

Establishing a New Framework for Paleontological Data Through an Evaluation of Current Data Sharing Practices

Holly Little ◆ littleh@si.edu

Department of Paleo Informatics Manager
Smithsonian National Museum of Natural History

TDWG 2018



Smithsonian Institution

NATIONAL
MUSEUM *of*
**NATURAL
HISTORY**

Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records

Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records



1970-present: Creation of digital records

Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records



1970-present: Creation of digital records

Global

10,042,739 GBIF Fossil Occurrences



2012: Darwin Core (DwC) Standard



Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records



1970-present: Creation of digital records

Global

10,042,739 GBIF Fossil Occurrences



2012: Darwin Core (DwC) Standard



Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records



1970-present: Creation of digital records

Global



10,042,739 GBIF Fossil Occurrences

2012: Darwin Core (DwC) Standard



Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records



1970-present: Creation of digital records

Global

10,042,739 GBIF Fossil Occurrences



2012: Darwin Core (DwC) Standard



Local: NMNH Paleo

662,353 (+-) Fossil Occurrence Records



1970-present: Creation of digital records

- ❖ What data do we share?
- ❖ What isn't being shared? Why?
- ❖ Is the data clean and standardized?
- ❖ Are there other data formats or standards we could be sharing in?
 - New terms needed?

Global



10,042,739 GBIF Fossil Occurrences

2012: Darwin Core (DwC) Standard



NMNH Paleo Data

81 DwC terms

Issues and flags	Count	
Taxon match none	153,080	<div style="width: 100%;"></div>
Taxon match higherrank	97,043	<div style="width: 95%;"></div>
Geodetic datum assumed WGS84	68,327	<div style="width: 90%;"></div>
Taxon match fuzzy	15,263	<div style="width: 60%;"></div>
Country derived from coordinates	12,692	<div style="width: 55%;"></div>
Recorded date unlikely	3,187	<div style="width: 25%;"></div>
Country invalid	1,317	<div style="width: 15%;"></div>
Country coordinate mismatch	617	<div style="width: 10%;"></div>
Recorded date invalid	207	<div style="width: 5%;"></div>
Depth min/max swapped	92	<div style="width: 2%;"></div>

type	endDayOfYear	decimalLatitude	member
references	year	decimalLongitude	identification
institutionID	month	geodeticDatum	typeStatus
institutionCode	day	coordinateUncertaintyInMeters	identifiedBy
collectionCode	verbatimEventDate	verbatimLatitude	scientificName
datasetName	habitat	verbatimLongitude	higherClassification
basisOfRecord	fieldNotes	verbatimCoordinateSystem	kingdom
occurrenceID	locationID	georeferenceProtocol	phylum
catalogNumber	higherGeography	georeferenceRemarks	class
recordNumber	continent	earliestEraOrLowestErathem	order
recordedBy	waterBody	latestEraOrHighestErathem	family
individualCount	islandGroup	earliestPeriodOrLowestSystem	genus
sex	island	latestPeriodOrHighestSystem	subgenus
lifeStage	country	earliestEpochOrLowestSeries	specificEpithet
preparations	stateProvince	latestEpochOrHighestSeries	infraspecific
associatedMedia	county	earliestAgeOrLowestStage	taxonRank
associatedSequences	locality	latestAgeOrHighestStage	scientificName
otherCatalogNumbers	minimumElevationInMeters	group	
occurrenceRemarks	maximumElevationInMeters	Formation	
associatedOccurrences	verbatimElevation	verbatimDepth	
fieldNumber	minimumDepthInMeters		
startDayOfYear	maximumDepthInMeters		

Micropaleo Data

How can we compare the micropaleo collections from the NMNH, NHM London, and AMNH?

- Giles Miller (NHM London)

Initial Query parameters:

- fossilSpecimen + Foraminifera
- fossilSpecimen + Ostracoda
- datasetName:PAL + Foraminifera
- Ostracoda
- Foraminifera
- datasetName: AMNH MICROFOSSIL

- ✓ institutionCode
- ✓ collectionCode
- ✓ basisOfRecord
- ✓ occurrenceID
- ✓ catalogNumber

- ✗ ✓ Event terms
- ✗ ✓ Locality terms

DwC Term	NMNH	NHM	AMNH
institutionCode	✓	✓	✓
collectionCode	✓	✓	✓
basisOfRecord	✓	✓	✓
occurrenceID	✓	✓	✓
catalogNumber	✓	✓	✓
Event terms	✓	✓	✓
Locality terms	✓	✓	✓
...
101 DwC Terms

101 DwC Terms used across 3 datasets

API Query

Datasets

Type: Occurrence

Metadata containing: Fossil

<http://api.gbif.org/v1/dataset/search?q=fossil&type=occurrence>

Occurrence Records

basisOfRecord: fossilSpecimen

[https://api.gbif.org/v1/occurrence/search?basis_of_record=FOSSIL SPECIMEN](https://api.gbif.org/v1/occurrence/search?basis_of_record=FOSSIL_SPECIMEN)

~~API Query~~ Download

Occurrence Records

DOWNLOAD | 9 AUGUST 2018

10,042,739 occurrences downloaded

DOI 10.15468/dl.r3tqg5

DOWNLOAD

FILTER APPLIED 9 AUGUST 2018

RERUN QUERY

Citation: GBIF.org (09 August 2018) GBIF Occurrence Download <https://doi.org/10.15468/dl.r3tqg5>

License: [CC BY 4.0](#)

File: 2 GB Darwin Core Archive

Involved datasets: [1,486](#)

Make sure to read the [data user agreement](#) and [citation guidelines](#).

fRecord: fossilSpecimen

https://api.gbif.org/v1/occurrence/reh?basis_of_record=FOSSIL_SPECIMEN

Datasets

Total datasets: 1646

Tool: Python*

Sources: API Query (160)
dataset list from
occurrence download
(1486)

**All scripts still need some debugging and cleaning*

API Request

1. Query GBIF for datasets with 'fossil'
2. Extract GBIF dataset Key
3. Query GBIF with dataset key
4. Find DwC-A endpoint
5. Request Zip file from endpoint
6. Unzip DwC-A
7. Extract metadata for dataset CSV report

Key elements:

- Filtered down to endpoints that were tagged as DWC_ARCHIVE
- Account for inactive links or bad zip files

```
dataset_api_short.py
9 URL_DS = "http://api.gbif.org/v1/dataset/search?" #GBIF API dataset endpoint
10 DS_parameters = {"q": "fossil", "type": "occurrence", "limit": 200} #query
11
12 r = requests.get(url = URL_DS, params = DS_parameters)
13 data = r.json()
14 num_results=data['count']
15
16 sum_list = [] #used to create final rows for each dataset in the output
17 sum_keys = ['key','title','publishingOrganizationTitle','license','publi
18
- {
  key: "96830f08-f762-11e1-a439-00145eb45e9a",
  title: "Collection Graptolithina fossil SMF",
  description: "fossil Graptolithina of the world",
  type: "OCCURRENCE",
  hostingOrganizationKey: "c76cf030-2a95-11da-9cc1-b8a03c50a862",
  hostingOrganizationTitle: "Senckenberg",
  countryCoverage: [ ],
  publishingCountry: "DE",
  publishingOrganizationKey: "c76cf030-2a95-11da-9cc1-b8a03c50a862",
  publishingOrganizationTitle: "Senckenberg",
  license: "http://creativecommons.org/licenses/by/4.0/legalcode",
  decades: [ ],
  keywords: [ ],
  recordCount: 1006
},
```

```
18
19 sum_list.append(summary)
20
21 if endpoints == "DWC_ARCHIVE":
22     try:
23         url_dwc_a = data_key['endpoints'][0]['url']
24         rqt = requests.get(url = url_dwc_a)
25         z = zipfile.ZipFile(io.BytesIO(rqt.content))
26         z.extractall(key)
27     except zipfile.BadZipFile:
28         url_error = "Broken url: "+url_dwc_a
29         print(url_error)
30     else:
31         other_endpoint = "Endpoint is " + endpoints
32         print(other_endpoint)
33
34 sum_keys.append('citation')
35 sum_keys.append('accessDate')
36 sum_keys.append('endpoint')
37
38 with open('gbif_datasets.csv', 'w') as f:
39     writer = csv.writer(f)
40     writer.writerow(sum_keys)
41     for summary in sum_list:
42
```

API Request

- ~~1. Query GBIF for datasets with 'fossil'~~
- ~~2. Extract GBIF dataset Key~~
3. Query GBIF with dataset key
4. Find DwC-A endpoint
5. Request Zip file from endpoint
6. Unzip DwC-A
7. Extract metadata for dataset CSV report

Key elements:

- Filtered down to endpoints that were tagged as `DWC_ARCHIVE`
- Account for inactive links or bad zip files

```
dataset_api.py  
9 URL_DS = "http://api.gbif.org/v1/dataset/search?" #GBIF API dataset endpoint  
10 DS_parameters = {"q": "fossil", "type": "occurrence", "limit": 200} #query  
11  
12 r = requests.get(url = URL_DS, params = DS_parameters)  
13 data = r.json()  
14 num_results=data['count']  
15  
16 sum_list = [] #used to create final rows for each dataset in the output  
17 sum_keys = ['key', 'title', 'publishingOrganizationTitle', 'license', 'publi
```

```
01cec130-c69f-497e-ab15-0534c50e003d.xml  
1cf7c55e-0258-44fe-a556-b8b010f7b0ac.xml  
1de4d9a0-7430-44ad-8367-cc0298906ad8.xml  
1eb5e969-4412-4f08-81ec-3de057e559a1.xml  
1ec9f790-5490-4656-8dac-396b49c7cd41.xml  
1f2cfb6f-c91b-498e-80f3-8eeec688292.xml  
1f45504c-460b-446d-a10d-c6841e98cb67.xml  
2a629a9a-38d1-496b-afb7-b4ff3b8fae60.xml  
2b5a4824-aa72-488b-9c13-118436f0d969.xml  
2ce84acf-0e82-4da9-9cda-730973a51d7f.xml  
2e4cc37b-302e-4f1b-bbbb-1f674ff90e14.xml  
2e6b7086-17a9-463c-a355-43738c0ad85b.xml  
02e39512-ccc8-4e30-aa5a-8e86b8952116.xml  
2ec89f62-a165-4f9f-80ff-ea475845996f.xml  
2f260519-acf1-420a-a4f0-f52b8b7e27d8.xml  
2f391085-e522-4662-afa6-953433267374.xml  
2fba9985-ac30-46cb-99bf-91ccde0d8d2f.xml  
2fd293d3-1df5-4b60-912e-40aac75e3ace.xml  
3a6c8f26-3987-42d0-899f-d1b6699b41b6.xml  
3a7d0a6f-e525-4ccc-892e-47008d0e88d2.xml  
3a68981d-33da-433e-9d0d-f7e7837e5767.xml  
3ba47554-1afd-482a-aea6-b79661b7724a.xml  
3be0f16e-fd4c-4513-a454-ab63373541bc.xml  
3c57ed85-e109-4e4f-92ea-8c007cde0903.xml  
3c001217-eea8-4f59-8b28-885699f8cd6c.xml  
3d725a18-dd38-42fe-bd4b-1aa81381930a.xml  
3e603477-0c14-4236-bf54-44840c59aba5.xml  
3ed891d2-ae3c-45a6-a27f-e9f0a9f4b8c1.xml  
3f42d91f-ccc2-46c2-8907-f117aab0b7c0.xml  
ts): #Loop through all datasets returned by api  
['key'] #key is used to access the full detail  
if.org/v1/dataset/" + key #the full json file  
URL_KEY)  
dpoints']][0]['type']  
y for storing summary data about each dataset  
ata to pull for data summary sheet  
sults']['ds][keys]  
Found"  
ta_key['citation']['text']  
dt.strftime(today, '%m/%d/%Y')  
dpoints  
4IVE':  
_key['endpoints']][0]['url']  
et(url = url_dwc_a)  
ile(io.BytesIO(rqt.content))  
pFile:  
ken url: "+url_dwc_a  
dpoint is "+endpoints  
57 with open('gbif_datasets.csv', 'w') as f:  
58     writer = csv.writer(f)  
59     writer.writerow(sum_keys)  
60     for summary in sum_list:
```

```
57 with open('gbif_datasets.csv', 'w') as f:  
58     writer = csv.writer(f)  
59     writer.writerow(sum_keys)  
60     for summary in sum_list:
```


