# Webinar Series

## Data Use Skills
## Featuring Data from Natural History Collections
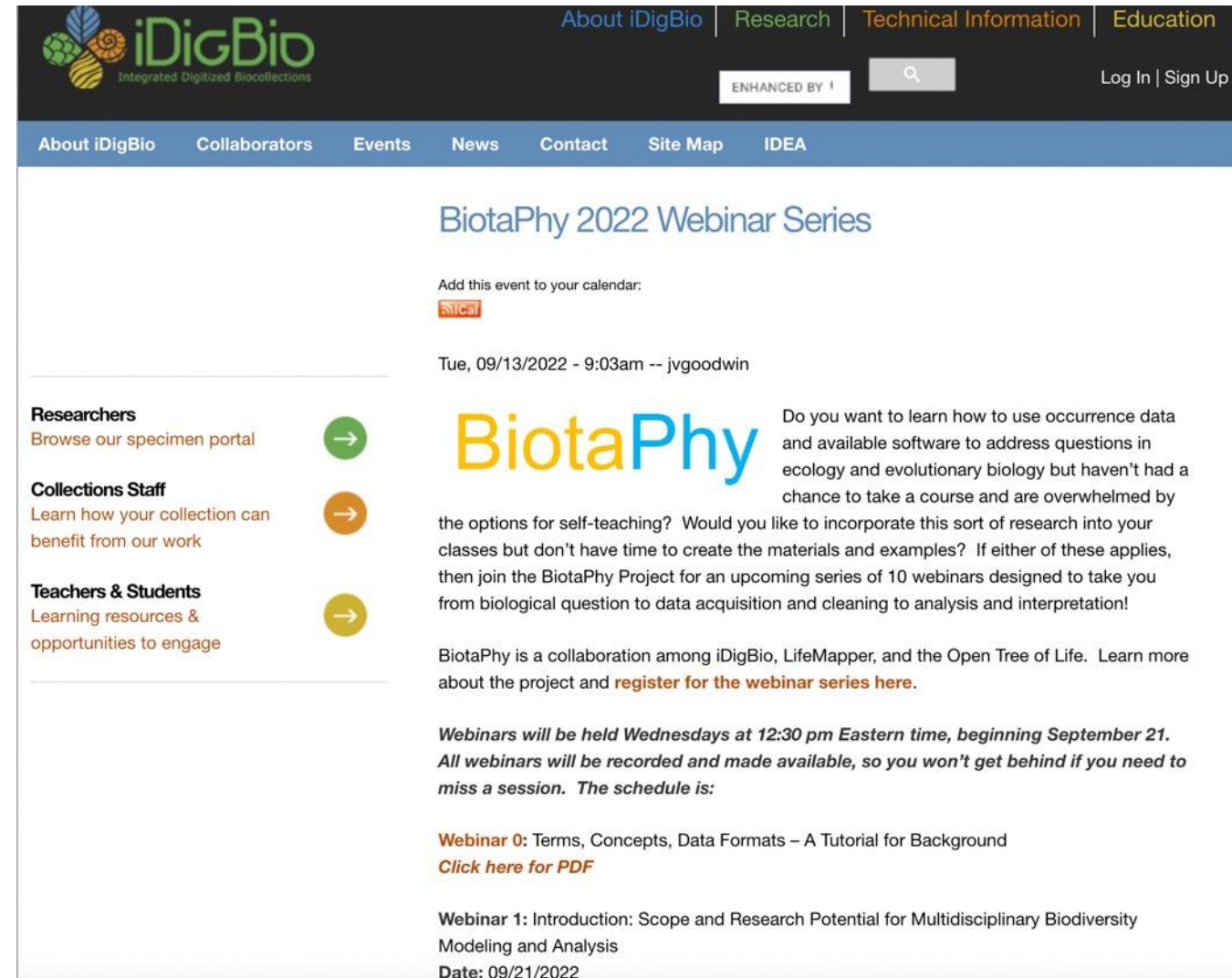
# September 21-November 30, 2022

https://www.idigbio.org/content/biotaphy-2022-webinar-series

# iDigBio:

https://www.idigbio.org/content/biotaphy-2022-webinar-series

# iDigBio.org

# Thank You

## Maria Cortez
## Aimee Stewart

## Jill Goodwin
## Gil Nelson

# Webinar 5

## Big Data Munging:
## Finding, Acquiring, and Preparing Species Occurrence Data and Tree Data

# Goals

**Learn how to find, manipulate, combine, and use occurrence and tree data in biodiversity analyses**

# Learning Objectives

**Biological Objectives:**

✓ **Introduce different data sources (iDigBio, GBIF, SpeciesLink, RainBio, OpenTree, etc.) and data types (occurrence and tree)**

✓ **Showcase options for downloading data directly from portals**

# Learning Objectives

**Technological Objectives**:

✓ **Group retrieved data by taxonomy or other field type**.


✓ **Learn how to combine records from heterogeneous sources**.

# Webinar organization

1. **Exploring Concepts**: why do we need occurrence and tree data, and where can we find them?

2. **Demonstrations**: how to download occurrence and tree data directly from portals.

