# R and Python clients for GBIF species occurrence data

Scott Chamberlain*,a, Carl Boettiger[b]

*a University of California, Berkeley, CA, USA*

*b University of California, Berkeley, CA, USA*

## Abstract

1. xxx

2. xxx

3. xxx

4. xxxx

## Introduction

Perhaps the most fundamental element in many fields of ecology is the individual organism. The number of individuals of each species in a given location forms the basis for many sub-fields of ecology and evolution. Some research questions necessitate collecting new data, while others can easily take advantage of existing data. In fact, some ecology fields are built largely on existing data, e.g., macro-ecology (Brown, 1995; Beck et al., 2012).

Data on individuals, including which species, and where they're found, can be used for a large number of research questions. Biodiversity records have been used for a suite of other use cases: validating habitat suitability models with real occurrence data (Ficetola et al., 2014); ancestral range reconstruction (Ferretti et al., 2015; María Mendoza et al., 2015); development of invasive species watch lists (Faulkner et al., 2014); evaluating risk of invasive species spread (Febbraro et al., 2013); and effects of climate change on future biodiversity (Brown et al., 2015).

In addition to wide utility, this data is important for conservation. Biodiversity loss is one of the greatest challenges of our time (Pimm et al., 2014), and some have called this the sixth great mass extinction (Ceballos et al., 2015). Given this challenge there is a great need for data on specimen records, whether collected from live sightings in the field or specimens in museums.

*Corresponding author

*Email addresses:* `scott(at)ropensci.org` (Scott Chamberlain), `carl(at)ropensci.org` (Carl Boettiger)

**Global Biodiversity Information Facility**

There are many online services that collect and maintain specimen records. However, Global Biodiversity Information Facility (hereafter, GBIF, http://www.gbif.org) is the largest collection of biodiversity records globally, currently with 643 million records, 1.6 million taxa, 15,450 datasets from 780 publishers (as of 2016-02-09). Many large biodiversity warehouses such as iNaturalist (http://www.inaturalist.org), VertNet (http://vertnet.org), and USGS's Biodiversity Information Serving Our Nation (BISON; http://bison.usgs.ornl.gov) all feed into GBIF.

The most important data organizational level in GBIF is the individual record. In R using `rgbif` you'll recognize this as a row in the data.frame that is returned from `occ_search()` or `occ_data()`. Going upstream, each record is part of a dataset, where each dataset is submitted by an organization, organizations are organized into nodes, datasets are published through institutions (which may be hosted at another organization), and a network is a group of datasets (managed by GBIF).

Each record has some taxonomic name associated with it, which itself is linked to a lot of other taxonomic data. GBIF maintains their own taxonomic data.

All this organization matters because you can navigate through GBIF data through all these organizational levels.

**The clients**

*The rgbif package*

Herein, we describe the `rgbif` software package (Chamberlain et al.) for working with GBIF data in the R programming environment (R Core Team, 2014). R is a widely used language in academia, as well as non-profit and private sectors. Importantly, R makes it easy to execute all steps of the research process, including data management, data manipulation and cleaning, statistics, and visualization. Thus, an R client for getting GBIF data is a powerful tool to facilitate reproducible research.

The `rgbif` package is nearly completely written in R (a small Javascript library is included for reading well known text (Herring, 2011)), uses an MIT license to maximize use everywhere. `rgbif` is developed publicly on GitHub at https://github.com/ropensci/rgbif, where development versions of the package can be installed, and bugs and feature requests reported. Stable versions of `rgbif` can be installed from CRAN, the distribution network for R packages. `rgbif` is part of the rOpenSci project (http://ropensci.org), a developer network making R software to facilitate reproducible research.

*The pygbif library*

`pygbif` (Chamberlain) is a Python library for working with GBIF data in the Python programming
environment (**???**). Python is a general purpose programming language used widely in all sectors,
and for all parts of software development including server and client side use cases. Python is used
exclusively in some scientific disciplines (e.g., astronomy), and has partial usage in other disciplines. A
Python client for GBIF data is an important tool given the even wider usage of Python than R, though
maybe slightly less than R for ecology/biology.

```
pip install pygbif
```

```
import pygbif
```

The `pygbif` library is less mature and complete than the R package. It also uses an MIT license to
maximize use everywhere. `pygbif` is developed publicly on GitHub at https://github.com/sckott/pygbif,
where development versions of the package can be installed, and bugs and feature requests reported.
Stable versions of `pygbif` can be installed from pypi, the distribution network for Python libraries.