Total datasets: 1476

key	title	publis license	publishingCc	recordCount	citation	accessDate	endpoint
7ca4f6dc-f762-11e1-a439-00145eb45e9a	(Table 2) Distribution of ice-rafted fossils in	PANG. http://creativecommons.org/licenses/by/4.0/	DE	48	Spiegler D (1988). (Table 2) Distribution of ice-raft	8/19/18	DWC_ARCHIVE
b275a4c1-9859-4f3c-8ead-d86dde82f0bc	The fossil collection (F) of the MusÃ©um r	MNHN http://creativecommons.org/licenses/by/4.0/	FR	Key Not Found	MNHN - Museum national d'Histoire naturelle (20	8/19/18	DWC_ARCHIVE
29334ea1-9b35-4370-9873-8d85ad56be6e	Fossil occurrence of planktonic foraminife	PANG. http://creativecommons.org/licenses/by/4.0/	DE	18	Minoshita M, Tobin H, Ashi J, Kimura G, Lallemand	8/19/18	DWC_ARCHIVE
815c3a6e-f762-11e1-a439-00145eb45e9a	(Table 1) Distribution of early late Eocene	PANG. http://creativecommons.org/licenses/by/4.0/	DE	2610	Nicora A, Premoli Silva I (1989). (Table 1) Distrib	8/19/18	DWC_ARCHIVE
680c12e3-6440-4559-8add-b263cb48f31a	(Table 2) Fossil record of foraminifera and	PANG. http://creativecommons.org/licenses/by/4.0/	DE	404	Wosziidlo H (1961). (Table 2) Fossil record of forar	8/19/18	DWC_ARCHIVE
2e6b7086-17a9-463c-a355-43738c0ad85b	(Table 1) Fossil record of foraminifera and	PANG. http://creativecommons.org/licenses/by/4.0/	DE	353	Wosziidlo H (1961). (Table 1) Fossil record of forar	8/19/18	DWC_ARCHIVE
d6618bc2-ce35-4c16-94e1-d58063736899	(Table 3) Minerals, and calcareous and silic	PANG. http://creativecommons.org/licenses/by/4.0/	DE	286	Gervais E (1986). (Table 3) Minerals, and calcareo	8/19/18	DWC_ARCHIVE
b201c8fb-0502-437f-ba77-8e092273328b	(Table 5) Record of mollusca and ostracod.	PANG. http://creativecommons.org/licenses/by/4.0/	DE	120	Krienke H, Strahl J, Frenzel P, Keding E (1998). (Ta	8/19/18	DWC_ARCHIVE
81ca34fe-f762-11e1-a439-00145eb45e9a	(Table 2) Distribution of early late Eocene	PANG. http://creativecommons.org/licenses/by/4.0/	DE	36	Nicora A, Premoli Silva I (1989). (Table 2) Distrib	8/19/18	DWC_ARCHIVE
f6a07b42-1d2c-11e2-8f64-00145eb45e9a	Collection Paleontology - GPIT	Sencki http://creativecommons.org/licenses/by/4.0/	DE	28350	Senckenberg. Collection Paleontology - GPIT. Occi	8/19/18	BIOCASE
1aecc602-8847-4980-b493-bb21bb4aace5	PalÃ¤ontologie Marburg	Sencki http://creativecommons.org/licenses/by/4.0/	DE	5166	Senckenberg. PalÃ¤ontologie Marburg. Occurren	8/19/18	BIOCASE
40275fc6-3e71-4d76-aafe-276bdeefcd86	SDSM Botany and Paleobotany Collections	South http://creativecommons.org/publicdomain/zi	US	Key Not Found	Hess G, Shelton S (2018). SDSM Botany and Paleo	8/19/18	DWC_ARCHIVE
71fa736a-f762-11e1-a439-00145eb45e9a	MfN - Glacial erratics	Musei http://creativecommons.org/licenses/by/4.0/	DE	2	Museum fÃ¼r Naturkunde Berlin. MfN - Glacial er	8/19/18	BIOCASE
21a0cb45-f0ec-4901-a868-cd45abe74601	Canadian Museum of Nature Palynology C	Canad http://creativecommons.org/licenses/by/4.0/	CA	14568	Shepherd K, Shorthouse D (2018). Canadian Muse	8/19/18	DWC_ARCHIVE
53e9eca5-6831-4a43-a109-130feecfcd9c	Cleveland Museum of Natural History Inve	Clevel http://creativecommons.org/publicdomain/zi	US	20585	Dunn D (2018). Cleveland Museum of Natural His	8/19/18	DWC_ARCHIVE
c8681cc2-9d0a-4c5f-b620-5c753abfe2bc	NMNH Paleobiology Specimen Records	Natio http://creativecommons.org/publicdomain/zi	US	662355	Orrell T, Hollowell T (2018). NMNH Paleobiology S	8/19/18	DWC_ARCHIVE
6720aee6-2aad-446d-bb97-ba009d1b5666	CM Vertebrate Paleontology Collection	Carnej http://creativecommons.org/publicdomain/zi	US	79254	Henrici A (2016). CM Vertebrate Paleontology Col	8/19/18	DWC_ARCHIVE
83216388-f762-11e1-a439-00145eb45e9a	ColecciÃ³n PaleontolÃ³gica Invertebrados	CCT Cí http://creativecommons.org/licenses/by/4.0/	AR	190	CCT CONICET-CENPAT Centro CientÃ­fico TecnolÃ	8/19/18	DIGIR
361d98ca-ba9d-4b85-a64f-7db8fb35c6a9	CMC Cincinnati Museum Center Invertebr	Cincin http://creativecommons.org/licenses/by-nc/4	US	60165	Hunda B, Kling A, Storrs G (2016). CMC Cincinnati	8/19/18	DWC_ARCHIVE
95da13d0-f762-11e1-a439-00145eb45e9a	ELM paleontology	Geocc http://creativecommons.org/licenses/by-nc/4	EE	24765	Geocollections of Estonia. ELM paleontology. Occ	8/19/18	BIOCASE
8130e5c6-f762-11e1-a439-00145eb45e9a	IG TUT paleontology	Geocc http://creativecommons.org/licenses/by-nc/4	EE	110964	Geocollections of Estonia. IG TUT paleontology. O	8/19/18	BIOCASE
8f7e3c45-4d76-4982-9478-651439e0cd4b	Natural History Museum, University of Tar	PlutoF http://creativecommons.org/licenses/by-nc/4	EE	Key Not Found	PlutoF (2018). Natural History Museum, Universit	8/19/18	DWC_ARCHIVE
561367c9-576b-4aa8-8fc9-1981155b11af	Museo Argentino de Ciencias Naturales "B	Musei http://creativecommons.org/licenses/by/4.0/	AR	Key Not Found	Gutierrez P R, RodrÃ­guez D (2018). Museo Argen	8/19/18	DWC_ARCHIVE
85530b4f-67af-4c8e-a176-0ea3c322d9fc	Museo Argentino de Ciencias Naturales "B	Musei http://creativecommons.org/licenses/by/4.0/	AR	Key Not Found	Genise J, RodrÃ­guez D (2018). Museo Argentino	8/19/18	DWC_ARCHIVE
f5783452-6ff5-4bef-8f30-9adfeba81bd7	Paleontological Research Institution Collec	Paleor http://creativecommons.org/publicdomain/zi	US	19541	Skibinski L (2018). Paleontological Research Instit	8/19/18	DWC_ARCHIVE
e1e16cf0-ada2-11e2-8fbc-00145eb45e9a	Neptune Deep-Sea Microfossil Occurrence	Musei http://creativecommons.org/licenses/by/4.0/	DE	500808	Museum fÃ¼r Naturkunde Berlin. Neptune Deep-	8/19/18	BIOCASE
3c001217-eea8-4f59-8b28-885699f8cd6c	Ohio University Invertebrate Paleontology	Ohio U http://creativecommons.org/publicdomain/zi	US	2412	Stigall A (2016). Ohio University Invertebrate Pale	8/19/18	DWC_ARCHIVE
bea28c6b-4282-4e0e-894d-7c65d050ffa9	NCSM Invertebrate Paleontology Collectio	North http://creativecommons.org/publicdomain/zi	US	11877	Norton B (2016). NCSM Invertebrate Paleontology	8/19/18	DWC_ARCHIVE
0ec927cf-325a-4d63-9499-d721c734463a	LACM Entomology Collection	Natur: http://creativecommons.org/publicdomain/zi	US	184364	Mertz B (2018). LACM Entomology Collection. Ver	8/19/18	DWC_ARCHIVE
b25c3e4d-4742-4d1a-a8cb-4b7648bace74	CSIRO Survey TT200801 - ROV Jason cruise	CSIRO http://creativecommons.org/licenses/by/4.0/	AU	1353	Althaus F, Watts D (2017). CSIRO Survey TT20080	8/19/18	DWC_ARCHIVE
b78dc638-1dde-43b4-a4cb-1e7e8ba2f185	ColecciÃ³n de moldes endocraneanos del	CCT Cí http://creativecommons.org/licenses/by/4.0/	AR	146	Dozo T (2015). ColecciÃ³n de moldes endocranea	8/19/18	DWC_ARCHIVE
ee789ae4-ef51-4ff2-931b-bc61b2dbe40e	Museo Argentino de Ciencias Naturales "B	Musei http://creativecommons.org/licenses/by/4.0/	AR	Key Not Found	Kramarz A, RodrÃ­guez D (2018). Museo Argentin	8/19/18	DWC_ARCHIVE
abe3af3c-07e8-4cea-b21c-93382834f490	Collection scientifique et pÃ©dagogique d	Ecole http://creativecommons.org/licenses/by/4.0/	FR	581	Moussus J (2014). Collection scientifique et pÃ©c	8/19/18	DWC_ARCHIVE
048d1f67-f86c-441d-b6f9-04b1b464e146	Museo Argentino de Ciencias Naturales "B	Musei http://creativecommons.org/licenses/by/4.0/	AR	Key Not Found	Kramarz A, RodrÃ­guez D (2018). Museo Argentin	8/19/18	DWC_ARCHIVE
463a3bd9-a9bf-4667-ae34-44b023307367	Museo Argentino de Ciencias Naturales "B	Musei http://creativecommons.org/licenses/by/4.0/	AR	Key Not Found	Totah V, RodrÃ­guez D (2018). Museo Argentino c	8/19/18	DWC_ARCHIVE
7814596d-f762-437d-abda-a2e5ad64df92	Museo Argentino de Ciencias Naturales "B	Musei http://creativecommons.org/licenses/by/4.0/	AR	Key Not Found	Del Fuevo G, RodrÃ­guez D (2018). Museo Argentin	8/19/18	DWC_ARCHIVE