3. **Exercises**: practice automated ways of downloading and treating occurrence and tree data.

4. **Session Summary, Q&A, and Discussion**

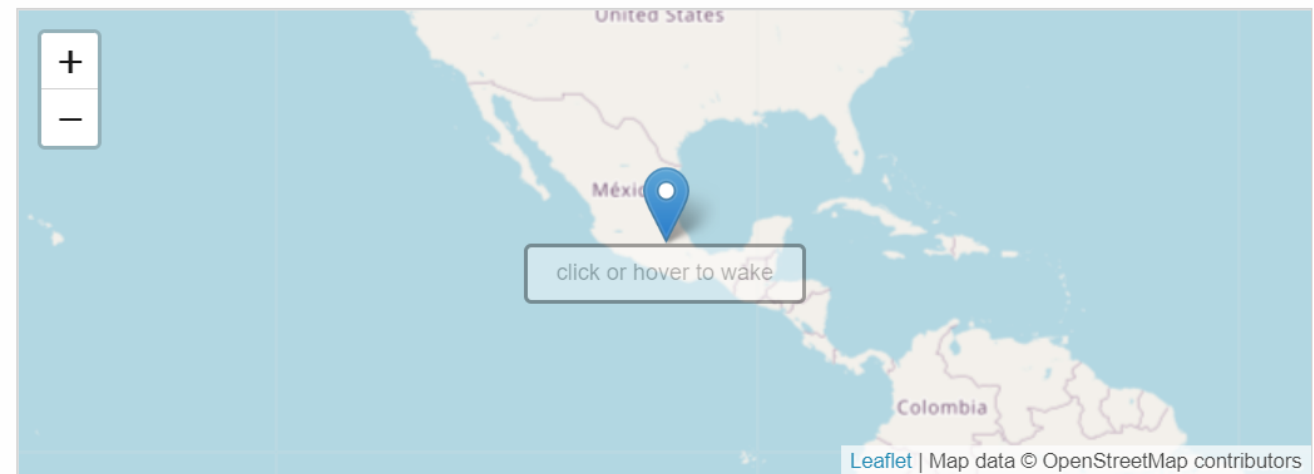**Why do we need occurrence and tree data and where can we find them?**

# Exploring Concepts

**Occurrence data are essential for producing species distribution models, estimating phylogenetic diversity, and more!**

### *Heuchera mexicana* W. Schaffn. ex Small & Rydb., 1905

From Computarización del Herbario ENCB, IPN. Fase IV. Base de datos de la familia Pinaceae y de distintas familias de la clase Magnoliopsida depositadas en el Herbario de la Escuela Nacional de Ciencias Biológicas, IPN

| | | | |
|---|---|---|---|
| Continent | North America | Institution Code | Encb-ipn |
| Country | Mexico | Collection Code | Encb |
| State/Province | Mexico | Catalog Number | No Disponible |
| County/Parish | Ecatzingo | Collected By | Rosario Vázquez |
| Locality | Campamento Tlamacas, Amecameca | Date Collected | 1962-05-29 |
| Latitude | 19.065 | | |
| Longitude | -98.6336111 | | |

**Lat: -19.065**

**Lon: -98.6336111**

click or hover to wake

Leaflet | Map data © OpenStreetMap contributors

![iDigBio - Integrated Digitized Biocollections]

**Species distribution models**

*Prunus geniculata*

now

2050

**Phylogenetic diversity**

min= 0.16  max= 0.64

Species diversity
856
57
● Sampling location

Level III - Ecoregions
A - Southeastern Plains
B - Southern Coastal Plains
C - Southern Florida Coastal Plains

Chronogram

Phylogram

# Exploring Concepts



Tree data are essential for performing spatial phylogenetic analyses!

# Exploring Concepts

**Species distribution models**

*Prunus geniculata*

now

2050

**Phylogenetic diversity**

min= 0.16  max= 0.64

B

Species diversity
856
57
● Sampling location

Level III - Ecoregions
A - Southeastern Plains
B - Southern Coastal Plains
C - Southern Florida Coastal Plains

Chronogram

Phylogram

**Demonstration: how to download data directly from portals!**

- **Useful for smaller datasets**

- **Useful for demonstrations**

- **Useful as a teaching tool**

- **Easy to visualize the data**

# Demonstration: Downloading data from portals

# Demonstration: Downloading data from portals

# Demonstration: Downloading data from portals



*"...speciesLink promotes free and open access to data, information, and tools..."*
**most databases included are Brazilian**

# Demonstration: Downloading data from portals

*"...database contains high quality georeferenced occurrences of vascular plants from sub-Saharan tropical Africa."*

## The RAINBIO mega database

The major output of RAINBIO is the RAINBO mega database.

The RAINBIO mega database contains high quality georeferenced occurrences of vascular plants from sub-Saharan tropical Africa. It is a compilation of thirteen public and non-public databases made available under the RAINBIO project funded by CESAB. The database was filtered, quality-checked and verified by the CESAB RAINBIO Consortium. The database holds 610 117 georeferenced occurrences for 25,356 species of vascular plants and 29,659 taxa (including subspecies and varieties), 3,158 genera and 273 families. The database follows the Darwin Core standard. The RAINBIO database is subject to be updated in the future.

