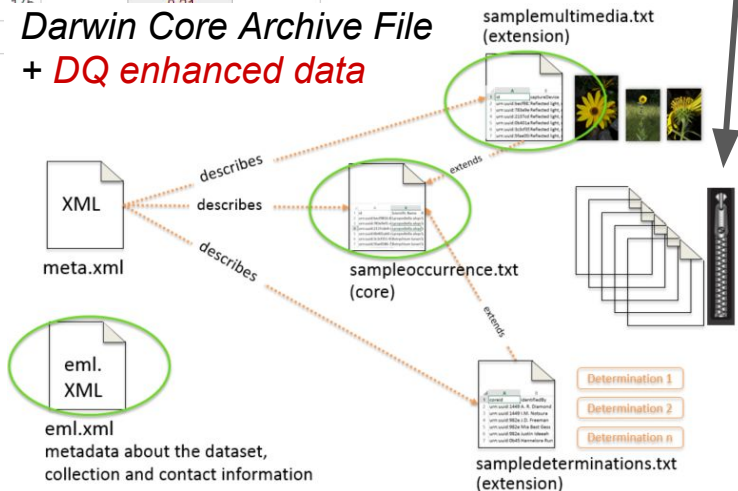
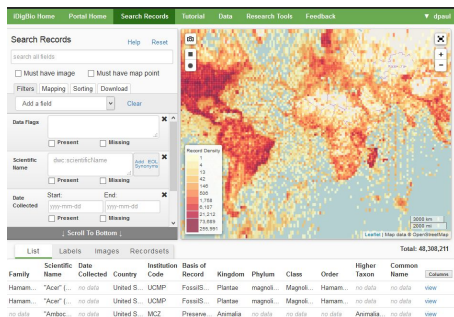
Family	Scientific Name	Date Collected	Country	Institution Code	Basis of Record	Kingdom	Phylum	Class	Order	Higher Taxon	Common Name	Columns
Hamam...	"Acer" (...)	no data	United S...	UCMP	FossilS...	Plantae	magnoli...	Magnoli...	Hamam...	no data	no data	view
Hamam...	"Acer" (...)	no data	United S...	UCMP	FossilS...	Plantae	magnoli...	Magnoli...	Hamam...	no data	no data	view
no data	"Amboc...	no data	United S...	MCZ	Preserve...	Animalia	no data	no data	no data	Animalia...	no data	view

# Data assurance practices



Flag	Records With This Flag	(%) Percent With This Flag
idigbio_isocountrycode_added	83287	99.74
dwc_continent_added	83258	99.705
dwc_phylum_added	82479	98.773
dwc_class_added	82471	98.763
dwc_order_added	82465	98.756
geopoint_datum_error	8223	9.847
geopoint_datum_missing	5137	6.152
rev_geocode_eez	691	0.828
rev_geocode_mismatch	249	0.298
rev_geocode_corrected	175	0.21
rev_geocode_lon_sign		
geopoint_low_precision		

Darwin Core Archive File  
+ *DQ enhanced data*



- Data Quality Flags *by recordset* at iDigBio
- Download Darwin Core Archive Files

- *raw data*
- *and a bonus file*
- *available to everyone*

- Annotations coming
- **feedback loop**
- **transparency**

# Contact the provider...

## Specimen Record

Animalia > Chordata > Amphibia > Anura > Hyperolidae

### *Afrixalus fulvovittatus*

From Museum of Comparative Zoology, Harvard University

Continent Africa  
Country Cameroon  
Locality Nkambe  
Latitude 6.5486166667  
Longitude 10.7599833333

Institution MCZ  
Code  
Collection Herp  
Code  
Catalog A 148086  
Number

### Contacts

**Name** MCZ Harvard University  
**Role** none  
**Email** none  
**Phone** none

**Name** Brendan Haley  
**Role** Senior Database Manager  
**Email** [bhaley@oeb.harvard.edu](mailto:bhaley@oeb.harvard.edu)  
**Phone** none



### Media



### Contents

Summary  
Map  
Media  
Attribution  
All Data

# Future DQ Work at iDigBio

- Full Taxonomic resolution against the GBIF Backbone, including: adding TaxonID (for all taxonomic levels) values, accepted names, and canonical names.
- Expanded geographic name/point validation/augmentation using GADM, GeoNames, or some similar place name database.
- Opening up the stage 1 corrections process to include externally provided corrections data sources. An API has been provisionally designed, but it hasn't undergone and testing or validation.
- Opening up the stage 1 corrections process to include annotations Since, from a technical perspective the process of applying an annotation to a record is the same as applying a general correction that only applies to specific single records or versions. This is even farther out than the last bullet, since it would need an entire workflow built around it to be effective.
- Harmonization of flag values, conditions and actions with other projects (such as this ALA/GBIF list: <http://bit.ly/evMJv5> )

