

# CYWG (CyberInfrastructure Working Group) October 2014

## iDigBio Data Ingestion

Dan Stoner

Advanced Computing and Information Systems Laboratory (ACIS)  
University of Florida

 [dstoner@acis.ufl.edu](mailto:dstoner@acis.ufl.edu)

 [@thatlinuxbox](https://twitter.com/thatlinuxbox)

Over 300 Data Providers...





... and many more.

## Data Ingestion Progress

November 2013 -

4.2 million specimen records

0.9 million media records

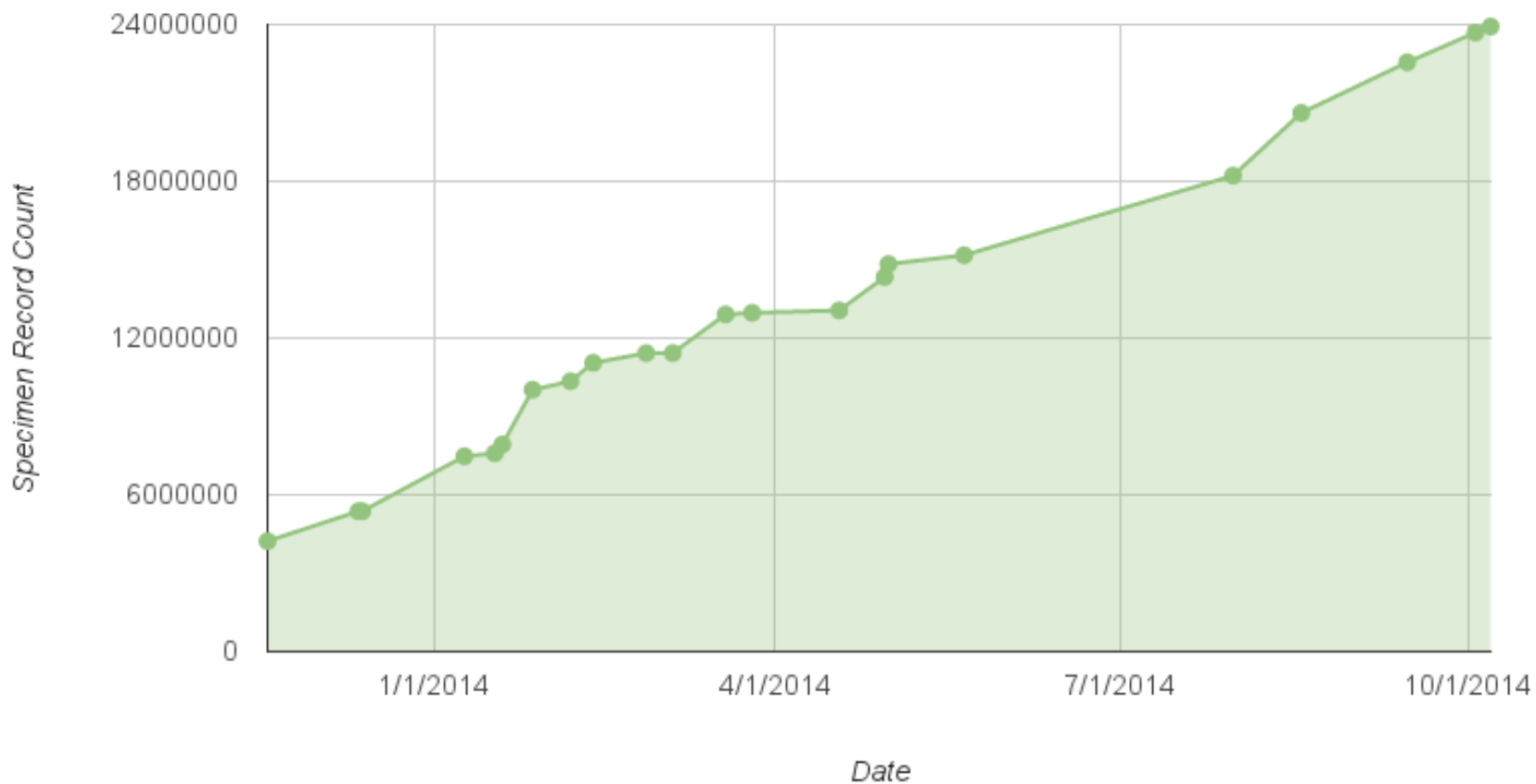


October 2014 -

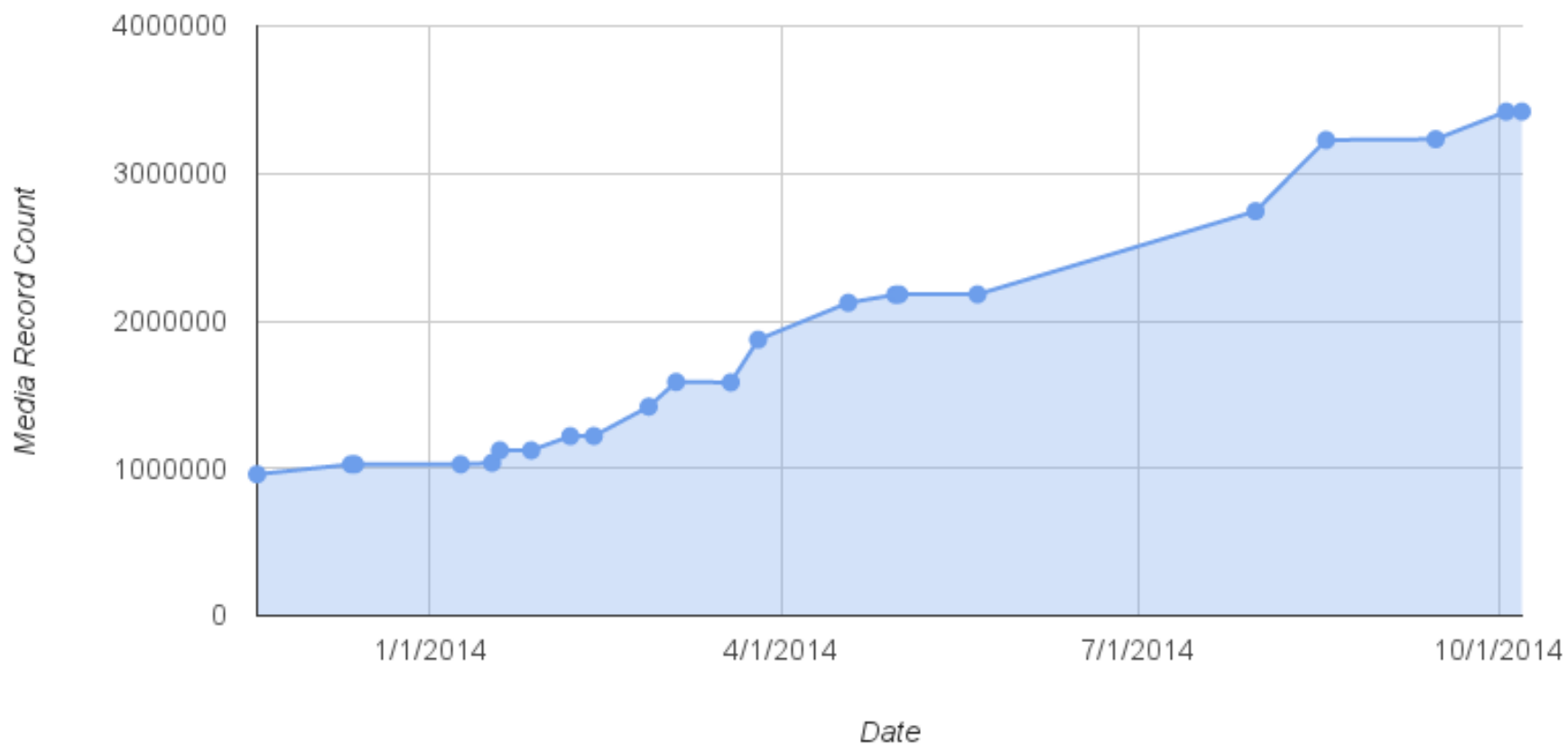
23.9 million specimen records

3.4 million media records

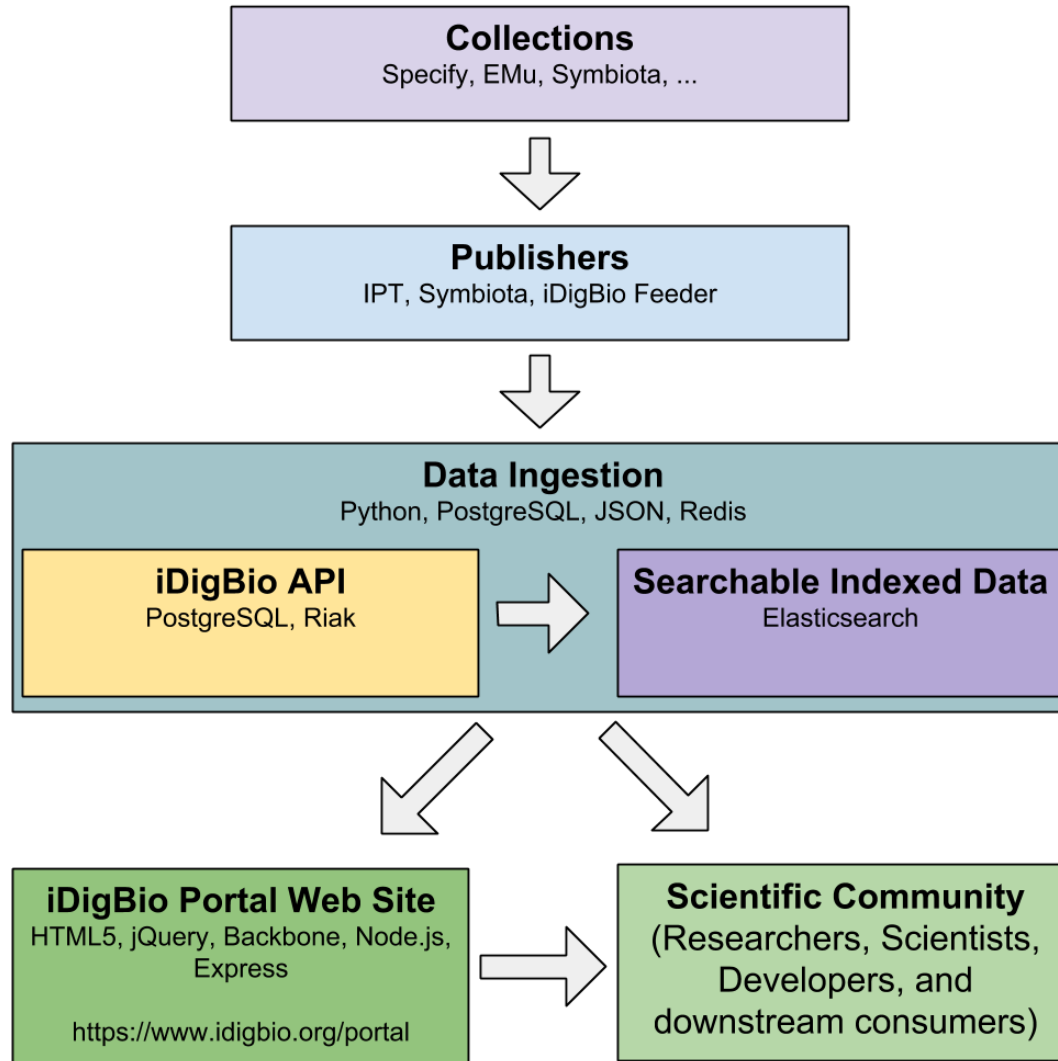
### iDigBio Data Ingestion - Specimen Records