The RAINBIO database is available here

The database comes in two formats: a .csv file and and R.data project file. You can open the R.data file directly in R and use the available custom functions to extract useful data and produce distribution maps.

# Demonstration: Downloading data from portals

# Time to Exercise!

What happens when there is a large dataset? Should we download and treat 40,000 records individually?

We use BiotaPhy tools to automate occurrence and tree data munging!

# Time to Exercise!

Let's put the automated framework developed by BiotaPhy to the test!

How to split and merge occurrence data by species:

**3 steps:**

- ✓ **Data Preparation**

- ✓ **Run Tutorial**

- ✓ **Inspect Output**

Input: **occurrence records**

Input: **Wrangler configuration file**

Input: **Script parameter file**

## Data Preparation

### Input: occurrence records

The split_occurrence_data tool accepts one or more datasets, each must be either a Darwin Core Archive (DwCA) file or a CSV file containing records for one or more taxa.

More information is in the **Occurrence Data** section of data_wrangle_occurrence.

**BiotaPhy**

## Specimen Occurrences: Data and Wrangling

### Occurrence Data

Several tools (split_occurrence_data, wrangle_occurrences) accept occurrence data. The filename must be specified in the script parameter file, described in each tool's documentation and linked above. Data can be in one of two formats:

1. Darwin Core Archive (DwCA) file. DwCA files may be downloaded from several places, including GBIF and iDigBio.
   1. To download from GBIF, choose your filters in the GBIF portal. For example, the example data was downloaded after selecting occurrences where genus='Heuchera L' Then choose the download link at the upper right column header.
   2. To download from iDigBio, instructions for querying and downloading from the command prompt are at idigbio_download.
   3. The tutorial example DwCA is at occurrence_idigbio.zip

## Specimen Occurrences: Data and Wrangling

### Occurrence Data

Several tools (split_occurrence_data, wrangle_occurrences) accept occurrence data. The filename must be specified in the script parameter file, described in each tool's documentation and linked above. Data can be in one of two formats:

1. Darwin Core Archive (DwCA) file. DwCA files may be downloaded from several places, including GBIF and iDigBio.
    1. To download from GBIF, choose your filters in the GBIF portal. For example, the example data was downloaded after selecting occurrences where genus='Heuchera L' Then choose the download link at the upper right column header.
    2. To download from iDigBio, instructions for querying and downloading from the command prompt are at idigbio_download.
    3. The tutorial example DwCA is at occurrence_idigbio.zip

# Data Preparation: occurrence records

**BiotaPhy**

## Specimen Occurrences: Data and Wrangling

## Occurrence Data

Several tools (split_occurrence_data, wrangle_occurrences) accept occurrence data. The filename must be specified in the script parameter file, described in each tool's documentation and linked above. Data can be in one of two formats:

1. Darwin Core Archive (DwCA) file. DwCA files may be downloaded from several places, including GBIF and iDigBio.

   1. To download from GBIF, choose your filters in the GBIF portal. For example, the example data was downloaded after selecting occurrences where genus='Heuchera L.' Then choose the download link at the upper right column header.

   2. To download from iDigBio, instructions for querying and downloading from the command prompt are at idigbio_download.

   3. The tutorial example DwCA is at occurrence_idigbio.zip

# Data Preparation: occurrence records



## Download occurrence data from iDigBio

To download from iDigBio, full instructions are at the Download API reference.

To pull data from the command prompt, use the `curl` command to pull text response directly to terminal with the example query_url:
Euphorbia

```
$ curl https://api.idigbio.org/v2/download/?rq=%7B%22genus%22%3A%22euphorbia%22%7D&email=donotreply%40idigbio.org
[1] 58979
astewart@murderbot:~/git/tutorials$ {
  "complete": false,
  "created": "2022-05-02T15:28:41.730968+00:00",
  "expires": "2022-06-01T15:28:41.628063+00:00",
  "hash": "18911492e8517926cb8693fc9f971cf066107016",
  "query": {
    "core_source": "indexterms",
    "core_type": "records",
    "form": "dwca-csv",
    "mediarecord_fields": null,
    "mq": null,
    "record_fields": null,
    "rq": {
      "genus": "euphorbia"
    }
  },
  "status_url": "https://api.idigbio.org/v2/download/d54c0ad7-6697-4096-9f11-b2a9a6041a38",
  "task_status": "PENDING"
}
```

# Data Preparation: occurrence records

## Specimen Occurrences: Data and Wrangling

### Occurrence Data

Several tools (split_occurrence_data, wrangle_occurrences) accept occurrence data. The filename must be specified in the script parameter file, described in each tool's documentation and linked above. Data can be in one of two formats:

1. Darwin Core Archive (DwCA) file. DwCA files may be downloaded from several places, including GBIF and iDigBio.
    1. To download from GBIF, choose your filters in the GBIF portal. For example, the example data was downloaded after selecting occurrences where genus='Heuchera L' Then choose the download link at the upper right column header.
    2. To download from iDigBio, instructions for querying and downloading from the command prompt are at idigbio_download.
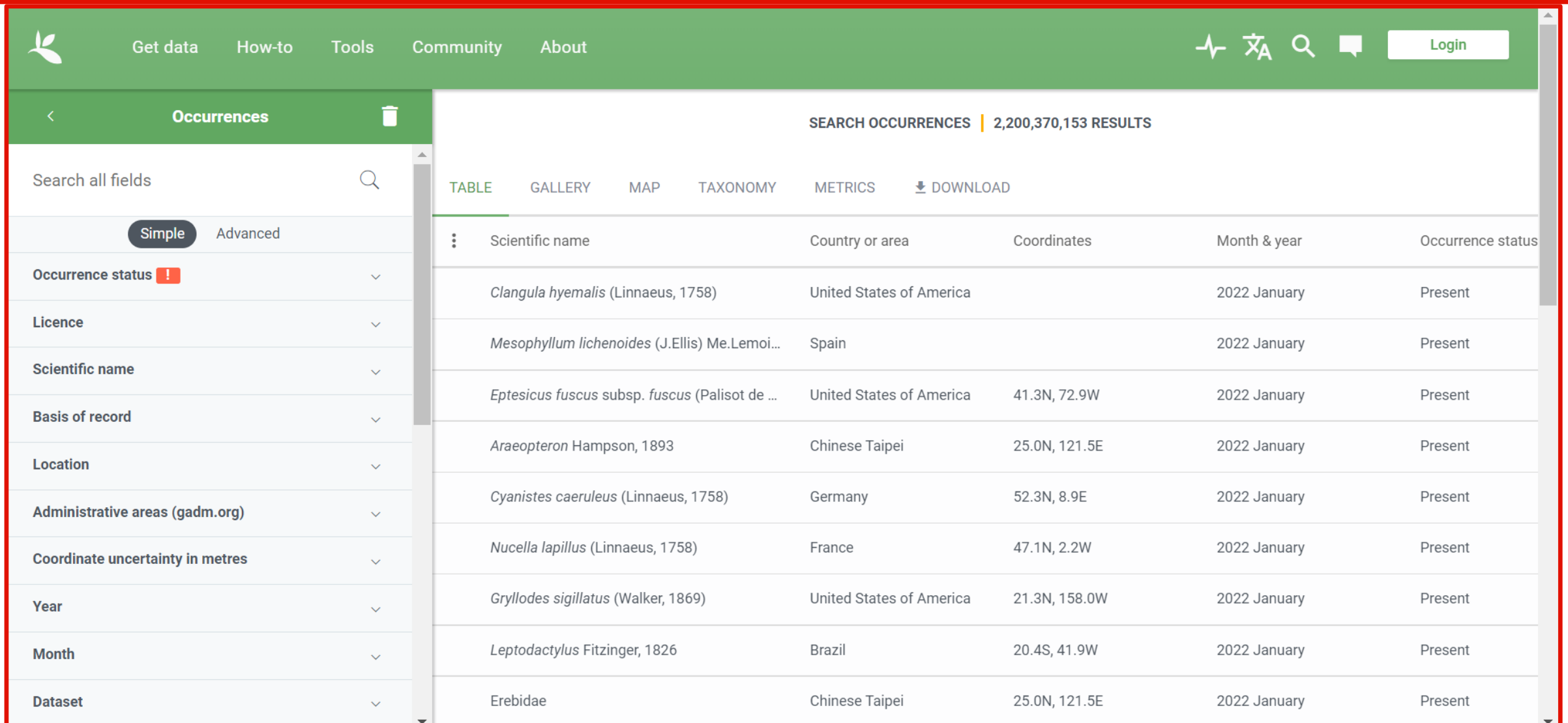    3. The tutorial example DwCA is at occurrence_idigbio.zip

# Data Preparation: occurrence records

| | | |
|---|---|---|
| grid.cpg | debugging | 3 months ago |
| grid.dbf | debugging | 3 months ago |
| grid.prj | debugging | 3 months ago |
| grid.qpj | debugging | 3 months ago |
| grid.shp | debugging | 3 months ago |
| grid.shx | debugging | 3 months ago |
| heuchera.csv | initial one-shot; unfinished | 5 months ago |
| heuchera.nex | data | last month |
| heuchera.txt | added pre-commit checks; unfinished | 4 months ago |
| heuchera3.nex | updates for matrix/tree match | 2 months ago |
| occurrence_gbif.csv | input data | 7 days ago |
| occurrence_idigbio.csv | input data | 7 days ago |
| occurrence_idigbio.zip | tutorial data | yesterday |
| subtree-ottol-1035588-Saxifragaceae.tre | added pre-commit checks; unfinished | 4 months ago |

# Data Preparation: occurrence records

2. CSV file containing records for one or more taxa.

   1. A CSV file is a text file with one species occurrence record per line. The file must be a delimited text file, and the first line must contain field names. Each record/line must contain a species (or other group) identifier, such as scientificName or species_name, and x and y coordinates indicating a geographic location. The field names for these 3 columns are specified in the script parameter file. One simple tutorial example occurrence datafile is heuchera.csv which contains different heuchera species, grouped by name, with x and y coordinates. Another tutorial example file is a CSV file containing many fields, downloaded from gbif, occurrence_gbif.csv.