# *NSF ADBC Digitization Thematic Collection Network*

Plants, Herbivores, and Parasitoids: A Model System for the study of Tri-Trophic Associations



Photo: S. Bauer



Illustration: W. J. Hooker



Photo: C. A. Johnson

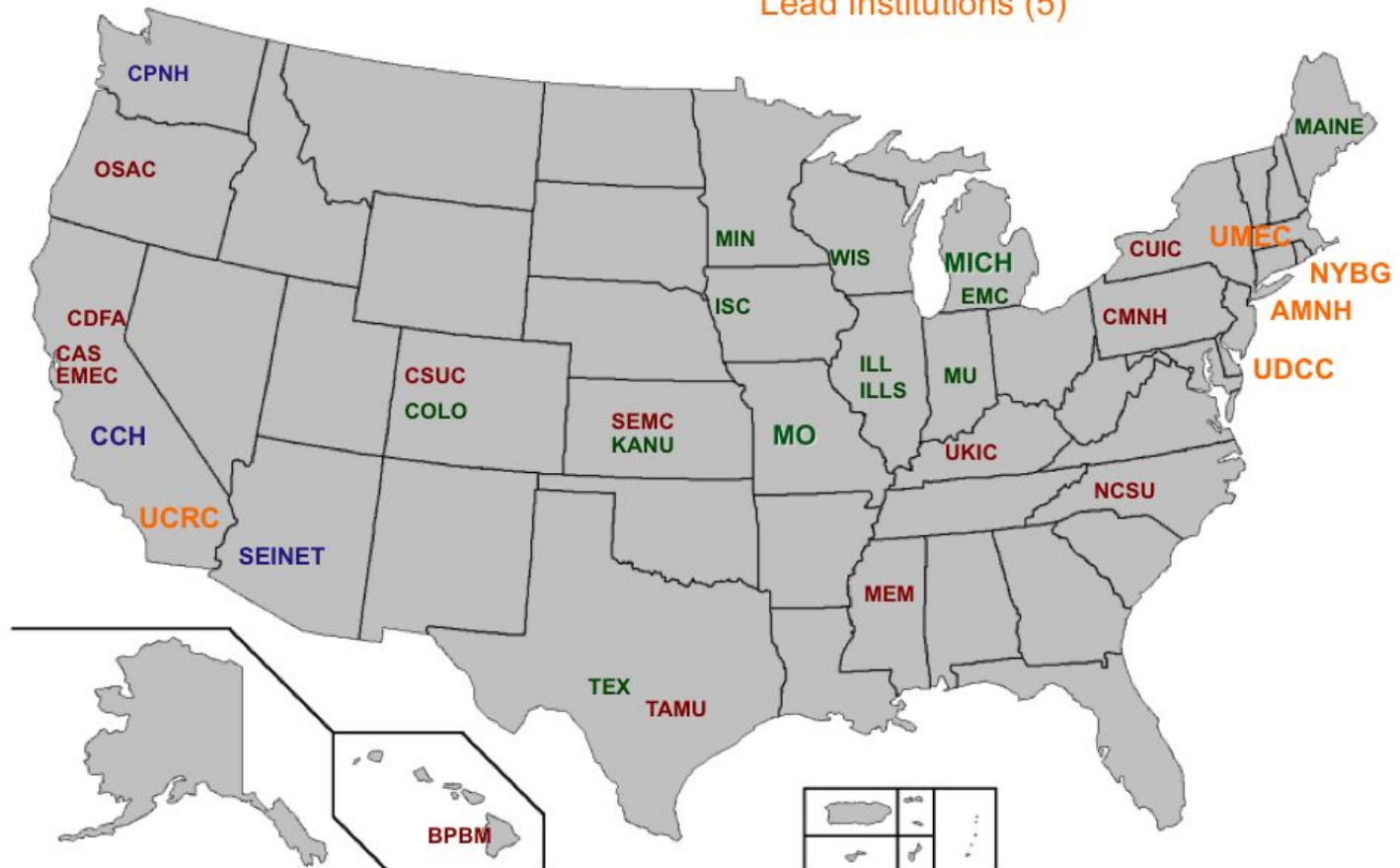
**3,645,975**  
new records in 4 years

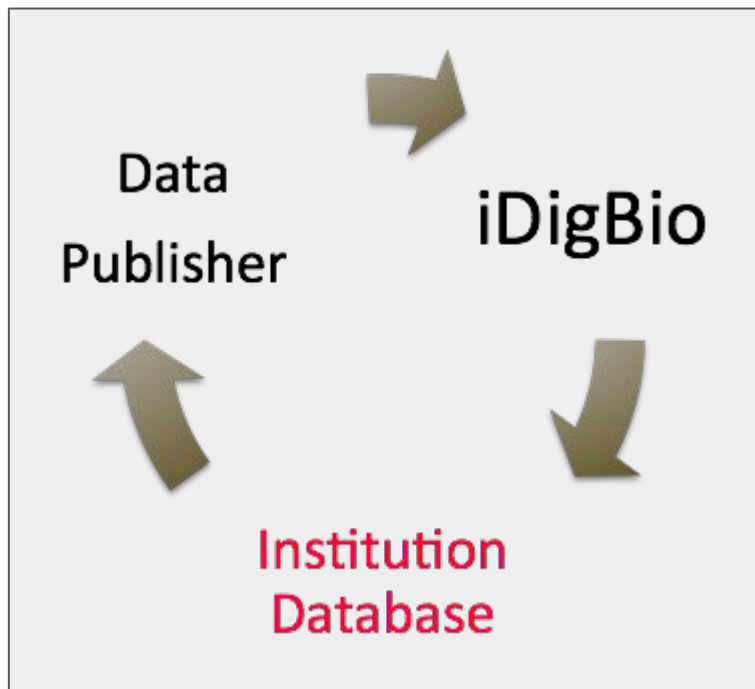
# TTD Network

Herbaria (14) Botanical Data Contributors (3)

Insect Collections (18)

Lead Institutions (5)





1. Collection manager commitment and training
2. QC tools (data entry, data evaluation)
3. **Research use of data**

# Acknowledgements and References

DRAFT FOR CONSULTATION: Best Practice Guide for Data Gap Analysis for biodiversity stakeholders [http://www.gbif.org/system/files\\_force/gbif\\_resource/resource-82566/DRAFT-Best-Practice-Guide-for-Data-Gap-Analysis-for-biodiversity-stakeholders.pdf](http://www.gbif.org/system/files_force/gbif_resource/resource-82566/DRAFT-Best-Practice-Guide-for-Data-Gap-Analysis-for-biodiversity-stakeholders.pdf) Arturo H. Ariño, Vishwas Chavan, Javier Otegui

VertNet Migrators <https://github.com/VertNet/toolkit>

VertNet Portal Spatial Quality Tab <http://www.vertnet.org/resources/spatialqualitytabguide.html>

VertNet GitHub Reference & Set Up for Data Issue Tracking <http://www.vertnet.org/resources/issuetrackingguide.html>

Improving Data Quality: iDigBio Recordset data cleaning method, tools, and data flags <https://www.idigbio.org/content/improving-data-quality-idigbio-recordset-data-cleaning-method-tools-and-data-flags>

A summer learning R to clean up data with the iDigBio portal recordset correction feature <https://www.idigbio.org/content/summer-learning-r-clean-data-idigbio-portal-recordset-correction-feature>