*Library interfaces*

`rgbif` and `pygbif` are designed following the GBIF Application Programming Interface, or API. The
GBIF API has four major components: registry, taxonomic names, occurrences, and maps. We also
include functions to interface with the OAI-PMH GBIF service; only dataset information is available
vis this service, however. We ignore maps in both libraries as it is concerned with generating maps
primarily for web applications. Both libraries have a suite of functions dealing with each of registry,
taxonomic names, and occurrences - we'll go through each in turn describing design and example usage.

*GBIF headers*

With each request `rgbif` and `pygbif` make to GBIF's API, we send request headers that tell GBIF what
client the request is coming from, including what version of the library. This helps GBIF know what
proportion of requests are coming from which client, and therefore from R vs. Python; this information
is helpful for GBIF in thinking about how people are using GBIF data.

*Registry*

The GBIF registry API services are spread across five sets of functions via the main GBIF API:

• Datasets

• Installations

• Networks

• Nodes

• Organizations

Dataset information in general is available via the OAI-PMH service, functions in `rgbif` prefixed with
`gbif_oai_`, but not available in `pygbif` yet.

Datasets are owned by organizations. Organizations are endorsed by nodes to share datasets with GBIF.
Datasets are published through institutions, which may be hosted at another organization. A network
is a group of datasets (managed by GBIF). Datasets are the units that matter the most with respect
to registry information, while installations, networks, nodes, and organizations are simply higher level
organizational structure.

*Datasets*

Dataset functions include search, dataset metadata retrieval, and dataset metrics. Searching for datasets
is an important part of the discovery process. One can search for datasets on the GBIF web portal.
However, programmatic searching using this package is much more powerful. Identifying datasets
appropriate for a research question is helpful as you can get metadata for each dataset, and track down
dataset specific problems, if any.

The `dataset_search()` function is one way to search for datasets. Here, we search for the query term
"oregon", which finds any datasets that have terms matching that term.

```
res <- dataset_search(query = "oregon")
res$data$datasetTitle[1:10]
#>  [1] "UCDavis - Western USA - Monarch Butterflies - 1892-2005"
#>  [2] "A geographic distribution database of Mononychellus mites (Acari: Tetranychidae) on cas
#>  [3] "Plantas Acuáticas de la Orinoquía Colombiana"
```

```
#>   [4] "USBombus, contemporary survey data of North American bumble bees (Hymenoptera,  Apidae,
#>   [5] "Wool carder bees of the genus Anthidium in the Western Hemisphere"
#>   [6] "CM Birds Collection"
#>   [7] "SDNHM Birds Collection"
#>   [8] "University of British Columbia Herbarium (UBC) - Bryophytes Collection"
#>   [9] "University of British Columbia Herbarium (UBC) - Vascular Plant Collection"
#> [10] "Bryophyte Collection - University of Washington Herbarium (WTU)"
```

See also `datasets()` and `dataset_suggest()` in `rgbif` for searching for datasets.

In Python, we can similary search for datasets. Here, search for datasets of type `OCCURRENCE`:

```python
from pygbif import registry
registry.datasets(type="OCCURRENCE")
```

*Dataset metrics.* Dataset metrics are another useful way of figuring out what datasets you may want to use. One drawback is that these metrics data are only available for datasets of type *checklist*, but there are quite a lot of them (2737).

Here, we search for dataset metrics for a single dataset, with uuid `ec93a739-1681-4b04-b62f-3a687127a17f`, a checklist of the ants (Hymenoptera: Formicidae) of the World.

```r
res <- dataset_metrics(uuid='ec93a739-1681-4b04-b62f-3a687127a17f')
data.frame(rank = names(res$countByRank),
           count = unname(unlist(res$countByRank)))
```

| rank | count |
|------|-------|
| SPECIES | 13710 |
| SUBSPECIES | 3234 |
| GENUS | 726 |
| TRIBE | 53 |
| SUBFAMILY | 20 |
| FAMILY | 2 |
| KINGDOM | 1 |

| rank | count |
|------|-------|
| PHYLUM | 1 |
| CLASS | 1 |
| ORDER | 1 |

And in Python, get metrics for the same dataset as above:

```python
from pygbif import registry
registry.dataset_metrics(uuid='ec93a739-1681-4b04-b62f-3a687127a17f')
```