Meta.xml

1. Open meta.xml file for each DwC-A
2. Read terms and assign "X" for term presence
3. Count total number of terms
4. Count number of extensions
5. Record extension names
6. Append results to CSV

Key elements:

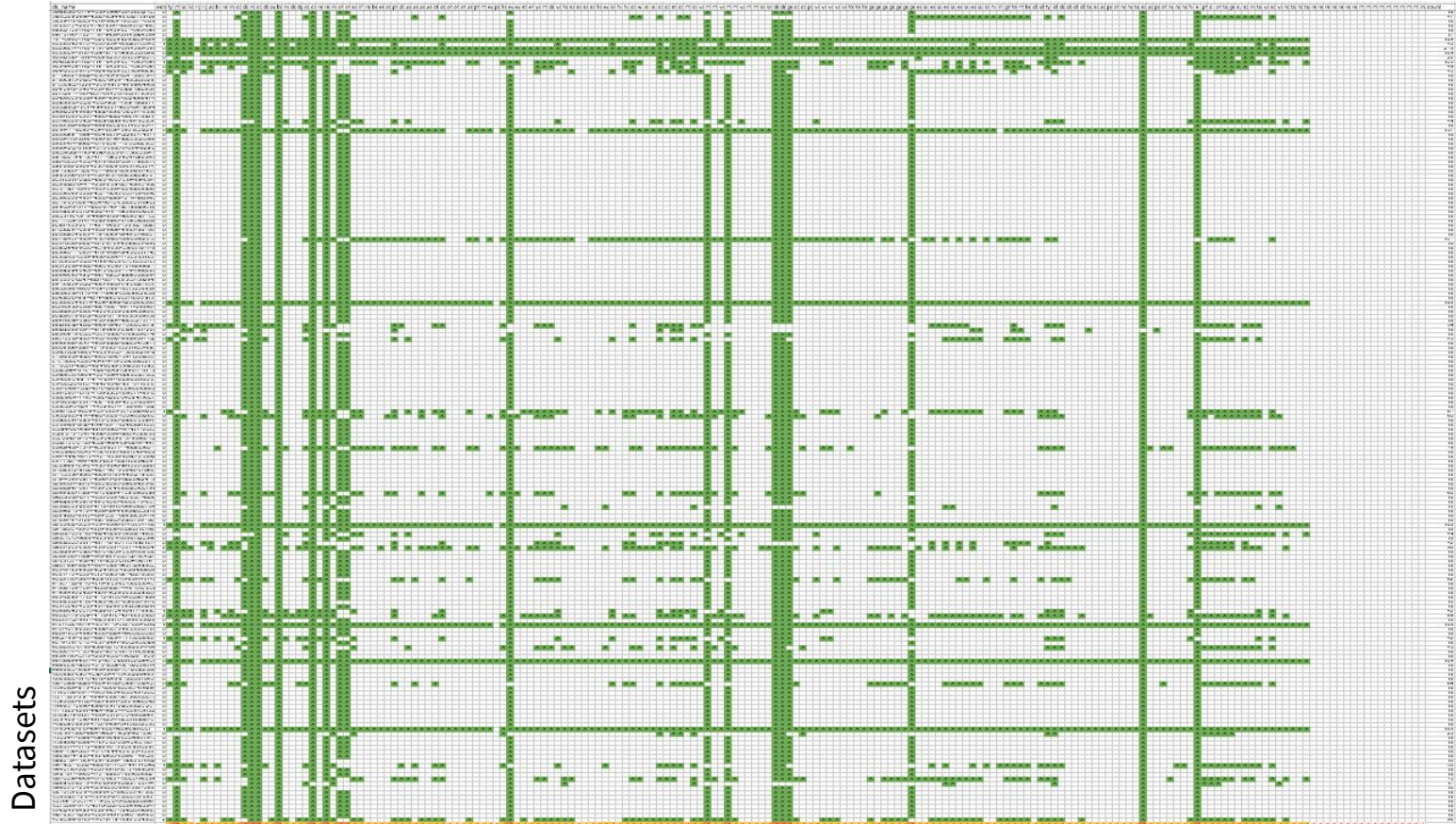
- Some terms not checked
- Lowercase all terms
- Does not account for poor formatting
- Removed ~900 rows for variations on the same dataset

```
<archive xmlns="http://rs.tdwg.org/dwc/text/" metadata="metadata.xml">
  <core encoding="UTF-8" fieldsTerminatedBy="\t" linesTerminatedBy="\n" fieldsEn
  <files>
    <location>occurrence.txt</location>
  </files>
  <id index="0" />
  <field index="0" term="http://rs.qbif.org/terms/1.0/qbifID"/>
  <field index="1" term="http://purl.org/dc/terms/abstract"/>
  <field index="2" term="http://purl.org/dc/terms/accessRights"/>
  <field index="3" term="http://purl.org/dc/terms/accrualMethod"/>
  <field index="4" term="http://purl.org/dc/terms/accrualPeriodicity"/>
  <field index="5" term="http://purl.org/dc/terms/accrualPolicy"/>
  <field index="6" term="http://purl.org/dc/terms/alternative"/>
  <field index="7" term="http://purl.org/dc/terms/audience"/>
  <field index="8" term="http://purl.org/dc/terms/available"/>
  <field index="9" term="http://purl.org/dc/terms/bibliographicCitation"/>
  <field index="10" term="http://purl.org/dc/terms/conformsTo"/>
```

```
dwc-a_r.p
11
12 FP_URI = 'dwc-a_xml2.csv'
13 rows = []
14 for folder in directory_list:
15     FP_XML = folder+ '/meta.xml'
16
17     tree = ET.parse(FP_XML)
18     root = tree.getroot()
19     root_len = len(root)
20     print(root_len)
21     terms = {}
22     for child in root[0]:
23         child = child.attrib
24         #each child is returned as a dictionary with index and term k
25         for key, value in child.items():
26             if key == "term":
27                 term = value.rsplit("/",1)[1].lower()
28
29
30
31
32     for col in uri_list[0]:
33         try:
34             newrow.append(terms[col.lower()])
35         except KeyError:
36             newrow.append('')
37     rows.append(newrow)
38
39 with open(FP_URI, 'a', newline = '') as f: #figure out what newl
40     writer = csv.writer(f)
41     for newrow in rows:
42         writer.writerow(newrow)
```