iDigBio Data recommendations for optimal searchability and applicability in the aggregate [https://www.idigbio.org/wiki/index.php/Data\\_Ingestion\\_Guidance#Data\\_recommendations\\_for\\_optimal\\_searchability\\_and\\_applicability\\_in\\_the\\_aggregate](https://www.idigbio.org/wiki/index.php/Data_Ingestion_Guidance#Data_recommendations_for_optimal_searchability_and_applicability_in_the_aggregate)

Exploring unique values in iDigBio using Apache Spark <https://www.idigbio.org/content/exploring-unique-values-idigbio-using-apache-spark>



# Acknowledgements and References, cont'd

Biocode Field Information Management System [ppt](#) [youtube](#). A Field Information Management System (FIMS) enables data collection at the source (in the field) by generating spreadsheet templates, validating data, and assigning persistent identifiers for every unique biological sample. The following diagram shows how the system works. The most typical functions are Generating Templates and Validating Data, both of which can be found under the Tools menu.

[Generate a Template](#)

[Validate data](#)

[How FIMS works](#)

## Contact us

Deb Paul [dpaul@fsu.edu](mailto:dpaul@fsu.edu)

Katja Seltmann [seltmann@ccber.ucsb.edu](mailto:seltmann@ccber.ucsb.edu)

Laura Russell [larussell@vertnet.org](mailto:larussell@vertnet.org)

David Bloom [dbloom@vertnet.org](mailto:dbloom@vertnet.org)

iDigBio is funded by a grant from the National Science Foundation's Advancing Digitization of Biodiversity Collections Program (Cooperative Agreement EF-1115210). Any opinions, findings, and conclusions or recommendations expressed in this material are those of the author(s) and do not necessarily reflect the views of the National Science Foundation.