| | species_name | x | y |
|---|---|---|---|
| 1 | | | |
| 2 | Bensoniella oregona | -123.751 | 42.802 |
| 3 | Bensoniella oregona | -123.7903 | 42.802 |
| 4 | Bensoniella oregona | -123.7903 | 42.802 |
| 5 | Bensoniella oregona | -123.7707 | 42.7873 |
| 6 | Bensoniella oregona | -123.751 | 42.9927 |
| 7 | Bensoniella oregona | -123.9646 | 42.7788 |
| 8 | Bensoniella oregona | -123.7117 | 42.9047 |

# Data Preparation: script file

## Input: Script parameter file

A JSON parameter file is required for this command. The parameter file in our first example is split_gbif.json. This one splits GBIF data, which already contains accepted names, so we can skip name resolution.

The parameter file in our second example is split_resolve.json. This one splits both iDigBio and GBIF data, and resolves to the canonical form of accepted names according to the GBIF taxonomy service.

```
1   {
2       "log_filename": "/volumes/output/split_gbif.log",
3       "log_console": true,
4       "report_filename": "/volumes/output/split_gbif.rpt",
5       "max_open_writers": 100,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8           "/volumes/data/wranglers/no_wrangle.json",
9           "acceptedScientificName",
10          "decimalLongitude",
11          "decimalLatitude"
12        ]
13      ],
14      "out_dir": "/volumes/output/split_gbif"
15  }
```

```
1   {
2       "max_open_writers": 100,
3       "report_filename": "/volumes/output/split_resolve.rpt",
4       "log_filename": "/volumes/output/split_resolve.log",
5       "log_console": true,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8           "/volumes/data/wranglers/occ_resolve.json",
9           "acceptedScientificName",
10          "decimalLongitude",
11          "decimalLatitude"
12        ]],
13      "dwca": [
14          ["/volumes/data/input/occurrence_idigbio.zip",
15            "/volumes/data/wranglers/occ_resolve.json"
16          ]],
17      "out_dir": "/volumes/output/split_resolve"
18  }
```

# Data Preparation: script file

**Input: Script parameter file**

A JSON parameter file is required for this command. The parameter file in our first example is split_gbif.json. This one splits GBIF data, which already contains accepted names, so we can skip name resolution.

The parameter file in our second example is split_resolve.json. This one splits both iDigBio and GBIF data, and resolves to the canonical form of accepted names according to the GBIF taxonomy service.

These are the required and optional parameters:

- Required:
  - **out_dir**: Directory where the output data should be written. If the directory does not exist, it will be created
- Optional:
  - **max_open_writers**: The maximum number of data writers to have open at once.

**Required parameter**

```json
1   {
2       "log_filename": "/volumes/output/split_gbif.log",
3       "log_console": true,
4       "report_filename": "/volumes/output/split_gbif.rpt",
5       "max_open_writers": 100,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8           "/volumes/data/wranglers/no_wrangle.json",
9           "acceptedScientificName",
10          "decimalLongitude",
11          "decimalLatitude"
12        ]
13      ],
14      "out_dir": "/volumes/output/split_gbif"
15  }
```

# Data Preparation: script file

**Input: Script parameter file**

A JSON parameter file is required for this command. The parameter file in our first example is split_gbif.json. This one splits GBIF data, which already contains accepted names, so we can skip name resolution.

The parameter file in our second example is split_resolve.json. This one splits both iDigBio and GBIF data, and resolves to the canonical form of accepted names according to the GBIF taxonomy service.

These are the required and optional parameters:

- Required:
    - **out_dir**: Directory where the output data should be written. If the directory does not exist, it will be created
- Optional:
    - **max_open_writers**: The maximum number of data writers to have open at once.

**Optional parameters used**

```json
1   {
2       "log_filename": "/volumes/output/split_gbif.log",
3       "log_console": true,
4       "report_filename": "/volumes/output/split_gbif.rpt",
5       "max_open_writers": 100,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8           "/volumes/data/wranglers/no_wrangle.json",
9           "acceptedScientificName",
10          "decimalLongitude",
11          "decimalLatitude"
12        ]
13      ],
14      "out_dir": "/volumes/output/split_gbif"
15  }
```

# Data Preparation: script file

- **out_field**: The field name or names of columns to be included in output CSV files. If this field is left out, all fields from the first successfully processed record will be included in outputs.
- **dwca**: This is an optional argument, but either this, or **csv**, must be provided. List of 0 or more lists, each containing 2 arguments
  - input DwCA file, and
  - occurrence data wrangler configuration file (described in the next section).
- **csv**: This is an optional argument, but either this, or **dwca**, must be provided. List of 0 or more lists, each containing 5 arguments
  - input CSV file
  - occurrence data wrangler configuration file (described in the next section).
  - fieldname for grouping data (often a taxonomic designation such as scientificName)
  - fieldname for the longitude/x coordinate
  - fieldname for the latitude/y coordinate
- **species_list_filename**: File location to write list of species seen (after wrangling).
- **log_filename**: Output filename to write logging data
- **log_console**: 'true' to write log to console
- **report_filename**: output filename with data modifications made by wranglers