Total datasets: 244

Average # terms: 32



DwC Terms (In all charts terms are ordered by class and as listed in DwC specifications)

occurrence.txt (~~.tab~~)

1. Load .txt file into dataframe
2. Count number of usages per term
3. Calculate percentage
4. Append results to CSV

Key elements:

- Cannot process .tab files
- Lowercase all terms
- Skip lines with poorly formatted data (negligible)
- Probably includes none paleo data
- Memory errors
 - Run datasets in smaller batches

```
dwc-a_r.py
dwc-a_xml.py
1 import pandas as pd
2 import numpy as np
3 import os
4 import csv
5
6
7 directory_list = []
8 path = r"Volumes/HL_64_USB3/Paleo DwC/"
9 for child in os.listdir(path):
10     if os.path.isdir(child):
11         dirpath = os.path.join(path, child)
12         directory_list.append(dirpath)
13
14 usages = []
15 for folder in directory_list:
16     print(folder)
17     usage = {}
18     try:
19         occurrence = pd.read_table(folder+"/occurrence.txt", dtype = 'object', error
20 row_total = occurrence.shape[0]
21 print(row_total)
22
23 usage['total_count'] = row_total
24 usage['set_name'] = os.path.basename(folder) #metadata to identify each
25 usage['institution_Code'] = occurrence['institutionCode'][1]
26
27 for column in occurrence: #loop to count number of times a term is used
28     counts = occurrence[column].count()
29     percentage = counts/row_total*100
30 # usage[column] = counts
31 usage[column.lower()] = percentage
32
33 usages.append(usage)
34
35 except FileNotFoundError:
36     usage['set_name'] = os.path.basename(folder)
37     usage['error'] = 'File Not Found'
38     usages.append(usage)
39
40 keys = []
41 keys = ['set_name', 'institution_Code', 'error', 'total_count']
42 with open('dwc_terms.csv', 'r') as f:
43     for line in f:
44         keys.append(line.strip().lower())
45
46 #with open('occurrence_count.csv', 'a', newline='') as f:
47 with open('occurrence_percent.csv', 'a', newline='') as f:
48     writer = csv.writer(f)
49 # writer.writerow(keys)
50 for usage in usages:
51     row = [usage.get(key, "") for key in keys]
52     writer.writerow(row)
```