*Networks, nodes, and installations*

Networks, nodes and installations are at a higher level of organization above datasets, but can be useful if you want to explore data from given organizations. Here, in R we search for the first 10 GBIF networks, returning just the title field.

```r
networks(limit = 10)$data$title
#>  [1] "GBIF Backbone Sources"
#>  [2] "Canadensys"
#>  [3] "Southwest Collections of Arthropods Network (SCAN)"
#>  [4] "VertNet"
#>  [5] "Dryad"
#>  [6] "GBIF Network"
#>  [7] "The Knowledge Network for Biocomplexity (KNB) "
#>  [8] "Online Zoological Collections of Australian Museums (OZCAM)"
#>  [9] "Catalogue of Life"
#> [10] "Ocean Biogeographic Information System (OBIS)"
```

And in Python:

```python
from pygbif import registry
registry.networks(limit = 10)
```

6

*Taxonomic names*

The GBIF taxonomic names API services are spread across five functions in `rgbif`:

- Search GBIF name backbone - `name_backbone()`
- Search across all checklists - `name_lookup()`
- Quick name lookup - `name_suggest()`
- Name usage of a name according to a checklist - `name_usage()`
- GBIF name parser - `parsenames()`

`pygbif` only has `name_backbone()` and `name_suggest()` at this time.

The goal of these name functions is often to settle on a taxonomic name known to GBIF's database. This serves two purposes: 1) when referring to a taxonomic name, you can point to a URI on the Internet, and 2) you can search for metadata on a taxon, and occurrences of that taxon in GBIF.

Taxonomic names are particularly tricky. Many different organizations have their own unique codes for the same taxonomic names, and some taxonomic groups have preferred sources for the definitive names for that group. That's why it's best to determine what name GBIF uses, and its associated identifier, for the taxon of interest instead of simply searching for occurrences with a taxonomic name.

When searching for occurrences (see below) you can search by taxonomic name (and other filters, e.g., taxonomic rank), but you're probably better off figuring out the taxonomic key in the GBIF backbone taxonomy, and using that to search for occurrences. The `taxonkey` parameter in the GBIF occurrences API expects a GBIF backbone taxon key.

*GBIF Backbone*

The GBIF backbone taxonomy is used in GBIF to have a consistent way to refer to taxonomic names throughout their services. The backbone has 4410899 unique names and 2497114 species names. The backbone taxonomy is also a dataset with key `d7dddbf4-2cf0-4f39-9b2a-bb099caae36c` (http://www.gbif.org/dataset/d7dddbf4-2cf0-4f39-9b2a-bb099caae36c).

We can search the backbone taxonomy with the function `name_backbone()` in both R and Python clients. Here, we're searching for the name *Poa*, restricting to genera, and the family *Poaceae*, in R

```
res <- name_backbone(name='Poa', rank='genus', family='Poaceae')
res[c('usageKey', 'kingdom')]
#> $usageKey
#> [1] 2704173
#>
#> $kingdom
#> [1] "Plantae"
```

and in Python

```
from pygbif import species
res = species.name_backbone(name='Poa', rank='genus', family='Poaceae')
[ res[x] for x in ['usageKey', 'kingdom'] ]
```

*Name searching*

One of the quickest ways to search for names is using `name_suggest()`, which does a very quick search and returns minimal data. Here, we're searching for the query term *Pum*, and we get back many names:

```
name_suggest(q='Pum', limit = 6)
```

| key | canonicalName | rank |
|---|---|---|
| 4848380 | Pumiliornis | GENUS |
| 2235094 | Pumilibranchipus | GENUS |
| 1795637 | Pumora | GENUS |
| 4581745 | Pumicia | GENUS |
| 4598790 | Pumilicopta | GENUS |
| 1593095 | Pumilomyia | GENUS |

The same in Python

```python
from pygbif import species
species.name_suggest(q='Pum', limit = 6)
```

<sup>144</sup> With these results, you can then proceed to search for occurrences with the taxon key(s), or drill down
<sup>145</sup> further with other name searching functions to get the exact taxon of interest.

<sup>146</sup> *Occurrences*

<sup>147</sup> GBIF provides two ways to get occurrence data: through the `/occurrence/search` route (see
<sup>148</sup> `occ_search` in `rgbif`, or `occurrences.search` in `pygbif`), or via the `/occurrence/download` route
<sup>149</sup> (many functions, see below). `occ_search()`/`occurrences.search` is the main function for the search
<sup>150</sup> route, and is more appropriate for smaller data, while the download functions are more appropriate for
<sup>151</sup> larger data requests.