**Optional parameters used**

```json
1   {
2       "log_filename": "/volumes/output/split_gbif.log",
3       "log_console": true,
4       "report_filename": "/volumes/output/split_gbif.rpt",
5       "max_open_writers": 100,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8          "/volumes/data/wranglers/no_wrangle.json",
9          "acceptedScientificName",
10         "decimalLongitude",
11         "decimalLatitude"
12        ]
13      ],
14      "out_dir": "/volumes/output/split_gbif"
15  }
```

# Data Preparation: script file



- **out_field**: The field name or names of columns to be included in output CSV files. If this field is left out, all fields from the first successfully processed record will be included in outputs.
- **dwca**: This is an optional argument, but either this, or **csv**, must be provided. List of 0 or more lists, each containing 2 arguments
  - input DwCA file, and
  - occurrence data wrangler configuration file (described in the next section).
- **csv**: This is an optional argument, but either this, or **dwca**, must be provided. List of 0 or more lists, each containing 5 arguments
  - input CSV file
  - occurrence data wrangler configuration file (described in the next section).
  - fieldname for grouping data (often a taxonomic designation such as scientificName)
  - fieldname for the longitude/x coordinate
  - fieldname for the latitude/y coordinate
- **species_list_filename**: File location to write list of species seen (after wrangling).
- **log_filename**: Output filename to write logging data
- **log_console**: 'true' to write log to console
- **report_filename**: output filename with data modifications made by wranglers

**Optional parameters used**

```
1   {
2       "log_filename": "/volumes/output/split_gbif.log",
3       "log_console": true,
4       "report_filename": "/volumes/output/split_gbif.rpt",
5       "max_open_writers": 100,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8          "/volumes/data/wranglers/no_wrangle.json",
9          "acceptedScientificName",
10         "decimalLongitude",
11         "decimalLatitude"
12        ]
13      ],
14      "out_dir": "/volumes/output/split_gbif"
15  }
```

**Optional parameters NOT used**

- **out_field:** The field name or names of columns to be included in output CSV files. If this field is left out, all fields from the first successfully processed record will be included in outputs.
- **dwca:** This is an optional argument, but either this, or **csv**, must be provided. List of 0 or more lists, each containing 2 arguments
  - input DwCA file, and
  - occurrence data wrangler configuration file (described in the next section).
- **csv:** This is an optional argument, but either this, or **dwca**, must be provided. List of 0 or more lists, each containing 5 arguments
  - input CSV file
  - occurrence data wrangler configuration file (described in the next section).
  - fieldname for grouping data (often a taxonomic designation such as scientificName)
  - fieldname for the longitude/x coordinate
  - fieldname for the latitude/y coordinate
- **species_list_filename:** File location to write list of species seen (after wrangling).
- **log_filename:** Output filename to write logging data
- **log_console:** 'true' to write log to console
- **report_filename:** output filename with data modifications made by wranglers

```
1   {
2       "log_filename": "/volumes/output/split_gbif.log",
3       "log_console": true,
4       "report_filename": "/volumes/output/split_gbif.rpt",
5       "max_open_writers": 100,
6       "csv": [
7         ["/volumes/data/input/occurrence_gbif.csv",
8          "/volumes/data/wranglers/no_wrangle.json",
9          "acceptedScientificName",
10         "decimalLongitude",
11         "decimalLatitude"
12        ]
13      ],
14      "out_dir": "/volumes/output/split_gbif"
15  }
```

# Data Preparation: wrangler file

## Input: Wrangler configuration file

A data wrangler configuration is a file containing a JSON list of zero or more wranglers - each performs a different operation, and each has its own parameters. More information on file format, available wrangler types, and the required and/or optional parameters for each are in the **Occurrence Wrangler Types** section of data_wrangle_occurrence. The file is specified in the Script parameter file described above.

An example wrangler configuration file occ_resolve.json resolves names with GBIF before grouping the data by name.

If more than one dataset is being processed, it is logical to apply the same wranglers to each.

```
9 lines (9 sloc)    243 Bytes
1  [
2      {
3          "wrangler_type": "AcceptedNameOccurrenceWrangler",
4          "name_resolver": "gbif",
5          "out_map_filename": "/volumes/output/occ_resolve.namemap",
6          "map_write_interval": 100,
7          "out_map_format": "json"
8      }
9  ]
```

# Let's run this tutorial!



## Update tutorial

Change directory to the top directory in your cloned tutorials repository on your local computer, then update the repository.

```
astewart:~/git/tutorials$ git pull
```

## Run tutorial

Initiate the split_occurrence_data process with the following:

For MacOSX or Linux systems: .. code-block:

```
./run_tutorial.sh split_occurrence_data data/config/split_resolve.json
```

For Windows systems:

```
./run_tutorial.bat split_occurrence_data data/config/split_resolve.json
```