Total datasets: 101

Percentages

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U	V	W	X	Y	Z	AA	AB	AC	AD	AE	AF	AG	AH	AI	AJ	AK	AL	AM	AN	AO	AP	AQ	AR	AS	AT	AU	AV	AW	AX	AY	AZ	BA	BB	BC	BD	BE	BF	BG	BH	BI	BJ	BK	BL	BM	BN	BO	BP	BQ	BR	BS	BT	BV	BW	BX	BY	BZ	CA	CB	CC	CD	CE	CF	CG	CH	CI	CJ	CK	CL	CM	CN	CO	CP	CQ	CR	CS	CT	CU	CV	CW	CX	CY	CZ	DA	DB	DC	DD	DE	DF	DG	DH	DI	DJ	DK	DL	DM	DN	DO	DP	DQ	DR	DS	DT	DV	DW	DX	DY	DZ	EA	EB	EC	ED	EE	EF	EG	EH	EI	EJ	EK	EL	EM	EN	EO	EP	EQ	ER	ES	ET	EJ	EV	EW	EX	EY	EZ	FA	FB	FC	FD	FE	FF	FG	FH	FI	FJ	FK	FL	FM	FN	FO	FP	FQ	FR	FS	FT	FV	FW	FX	FY	FZ	GA	GB	GC	GD	GE	GF	GG	GH	GI	GJ	GK	GL	GM	GN	GO	GP	GQ	GR	GS	GT	GV	GW	GX	GY	GZ	HA	HB	HC	HD	HE	HF	HG	HH	HI	HJ	HK	HL	HM	HN	HO	HP	HQ	HR	HS	HT	HV	HW	HX	HY	HZ	IA	IB	IC	ID	IE	IF	IG	IH	II	IJ	IK	IL	IM	IN	IO	IP	IQ	IR	IS	IT	IV	IW	IX	IY	IZ	JA	JB	JC	JD	JE	JF	JG	JH	JI	IJ	JK	JL	JM	JN	JO	JP	JQ	JR	JS	JT	JV	JW	JX	JY	JZ	KA	KB	KC	KD	KE	KF	KG	KH	KI	KJ	KK	KL	KM	KN	KO	KP	KQ	KR	KS	KT	KV	KW	KX	KY	KZ	LA	LB	LC	LD	LE	LF	LG	LH	LI	LJ	LK	LL	LM	LN	LO	LP	LQ	LR	LS	LT	LV	LW	LX	LY	LZ	MA	MB	MC	MD	ME	MF	MG	MH	MI	MJ	MK	ML	MM	MN	MO	MP	MQ	MR	MS	MT	MV	MW	MX	MY	MZ	NA	NB	NC	ND	NE	NF	NG	NH	NI	NJ	NK	NL	NM	NN	NO	NP	NQ	NR	NS	NT	NV	NW	NX	NY	NZ	OA	OB	OC	OD	OE	OF	OG	OH	OI	OJ	OK	OL	OM	ON	OO	OP	OQ	OR	OS	OT	OV	OW	OX	OY	OZ	PA	PB	PC	PD	PE	PF	PG	PH	PI	PJ	PK	PL	PM	PN	PO	PQ	PR	PS	PT	PV	PW	PX	PY	PZ	QA	QB	QC	QD	QE	QF	QG	QH	QI	QJ	QK	QL	QM	QN	QO	QP	QQ	QR	QS	QT	QV	QW	QX	QY	QZ	RA	RB	RC	RD	RE	RF	RG	RH	RI	RJ	RK	RL	RM	RN	RO	RP	RQ	RR	RS	RT	RV	RW	RX	RY	RZ	SA	SB	SC	SD	SE	SF	SG	SH	SI	SJ	SK	SL	SM	SN	SO	SP	SQ	SR	SS	ST	SV	SW	SX	SY	SZ	TA	TB	TC	TD	TE	TF	TG	TH	TI	TJ	TK	TL	TM	TN	TO	TP	TQ	TR	TS	TV	TW	TX	TY	TZ	UA	UB	UC	UD	UE	UF	UG	UH	UI	UJ	UK	UL	UM	UN	UO	UP	UQ	UR	US	UT	UV	UW	UX	UY	UZ	VA	VB	VC	VD	VE	VF	VG	VH	VI	VJ	VK	VL	VM	VN	VO	VP	VQ	VR	VS	VT	VV	VW	VX	VY	VZ	WA	WB	WC	WD	WE	WF	WG	WH	WI	WJ	WK	WL	WM	WN	WO	WP	WQ	WR	WS	WT	WV	WW	WX	WY	WZ	XA	XB	XC	XD	XE	XF	XG	XH	XI	XJ	XK	XL	XM	XN	XO	XP	XQ	XR	XS	XT	XV	XW	XX	XY	XZ	YA	YB	YC	YD	YE	YF	YG	YH	YI	YJ	YK	YL	YM	YN	YO	YP	YQ	YR	YS	YT	YV	YW	YX	YZ	ZA	ZB	ZC	ZD	ZE	ZF	ZG	ZH	ZI	ZJ	ZK	ZL	ZM	ZN	ZO	ZP	ZQ	ZR	ZS	ZT	ZV	ZW	ZX	ZY	ZZ
--	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----	----

Total datasets: 101

Percentages



Event

Location

GeoRef

Geo Context

Taxonomy

Presence

Top 10% (~200 usages)

modified
datasetID
institutionCode
collectionCode
basisOfRecord
catalogNumber
recordedBy
organismQuantity
organismQuantityType
eventDate
minimumElevationInMeters
minimumDepthInMeters
decimalLatitude
decimalLongitude
geodeticDatum
geologicalContextID
scientificName
kingdom

Percentages

Top 10% (>75%)

type
modified
language
institutionCode
collectionCode
basisOfRecord
occurrenceID
catalogNumber
countryCode
stateProvince
verbatimLocality
minimumElevationInMeters
MinimumDepthInMeters
namePublishedIn
subgenus
relatedResourceID

Counts

Top 10% (>950,000)

modified
institutionCode
collectionCode
basisOfRecord
occurrenceID
catalogNumber
stateProvince
verbatimLocality
minimumElevationInMeters
MinimumDepthInMeters
namePublishedIn
subgenus
specificEpithet
infraspecificEpithet
taxonRank
verbatimTaxonRank
scientificNameAuthorship

Presence

Top 10% (~200 usages)

modified
datasetID
institutionCode
collectionCode
basisOfRecord
catalogNumber
recordedBy
organismQuantity
organismQuantityType
eventDate
minimumElevationInMeters
minimumDepthInMeters
decimalLatitude
decimalLongitude
geodeticDatum
geologicalContextID
scientificName
kingdom

Percentages

Top 10% (>75%)

type
modified
language
institutionCode
collectionCode
basisOfRecord
occurrenceID
catalogNumber
countryCode
stateProvince
verbatimLocality
minimumElevationInMeters
MinimumDepthInMeters
namePublishedIn
subgenus
relatedResourceID

Percent with Fossil filter

Top 10% (of 6,152,144 records)

type
modified
institutionCode
collectionCode
basisOfRecord
occurrenceID
catalogNumber
stateProvince
verbatimLocality
minimumElevationInMeters
MinimumDepthInMeters
namePublishedIn
Subgenus
specificEpithet
infraspecificEpithet
scientificNameAuthorship

Occurrence Records

Occurrences: 10,042,739

Tool: SQL*
BigQuery*

Sources: Occurrence download

**Processing by Luis Villanueva, SI DPO*

Usage Counts

```
filename="columns_verbatim.csv"
while read -r line
do
    echo "$line"
    echo "$line" >> results_verbatim.txt

    echo "insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select '${line//[\'\t\r\n ']}', 'is_null', cast(count(gbifID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE
    bq query --use_legacy_sql=false "insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select '${line//[\'\t\r\n ']}', 'is_null', cast(count(gbifID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE
    echo "insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select '${line//[\'\t\r\n ']}', 'no_distinct', cast(count(distinct ${line//[\'\t\r\n ']}') as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE
    bq query --use_legacy_sql=false "insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select '${line//[\'\t\r\n ']}', 'no_distinct', cast(count(distinct ${line//[\'\t\r\n ']}') as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE

done < "$filename"
```

```
eventID
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'eventID', 'is_null', cast(count(gbifID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE eventID IS NULL)
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'eventID', 'no_distinct', cast(count(distinct eventID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE eventID IS NOT NULL)
parentEventID
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'parentEventID', 'is_null', cast(count(gbifID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE parentEventID IS NULL)
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'parentEventID', 'no_distinct', cast(count(distinct parentEventID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE parentEventID IS NOT NULL)
fieldNumber
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'fieldNumber', 'is_null', cast(count(gbifID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE fieldNumber IS NULL)
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'fieldNumber', 'no_distinct', cast(count(distinct fieldNumber) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE fieldNumber IS NOT NULL)
eventDate
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'eventDate', 'is_null', cast(count(gbifID) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE eventDate IS NULL)
insert into gbif_paleo_20180718.gbif_results_verbatim (colname, querytype, value) (select 'eventDate', 'no_distinct', cast(count(distinct eventDate) as string) FROM gbif_paleo_20180718.gbif_verbatim WHERE eventDate IS NOT NULL)
```