### iDigBio Data Ingestion - Media Records



## iDigBio Data Flow Diagram



Three types of data publishing technologies currently being consumed by iDigBio:

**GBIF Integrated Publishing Toolkit (IPT)** - a Java tool used to publish and share biodiversity datasets

<http://www.gbif.org/ipt/>

**Symbiota** – web-based collection management software

<http://symbiota.org/>

**iDigBio RSS Feeder** – data sharing service for providers who do not run infrastructure



# Dataset Formats Consumable by iDigBio

- IPT – DwC-A
- Symbiota portals – DwC-A
- iDigBio Feeder – DwC-A, CSV, ...

**If you can export specimen data from your system / database / spreadsheet into DwC-A (or even CSV), then you can share data with iDigBio.**

iDigBio RSS Feeder facilitates the sharing of over 1.5 million specimen records and 200 thousand media records from providers who do not need to run “servers”.

# Data Source Types Providing Data to iDigBio

The following collection systems, databases, applications are known to have a capability to be a data source for iDigBio.

- Specify Software Project
- EMu Museum Management System
- Symbiota
- Arctos
- Excel
- ...

The iDigBio Mobilization Team ([data@idigbio.org](mailto:data@idigbio.org)) assist with the preparation of data sets prior to Data Ingestion and are available to answer questions about sharing data with iDigBio.

See Also:

[https://www.idigbio.org/wiki/index.php/Digitization\\_Resources](https://www.idigbio.org/wiki/index.php/Digitization_Resources)

## Darwin Core Archive / DwC-A

<http://rs.tdwg.org/dwc/terms/guides/text/>

A Darwin Core Archive is a zip file that includes metadata about the dataset, the data itself, and any optional extension data.