**Goal: (usually) produce a file per species containing occurrence data. SPLITTING FILES IS ESSENTIAL FOR FACILITATING SDMS!**

```
1   {
2       "max_open_writers": 100,
3       "report_filename": "/volumes/output/split_resolve.rpt",
4       "log_filename": "/volumes/output/split_resolve.log",
5       "log_console": true,
6       "csv": [
7           ["/volumes/data/input/occurrence_gbif.csv",
8            "/volumes/data/wranglers/occ_resolve.json",
9            "acceptedScientificName",
10           "decimalLongitude",
11           "decimalLatitude"
12           ]],
13      "dwca": [
14          ["/volumes/data/input/occurrence_idigbio.zip",
15           "/volumes/data/wranglers/occ_resolve.json"
16           ]],
17      "out_dir": "/volumes/output/split_resolve"
18  }
```

Most outputs are configured in the script parameter file, and may include:

1. A "report_filename" named in the script parameter file, a summary of point manipulations by each wrangler will be written to this file. split_resolve.rpt
2. A "log_filename" named in the script parameter file, that will be created. split_resolve.log
3. **A "log_console" named in the script parameter file, logs will be written to the**

   command prompt during execution.

4. A directory, named in the out_dir parameter, of output CSV files, one per species (or other grouping field). The basename of each CSV file will be named by the value in the grouping field. The tutorial example outputs for this command have been moved to the `data/input` directory, since we will use them in a later exercise. split_resolve

The process also produces outputs according to the wrangler configuration file:

1. If the AcceptedNameOccurrenceWrangler is included, and there is a name-map file named in out_map_filename parameter, this file will be output. The name-map is a JSON file with pairs of names - the original name to the accepted name according to the specified authority. This name-map is suitable to use for input when resolving another dataset containing a subset of the same original names. A sample output name-map is occ_resolve.namemap.

**WE SAW THIS AS AN OUTPUT WHEN RESOLVING NAMES!**

# Let's look at the output!

| | | |
|---|---|---|
| 📄 | Bensoniella oregona.csv | add heuchera data for SDM |
| 📄 | Conimitella williamsii.csv | add heuchera data for SDM |
| 📄 | Elmera racemosa.csv | add heuchera data for SDM |
| 📄 | Heuchera abramsii.csv | add heuchera data for SDM |
| 📄 | Heuchera acutifolia.csv | add heuchera data for SDM |
| 📄 | Heuchera alba.csv | add heuchera data for SDM |
| 📄 | Heuchera alpestris.csv | add heuchera data for SDM |
| 📄 | Heuchera americana.csv | add heuchera data for SDM |
| 📄 | Heuchera amoena.csv | add heuchera data for SDM |
| 📄 | Heuchera bracteata.csv | add heuchera data for SDM |
| 📄 | Heuchera brevistaminea.csv | add heuchera data for SDM |
| 📄 | Heuchera caespitosa.csv | add heuchera data for SDM |
| 📄 | Heuchera caroliniana.csv | add heuchera data for SDM |
| 📄 | Heuchera cespitosa.csv | add heuchera data for SDM |
| 📄 | Heuchera chlorantha.csv | add heuchera data for SDM |
| 📄 | Heuchera cuneata.csv | add heuchera data for SDM |

**Directory containing the output!**

```
26 lines (26 sloc)    800 Bytes

 1   {
 2       "max_open_writers": 100,
 3       "report_filename": "/volumes/output/split_wrangle_occurrence_data.rpt",
 4       "log_filename": "/volumes/output/split_wrangle_occurrence_data.log",
 5       "log_output": false,
 6       "csv": [
 7         ["/volumes/data/input/heuchera.csv",
 8          "/volumes/data/wranglers/occ_wrangler_resolve.json",
            "species_name",
            "x",
            "y"
         ],
         ["/volumes/data/input/occ_heuchera_gbif.csv",
           "/volumes/data/wranglers/occ_wrangler_resolve.json",
           "acceptedScientificName",
16          "decimalLongitude",
17          "decimalLatitude"
18        ]
19       ],
20       "dwca": [
21          ["/volumes/data/input/occ_heuchera_gbif.zip",
22            "/volumes/data/wranglers/occ_wrangler_resolve.json"
23          ]
24        ],
25       "out_dir": "/volumes/output/heuchera_split_resolve"
26   }
```

# Let's look at the output!



**These columns are found in .dwca files**

**Opening one of the folders contained in the output directory!**

# Let's look at the output!

There are many columns in .dwca files, but the ones that will be relevant for now are the last three!