Top 10%

basisOfRecord	is_null	0	0.00%	no_dist	7
collectionCode	is_null	1,852	0.02%	no_dist	182,414
institutionCode	is_null	432,376	3.95%	no_dist	119
scientificName	is_null	609,520	5.57%	no_dist	797,610
identifier	is_null	1,078,133	9.85%	no_dist	7,542,714
catalogNumber	is_null	2,275,400	20.78%	no_dist	6,321,787
country	is_null	2,797,005	25.55%	no_dist	1,644
kingdom	is_null	2,891,281	26.41%	no_dist	52
occurrenceID	is_null	3,298,333	30.13%	no_dist	6,734,061
recordedBy	is_null	3,783,676	34.56%	no_dist	80,189
modified	is_null	3,810,201	34.80%	no_dist	662,352
genus	is_null	4,196,477	38.33%	no_dist	108,006
decimalLatitude	is_null	4,488,154	41.00%	no_dist	117,132
decimalLongitude	is_null	4,488,191	41.00%	no_dist	127,247
class	is_null	4,530,641	41.38%	no_dist	845
phylum	is_null	4,579,874	41.83%	no_dist	392
continent	is_null	4,777,000	43.63%	no_dist	95
stateProvince	is_null	4,780,674	43.67%	no_dist	8,055
locality	is_null	4,968,186	45.38%	no_dist	353,130
geodeticDatum	is_null	5,153,474	47.07%	no_dist	66
order	is_null	5,188,502	47.39%	no_dist	3,173
family	is_null	5,225,724	47.73%	no_dist	15,098
type	is_null	5,553,827	50.73%	no_dist	8
specificEpithet	is_null	5,795,788	52.94%	no_dist	149,375
higherGeography	is_null	6,120,647	55.91%	no_dist	32,384
language	is_null	6,254,868	57.13%	no_dist	7
formation	is_null	6,520,756	59.56%	no_dist	31,500
higherClassification	is_null	6,673,749	60.96%	no_dist	123,029
nomenclaturalCode	is_null	6,833,749	62.42%	no_dist	6
taxonRank	is_null	6,868,427	62.74%	no_dist	93
county	is_null	6,975,852	63.72%	no_dist	11,063
datasetName	is_null	7,082,686	64.69%	no_dist	176,268
license	is_null	7,092,843	64.79%	no_dist	106
accessRights	is_null	7,150,863	65.32%	no_dist	9
earliestPeriodOrLowest	is_null	7,354,179	67.17%	no_dist	473

earliestAgeOrLowestS	is_null	7,671,117	70.07%	no_dist	2,576
individualCount	is_null	7,921,337	72.35%	no_dist	865
preparations	is_null	7,933,821	72.47%	no_dist	87,058
verbatimLatitude	is_null	7,944,984	72.57%	no_dist	62,871
verbatimLongitude	is_null	7,944,984	72.57%	no_dist	74,637
previousIdentification	is_null	8,004,903	73.12%	no_dist	438,445
eventDate	is_null	8,015,021	73.21%	no_dist	27,157
earliestEpochOrLowest	is_null	8,081,951	73.82%	no_dist	962
dynamicProperties	is_null	8,178,769	74.71%	no_dist	190,553
verbatimLocality	is_null	8,234,354	75.21%	no_dist	49,462
countryCode	is_null	8,262,288	75.47%	no_dist	450
georeferenceRemarks	is_null	8,286,071	75.69%	no_dist	14,702
verbatimTaxonRank	is_null	8,302,772	75.84%	no_dist	29
associatedReferences	is_null	8,334,040	76.12%	no_dist	56,143
lithostratigraphicTerm	is_null	8,366,550	76.42%	no_dist	82,766
occurrenceStatus	is_null	8,393,663	76.67%	no_dist	5
collectionID	is_null	8,449,259	77.18%	no_dist	180,652
georeferenceVerificat	is_null	8,594,634	78.50%	no_dist	11
taxonID	is_null	8,649,671	79.01%	no_dist	185,664
samplingProtocol	is_null	8,701,708	79.48%	no_dist	3,036
coordinatePrecision	is_null	8,714,138	79.60%	no_dist	30
namePublishedInYear	is_null	8,721,800	79.67%	no_dist	217
taxonConceptID	is_null	8,794,612	80.33%	no_dist	157,772
institutionID	is_null	8,954,946	81.80%	no_dist	23
occurrenceRemarks	is_null	8,990,319	82.12%	no_dist	772,952
earliestEraOrLowestE	is_null	9,035,302	82.53%	no_dist	100
rightsHolder	is_null	9,252,032	84.51%	no_dist	25
coordinateUncertain	is_null	9,252,925	84.52%	no_dist	10,163
datasetID	is_null	9,331,943	85.24%	no_dist	1,329
references	is_null	9,397,213	85.84%	no_dist	1,543,006
locationID	is_null	9,419,653	86.04%	no_dist	119,272
group	is_null	9,421,710	86.06%	no_dist	3,534
year	is_null	9,510,404	86.87%	no_dist	315
identifiedBy	is_null	9,580,177	87.51%	no_dist	10,404
latestPeriodOrHighest	is_null	9,746,247	89.02%	no_dist	109

Distinct Values

```
filename="distinct_terms.txt"
while read -r line
do
  echo "$line"
  echo "$line" >> results_verbatim.txt

  echo "insert into gbif_paleo_20180718.gbif_results_verbatim_counts_datasets (colname, value, no_records, datasetName, title, InstitutionC
  bq query --use_legacy_sql=false "insert into gbif_paleo_20180718.gbif_results_verbatim_counts_datasets (colname, value, no_records, data
done < "$filename"
```

type
language
license
basisOfRecord
sex
lifeStage
reproductiveCondition
behavior
establishmentMeans
occurrenceStatus
preparations
disposition
continent
waterBody
islandGroup
island
countryCode
verbatimSRS
geodeticDatum
georeferenceVerificationStatus
typeStatus
taxonRank
nomenclaturalCode
taxonomicStatus

Controlled Vocabularies

colname	querytype	value	% of rows		
basisOfRecord	is_null	0	0.00%	no_distinct	7
country	is_null	2,797,005	25.55%	no_distinct	1,644
kingdom	is_null	2,891,281	26.41%	no_distinct	52
class	is_null	4,530,641	41.38%	no_distinct	845
phylum	is_null	4,579,874	41.83%	no_distinct	392
continent	is_null	4,777,000	43.63%	no_distinct	95
stateProvince	is_null	4,780,674	43.67%	no_distinct	8,055
geodeticDatum	is_null	5,153,474	47.07%	no_distinct	66
family	is_null	5,225,724	47.73%	no_distinct	15,098
type	is_null	5,553,827	50.73%	no_distinct	8
taxonRank	is_null	6,868,427	62.74%	no_distinct	93
county	is_null	6,975,852	63.72%	no_distinct	11,063
preparations	is_null	7,933,821	72.47%	no_distinct	87,058
earliestEpochOrLowestSeries	is_null	8,081,951	73.82%	no_distinct	962