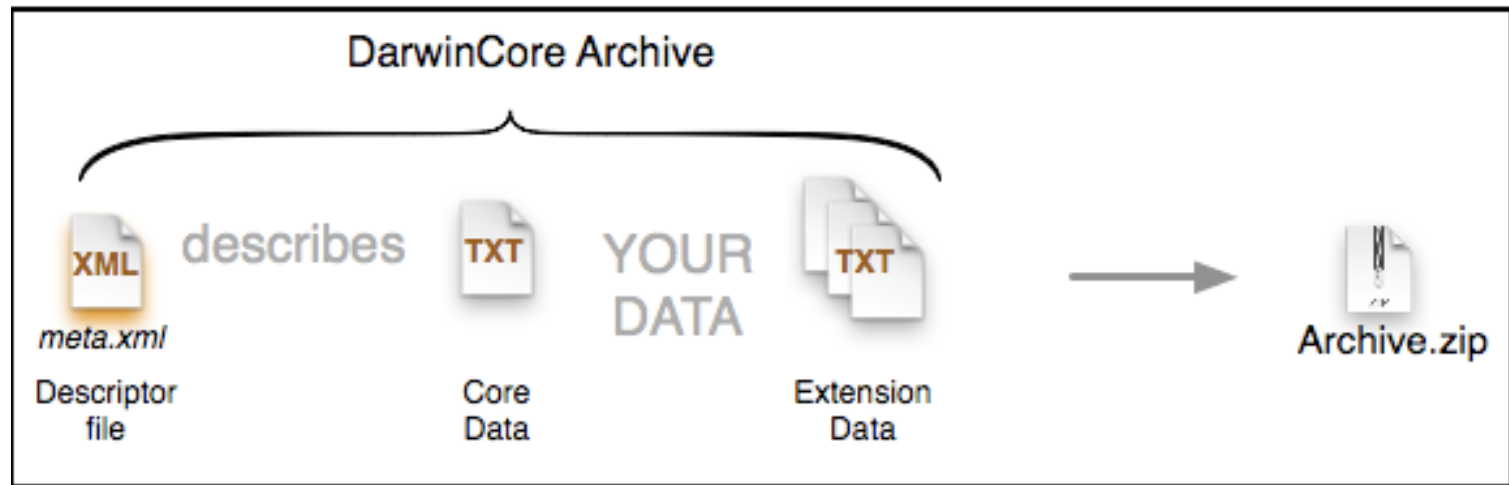


Image source: <http://tools.gbif.org/dwca-assistant/>

# Specimen Data – Darwin Core Standard

<http://rs.tdwg.org/dwc/terms/>

Field	Records With This Field	(%) Percent Used
<b>Institution Code</b> (dwc:institutionCode)	41,262	100
<b>Catalog Number</b> (dwc:catalogNumber)	41,262	100
<b>Collection Code</b> (dwc:collectionCode)	41,262	100
<b>Occurrence ID</b> (dwc:occurrenceID)	41,262	100
<b>Basis of Record</b> (dwc:basisOfRecord)	41,262	100
<b>Kingdom</b> (dwc:kingdom)	41,261	99.998
<b>Phylum</b> (dwc:phylum)	41,261	99.998
<b>Class</b> (dwc:class)	41,261	99.998
<b>Order</b> (dwc:order)	41,261	99.998
<b>Family</b> (dwc:family)	41,261	99.998
<b>Scientific Name</b> (dwc:scientificName)	41,261	99.998
<b>Locality</b> (dwc:locality)	41,248	99.966
<b>Specific Epithet</b> (dwc:specificEpithet)	41,157	99.746
<b>Genus</b> (dwc:genus)	41,124	99.666
<b>Continent</b> (dwc:continent)	40,963	99.275

## Recommended minimum fields for iDigBio Ingestion:

Record ID (recordId) - unique identifier for the digital record

Occurrence ID (occurrenceID) - unique identifier for the physical object or establishment of an Occurrence

Scientific Name (scientificName) - the full scientific name

Event Date (eventDate) - date-time, preferably in ISO 8601

Collector (recordedBy) - collector name, number, or field number

Locality Data (...) - verbatim and decimal locality fields, continent, country, water body, state/province, ....

Catalog Number (catalogNumber) - Barcode, catalog number, accession id or collection number

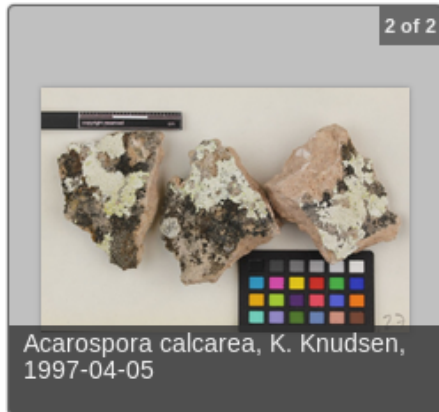
Institution Code (institutionID) - institution identifier

Collection Code (collectionID) - collection identifier

Paleo specimens should also include Geological Context fields.

## Media Data – Audubon Core / AC

[http://terms.tdwg.org/wiki/Audubon\\_Core\\_Term\\_List](http://terms.tdwg.org/wiki/Audubon_Core_Term_List)



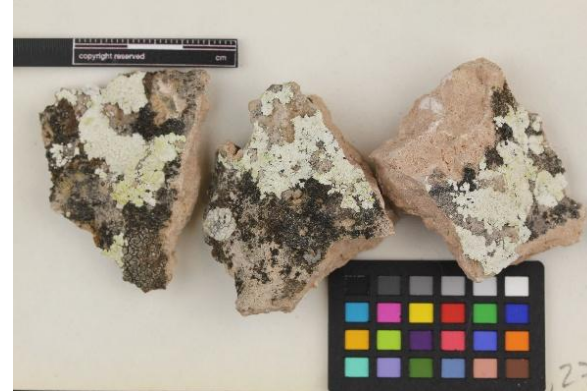
*Images Source: Arizona State University Lichen Herbarium (Accessed through iDigBio Specimen Data Portal, <https://www.idigbio.org/portal>, 2014-09-18)*

GBIF has a nice write-up on the benefits of AC over dwc:associatedMedia:

<http://gbif.blogspot.com/2014/05/multimedia-in-gbif.html>

## Audubon Core vocabularies address such concerns as:

- the management of the media and collections
- descriptions of their content
- their taxonomic, geographic, and temporal coverage
- appropriate ways to retrieve, attribute and reproduce them



### Media Metadata

<b>Associated Specimen Reference</b>	<a href="http://lichenportal.org/portal/collections/individual/index.php?occid=1374628">http://lichenportal.org/portal/collections/individual/index.php?occid=1374628</a>
<b>Type of Resource Subtype</b>	StillImage Photograph
<b>Metadata Date</b>	2013-04-24 02:00:19
<b>Provider-managed ID</b>	urn:uuid:9a77ed32-7fa4-4831-938e-a499078058a8
<b>Credit</b>	Arizona State University Lichen Herbarium (ASU)
<b>License Terms</b>	CC BY-NC-SA (Attribution-NonCommercial-ShareAlike)
<b>License URL</b>	<a href="http://creativecommons.org/licenses/by-nc-sa/3.0/">http://creativecommons.org/licenses/by-nc-sa/3.0/</a>
<b>Access URI</b>	<a href="http://storage.idigbio.org/asu/lichens/ASU0068/ASU0068021a_lg.jpg">http://storage.idigbio.org/asu/lichens/ASU0068/ASU0068021a_lg.jpg</a>
<b>Format</b>	image/jpeg

## Audubon Core can support “new” media types

### Specimen Data

dwc:catalogNumber: UF 105199

dwc:scientificName:

*Carcharocles megalodon*

dwc:stateProvince: Florida

dwc:county: Duval

dwc:latestPeriodOrHighestSystem:

Late Miocene

dwc:decimalLatitude: 30.39211

### Media Data

dwc:scientificName:

*Carcharocles megalodon*

dc:type: image

ac:subtype: <http://www.fabbers.com/StL.asp>

ac:subtypeLiteral: 3dModel

ac:tag: tooth

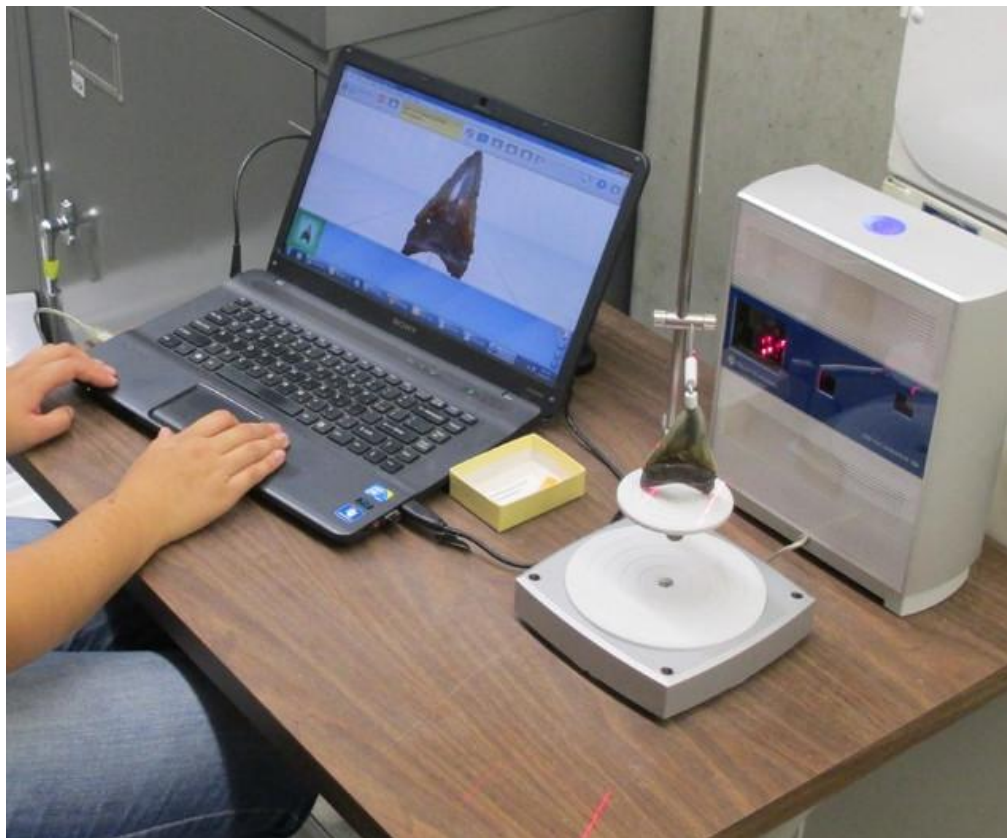


Image source: Aaron Wood, Florida Museum of Natural History



## Darwin Core Archive - Extension to link data between the Occurrence and Audubon Core record

### Specimen Data

dwc:occurrenceID: **3bca767a-5a25-42c0-12...**

dwc:scientificName: Carcharocles megalodon

dwc:catalogNumber: UF 105199

dwc:stateProvince: Florida

dwc:county: Duval

dwc:latestPeriodOrHighestSystem: Miocene

dwc:decimalLatitude: 30.39211

...

### Media Data

dcterms:identifier: 9a3025b1-f686-4e43-915f-...

coreid: **3bca767a-5a25-42c0-12...**

dwc:scientificName: Carcharocles megalodon

ac:associatedSpecimenReference: <http://museum...>

dc:type: image

ac:subtype: <http://www.fabbers.com/StL.asp>

ac:subtypeLiteral: 3dModel

ac:tag: tooth

...

In the wild, ac:associatedSpecimenReference tends NOT to provide the bare occurrence id of the related specimen, so instead we use the implicit relationship via coreid in the DwC-A.



Image source: [http://commons.wikimedia.org/wiki/File:Carcharocles\\_megalodon\\_tooth.JPG](http://commons.wikimedia.org/wiki/File:Carcharocles_megalodon_tooth.JPG)

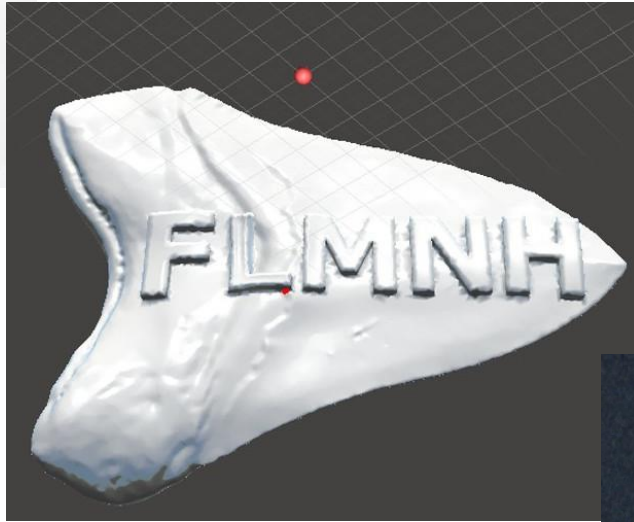


Image source: Aaron Wood



Image source: Aaron Wood  
3D Model printing by Robert Burns

# Recommended minimum Audubon Core fields for iDigBio Data Ingestion:

Access URI

Rights

Provider

Scientific name

Title

Description

Tags

# Practical Details

## Data Formats

- ISO 8601 Dates

- WGS84 Decimal Lat/Long

## Controlled Vocabularies

- ISO Country Names and Codes

- State/Province names

- Identifier Formats (UUID, ARK, URN, DOI, URI, URL, LSID, ...)