| teProvince | taxonID | taxonomicStatus | taxonRank | typeStatus | uuid | verbatimEventDate | verbatimLocality | version | waterBody | species_name | x | y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| fornia | 3032647 | accepted | species | | 0208288e-4c4a-4d81-aa8f-120ec6327f6b | | | | | Heuchera abramsii | -117.6461 | 34.28921 |
| fornia | 3032647 | accepted | species | | e6f75926-c0ee-42e0-9203-2d6af3f489f5 | | | | | Heuchera abramsii | -117.65052 | 34.2878 |
| fornia | 3032645 | accepted | genus | | 3a393342-9941-4cf0-9217-875a584f9b78 | | | | | Heuchera abramsii | -117.64617 | 34.28917 |
| fornia | 3032645 | accepted | genus | | 222f2377-aeaa-4657-8db0-328a4d66cdb0 | | | | | Heuchera abramsii | -117.655 | 34.28694 |
| fornia | 3032645 | accepted | genus | | b87e34f8-6cce-4221-8890-8efe6ddc170c | | | | | Heuchera abramsii | -117.64444 | 34.31361 |
| fornia | 3032645 | accepted | genus | | f71d5f6a-fa5b-44c3-aa7f-54eb4e98b891 | | | | | Heuchera abramsii | -117.64389 | 34.28861 |
| fornia | 3032647 | accepted | species | | 3dced37c-6821-46f8-83a5-dbf4e295b208 | | | | | Heuchera abramsii | -117.6461 | 34.28921 |
| fornia | 3032647 | accepted | species | | 66daffa9-e2eb-4688-991e-7529688a39e8 | | | | | Heuchera abramsii | -117.6461 | 34.28921 |
| fornia | 3032647 | accepted | species | | 4bd2eaa8-efb9-4dbe-b762-50ec2a1e7402 | 9 jul 1967 | | | | Heuchera abramsii | -117.9288 | 34.3493 |
| fornia | 3032647 | accepted | species | | ad455030-61b1-473c-aa5d-154f3ac74c8b | 2000-6-22 | | | | Heuchera abramsii | -117.5725 | 34.22833 |
| fornia | 3032645 | accepted | genus | | ee7c56eb-4107-4680-9b12-672709b2679f | | | | | Heuchera abramsii | -117.64611 | 34.28944 |
| fornia | 3032647 | accepted | species | | 3e4b21e4-7522-4c6c-8d1b-70baa438f67a | | | | | Heuchera abramsii | -117.6461 | 34.28921 |
| fornia | 3032647 | accepted | species | | 14f1dde4-9b83-4790-b298-6e39e2b305e1 | 7 sep 1998 | | | | Heuchera abramsii | -117.64444 | 34.31361 |
| fornia | 3032645 | accepted | genus | | 0358a4cb-924c-443f-b040-589854ba5b6b | | | | | Heuchera abramsii | -118.05186 | 34.2705 |
| fornia | 3032647 | accepted | species | | 04450e62-be59-473d-9a5b-222089afe3d2 | 1918-7-6 | | | | Heuchera abramsii | -117.64617 | 34.28917 |

| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| fornia | 3032647 | accepted | species | 5c4f43ab-bf5a-4db6-b9f5-fbb3e21c12b8 | | | | Heuchera abramsii | -117.6349 | 34.29494 |
| fornia | 3032647 | accepted | species | 7716404e-fd72-48df-bd50-b1fe23465bfa | 11 jul 1997 | | | Heuchera abramsii | -118.99 | 37.62 |
| fornia | 3032647 | accepted | species | 5873ea11-fa8f-4313-a75a-40da11a963dc | | | | Heuchera abramsii | -117.6444444 | 34.3136111 |
| fornia | 3032647 | accepted | species | 1d209b6c-9abd-4a7a-b031-a623da9ecb24 | 1966-6-23 | | | Heuchera abramsii | -118.05194 | 34.27056 |
| fornia | 3032647 | accepted | species | bb69bf29-8626-406a-aabb-019b169216eb | | | | Heuchera abramsii | -117.65452 | 34.28682 |
| fornia | 3032647 | accepted | species | d32828a4-fc80-4f0c-b58c-ec4060dd2a2d | 1992-7-28 | | | Heuchera abramsii | -117.63583 | 34.30333 |
| fornia | 3032647 | accepted | species | c6ad9258-3a5f-440a-9a0f-bf23177043b2 | 1998-8-27 | | | Heuchera abramsii | -117.65333 | 34.28694 |
| | | | | | | | | Heuchera abramsii | -117.6684 | 34.2881 |
| | | | | | | | | Heuchera abramsii | -117.6684 | 34.2881 |
| | | | | | | | | Heuchera abramsii | -117.655 | 34.28694 |
| | | | | | | | | Heuchera abramsii | -117.655 | 34.28694 |
| | | | | | | | | Heuchera abramsii | -117.65452 | 34.28682 |
| | | | | | | | | Heuchera abramsii | -117.65452 | 34.28682 |
| | | | | | | | | Heuchera abramsii | -117.65452 | 34.28682 |
| | | | | | | | | Heuchera abramsii | -117.65452 | 34.28682 |
| | | | | | | | | Heuchera abramsii | -117.6527 | 34.2869 |
| | | | | | | | | Heuchera abramsii | -117.6527 | 34.2869 |
| | | | | | | | | Heuchera abramsii | -117.6514 | 34.287 |
| | | | | | | | | Heuchera abramsii | -117.6514 | 34.287 |

**REMEMBER! This output merged .dwca and .csv files. The .csv files will show as the entries containing only the last three columns!**

- ✓ **There are multiple ways and repositories to access occurrence data**.

- ✓ **Phylogenetic trees can be obtained from the Open Tree of Life, as well as other sources**.

- ✓ **BiotaPhy tools enable automated data downloads of occurrences and trees.**

- ✓ **BiotaPhy tools split large data sets into species-specific data sets for SDM and other uses.**

# Any questions??

**Please use the chat to ask your questions!**