Preparations: 87,058 values

Important Paleo information, but we use it for multiple purposes:

- Material Type
- Prep work done
 - By % or action
- Anatomy/Morphology
- Internal notations of work done

value	no_records	datasetName
fossilized	421324	
default - 1	358025	
unknown (fossil)	183683	
PREP - 1	139514	
fossil	86342	cuinvert
Fossil - 1	84565	University of Kansas Biodive
Dry - 1	63153	
Permineralization - 1	61484	University of Kansas Biodive
Cataloged - 1	55998	
fossil	49369	Fossil Vertebrate
fossil - 1	44262	
PREP - 2	40680	
Secondary microslides	33016	NMNH Paleobiology
shell	31780	
fossil - 1	28795	
Microslide	23021	NMNH Paleobiology
PREP - 3	18047	
Primary microslides	17088	NMNH Paleobiology
Compression/Impression - 1	16514	University of Kansas Biodive
100% prepped - 2	1741	
shell	1727	
tooth, molar	1721	
Dry - 5	1720	
humerus	1686	
Embedded	698	
SHELL FRAGS	697	
CALCANEUM	697	
?	697	
Tooth; incomplete	693	
tooth fragment	685	
HÄmer D	685	
Dry - 20	679	

Presence

Top 10% (~200 usages)

modified
datasetID
institutionCode
collectionCode
basisOfRecord
catalogNumber
recordedBy
organismQuantity
organismQuantityType
eventDate
minimumElevationInMeters
minimumDepthInMeters
decimalLatitude
decimalLongitude
geodeticDatum
geologicalContextID
scientificName
kingdom

Percentages

Top 10% (>75%)

type
modified
language
institutionCode
collectionCode
basisOfRecord
occurrenceID
catalogNumber
countryCode
stateProvince
verbatimLocality
minimumElevationInMeters
MinimumDepthInMeters
namePublishedIn
subgenus
relatedResourceID

Counts

Top 10% (>950,000)

modified
institutionCode
collectionCode
basisOfRecord
occurrenceID
catalogNumber
stateProvince
verbatimLocality
minimumElevationInMeters
MinimumDepthInMeters
namePublishedIn
subgenus
specificEpithet
infraspecificEpithet
taxonRank
verbatimTaxonRank
scientificNameAuthorship

Occurrence Counts

Top 10%

basisOfRecord
collectionCode
institutionCode
scientificName
identifier
catalogNumber
country
kingdom
occurrenceID
recordedBy
modified
genus
decimalLatitude
decimalLongitude
class
phylum
continent
stateProvince
Locality
geodeticDatum
order

Problems

higherTaxonomy - only 30 datasets or 40% of records - Important for Paleo data (CLADES)

taxonRank - only 41 datasets or 38% of records (verbatimTaxonRank = 25%) with 93 distinct values

Preparations - only 47 datasets or 30% of records with 87,058 distinct values

We talk about it a lot, but aren't serving a lot of data and the data we do serve is very variable. Is it because we're not sure what to share?

Extensions - Barely using extensions (14 with multimedia, only 9 other extensions across all datasets)

Fixing NMNH Paleo

DwC Terms	USN	USNM	All datasets	percentage	All occurrences		
eventDate			99.499221	is_ 8,015,021	73.21%	no_ 27,157	
eventTime			#DIV/0!	is_ 10,903,619	99.60%	no_ 690	
startDayOfYear	X	0.323	47.408504	is_ 10,590,617	96.74%	no_ 367	
endDayOfYear	X		64.269201	is_ 10,746,718	98.16%	no_ 366	
year	X		53.19237	is_ 9,510,404	86.87%	no_ 315	
month	X	21.57	28.515609	is_ 10,097,654	92.23%	no_ 30	
day	X	21.57	34.120706	is_ 4,488,154	41.00%	no_ 117,132	
verbatimEventDate	X	37.61	59.913649	is_ 9,995,065	91.30%	no_ 33,090	
habitat	X	21.69	48.135115	is_ 10,945,077	99.97%	no_ 73	
fieldNumber	X	18.49	45.138677	is_ 10,454,413	95.49%	no_ 147,734	
fieldNotes	X	36.94	47.701501	is_ 10,947,675	100.00%	no_ 1	
eventRemarks			43.295624	is_ 10,915,548	99.70%	no_ 2,090	
locationID	X		58.23435	is_ 9,419,653	86.04%	no_ 119,272	
higherGeographyID			#DIV/0!	is_ 10,945,961	99.98%	no_ 3	
higherGeography	X		#DIV/0!	is_ 6,120,647	55.91%	no_ 32,384	
continent	X		0.9233132	is_ 4,777,000	43.63%	no_ 95	
waterBody	X		66.666743	is_ 10,844,516	99.06%	no_ 282	
islandGroup	X		29.616996	is_ 10,902,702	99.59%	no_ 65	
island	X	50.77	77.392119	is_ 10,897,567	99.54%	no_ 319	
country	X		100	is_ 2,797,005	25.55%	no_ 1,644	
countryCode		78.29	94.347025	is_ 8,262,288	75.47%	no_ 450	
stateProvince	X	0.672	92.594589	is_ 4,780,674	43.67%	no_ 8,055	
county	X	3.682	22.211235	is_ 6,975,852	63.72%	no_ 11,063	
municipality		0.008	10.96424	is_ 10,409,285	95.08%	no_ 17,303	
locality	X	0.006	2.1544214	is_ 4,968,186	45.38%	no_ 353,130	
verbatimLocality		74.91	93.117818	is_ 8,234,354	75.21%	no_ 49,462	

Conclusions: Community Effort

We all have data problems

There is A LOT of variation in how we map our data

We are underutilizing what is there

Definitely new terms needed

Some easy fixes/cleanup - documentation with the recommended values and formats

Citations

Wieczorek J, Bloom D, Guralnick R, Blum S, Döring M, Giovanni R, et al. (2012) Darwin Core: An Evolving Community-Developed Biodiversity Data Standard. PLoS ONE 7(1): e29715.

<https://doi.org/10.1371/journal.pone.0029715>

GBIF.org (09 August 2018) GBIF Occurrence Download <https://doi.org/10.15468/dl.r3tqg5>

GBIF.org (18 July 2018) GBIF Occurrence Download <https://doi.org/10.15468/dl.k61spa>

Ask me if you'd like the full citation list for all datasets!

Thank you

Luis Villanueva, Smithsonian
Digitization Program Office

Adam Mansur, Smithsonian NMNH
Department of Mineral Sciences

Stack Overflow

GBIF

All of the data providers!



EPICC collections
network working on an
Initial Paleo Guideline