- Copyright and Standard Licenses

- Apple Core guidelines for herbaria

  - <http://code.google.com/p/applecore/wiki/Introduction>

# Ingestion Process Changes Over the Past Year

- New Staff (Dan Stoner... that's me!)
- Improved parallelization of ingestion tasks
- Incremental Indexing
- Database Tuning
- Ingestion Reporting

# Ingestion Reporting

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

### Publisher Summary

Publisher Name	Record Count			Media Record Count		
	Digest	API	Index	Digest	API	Index
<a href="#">Berkeley Natural History Museums IPT</a>	1,860,584	1,859,985	1,859,985	0	0	0
<a href="#">Florida Museum of Natural History IPT Service</a>	1,047,587	1,047,587	1,047,587	0	0	0
<a href="#">MyCoPortal Darwin Core Archive rss feed</a>	1,679,459	1,679,458	1,679,458	371,346	371,346	371,346
<a href="#">Northern Great Plains Herbaria Darwin Core Archive rss feed</a>	43,012	43,012	43,012	0	0	0
<a href="#">KU Biodiversity Institute IPT</a>	2,010,071	2,011,170	2,011,170	0	0	0
<a href="#">The University of Connecticut Biological Collections</a>	172,098	172,102	172,102	166,689	166,707	166,707
<a href="#">xBioD IPT in the Museum of Biological Diversity at the Ohio State University</a>	521,710	521,782	521,782	2,593	2,593	2,593
<a href="#">CMC_specify</a>	9,131	9,131	9,131	0	0	0
<a href="#">Consortium of North American Bryophyte Herbaria Darwin Core Archive rss feed</a>	1,690,014	1,690,014	1,690,014	816,932	816,932	816,932
<a href="#">Museum of Comparative Zoology, Harvard University</a>	1,736,357	1,736,471	1,736,471	0	0	0
<a href="#">CNALH Darwin Core Archive rss feed</a>	1,232,891	1,232,891	1,232,891	649,241	649,241	649,241
<a href="#">SCAN Darwin Core Archive rss feed</a>	873,024	873,160	873,160	68,696	68,718	68,718
<a href="#">iDigBio Feeder RSS Feed</a>	1,316,574	1,316,574	1,316,574	19,024	19,024	19,024
<a href="#">Consortium of Intermountain Herbaria Darwin Core Archive rss feed</a>	204,129	204,131	204,131	74,014	74,015	74,015
<a href="#">CAS-IPT</a>	1,875,928	1,875,979	1,875,979	0	0	0
<a href="#">Macroalgal Herbarium Portal Darwin Core Archive rss feed</a>	2,145	2,145	2,145	1,937	1,937	1,937
<a href="#">CNH portal Darwin Core Archive rss feed</a>	89,199	89,199	89,199	56,557	56,557	56,557
<a href="#">IPT - Hosted by VertNet</a>	5,070,222	5,070,222	5,070,222	479,440	479,440	479,440
<a href="#">North American Network of Small Herbaria Darwin Core Archive rss feed</a>	4,162	4,162	4,162	4,273	4,273	4,273
<a href="#">Harvard University Herbaria IPT installation</a>	412,331	412,331	412,331	295,055	295,055	295,055
<a href="#">SNOMNH IPT</a>	310,328	310,328	310,328	0	0	0
<a href="#">Morphbank IPT Feed</a>	48,567	97,127	97,127	0	65,167	65,167
<a href="#">SEINet Darwin Core Archive rss feed</a>	347,210	347,216	347,216	161,533	161,627	161,627

## Planned Future Changes

- Parallelize more parts of Ingestion process (such as media processing)
- Support for additional publisher types (beyond IPT, Symbiota, iDigBio RSS Feeder)
- Improved Ingestion logging and error detection
- Support for additional media types (audio, 3D scans, ...)
- Data Quality

# Thank You!



[www.idigbio.org](http://www.idigbio.org)



[facebook.com/iDigBio](https://facebook.com/iDigBio)



[twitter.com/iDigBio](https://twitter.com/iDigBio)



[vimeo.com/idigbio](https://vimeo.com/idigbio)



[idigbio.org/rss-feed.xml](http://idigbio.org/rss-feed.xml)



[webcal://www.idigbio.org/events-calendar/export.ics](http://webcal://www.idigbio.org/events-calendar/export.ics)



End.