<sup>152</sup> GBIF imposes for any given search a limit of 200,000 records in the search service, after which point
<sup>153</sup> you can't download any more records for that search. However, you can download more records for
<sup>154</sup> different searches.

<sup>155</sup> We think the search service is still quite useful for many people even given the 200,000 limit. For those
<sup>156</sup> that need more data, we have created a similar interface in the download functions, that should be easy
<sup>157</sup> to use. Users should take note that using the download service has a few extra steps to get data into R,
<sup>158</sup> but is straight-forward.

<sup>159</sup> The download service, like the occurrence search service, is rate-limited. That is, you can only have one
<sup>160</sup> to three downloads running simultaneously for your user credentials. However, simply check when a
<sup>161</sup> download job is complete, then you should be able to start a new download request.

<sup>162</sup> *Download API*

<sup>163</sup> The download API syntax is similar to the occurrence search API in that the same parameters are
<sup>164</sup> used, but the way in which the query is defined is different. For example, in the download API you can
<sup>165</sup> do greater than searches (i.e., `latitude > 50`), whereas you cannot do that in the occurrence search
<sup>166</sup> API. Thus, unfortunately, we couldn't make the query interface exactly the same for both search and
<sup>167</sup> download functions.

<sup>168</sup> Using the download service can consist of as few as three steps: 1) Request data via a search; 2)
<sup>169</sup> Download data; 3) Import data into R.

Request data download given a query. Here, we search for the taxon key `3119195`, which is the key for *Helianthus annuus* (http://www.gbif.org/species/3119195).

```
occ_download('taxonKey = 3119195')
#> <<gbif download>>
#>    Username: xxxx
#>    E-mail: xxxx
#>    Download key: 0000840-150615163101818
```

You can check on when the download is ready using the functions `occ_download_list()` and `occ_download_meta()`. When it's ready use `occ_download_get()` to download the dataset to your computer.

```
(res <- occ_download_get("0000840-150615163101818", overwrite = TRUE))
#> <<gbif downloaded get>>
#>    Path: ./0000840-150615163101818.zip
#>    File size: 3.19 MB
```

What's printed out above is a very brief summary of what was downloaded, the path to the file, and its size (in human readable form).

Next, read the data in to R using the function `occ_download_import()`.

```
library("dplyr")
dat <- occ_download_import(res)
dat %>%
  select(gbifID, decimalLatitude, decimalLongitude)
#>       gbifID abstract accessRights accrualMethod accrualPeriodicity accrualPolicy alternative
#> 1  725767384       NA                        NA                 NA            NA            NA
#> 2  725767447       NA                        NA                 NA            NA            NA
#> 3  725767450       NA                        NA                 NA            NA            NA
#> 4  725767513       NA                        NA                 NA            NA            NA
#> 5  725767546       NA                        NA                 NA            NA            NA
#> 6  725767579       NA                        NA                 NA            NA            NA
```

```
#> 7   725767609      NA                          NA              NA        NA        NA
#> 8   725767645      NA                          NA              NA        NA        NA
#> 9   725767678      NA                          NA              NA        NA        NA
#> 10  725767681      NA                          NA              NA        NA        NA
#> ..        ...     ...           ...            ...              ...       ...        ...
#> Variables not shown: available (lgl), bibliographicCitation (chr), conformsTo (lgl), contribu
#>       coverage (lgl), created (chr), creator (lgl), date (lgl), dateAccepted (lgl), dateCopyri
#>       (lgl), dateSubmitted (lgl), description (lgl), educationLevel (lgl), extent (lgl), forma
#>       hasFormat (lgl), hasPart (lgl), hasVersion (lgl), identifier (chr), instructionalMethod
```

178 In Python

```python
from pygbif import occurrences as occ
occ.download('taxonKey = 3119195')
(res = occ.download_get("0000840-150615163101818", overwrite = True))
```

179 We don't have `pygbif` functionality at the moment for importing data, but it's coming soon.

180 *Downloaded data format.* The downloaded dataset from GBIF is a Darwin Core Archive (DwC-A), an

181 internationally recognized biodiversity informatics standard (http://rs.tdwg.org/dwc/). The DwC-A

182 downloaded is a compressed folder with a number of files, including metadata, citations for each of the

183 datasets included in the download, and the data itself, in separate files for each dataset as well as one

184 single `.txt` file. In `rgbif::occ_download_import()`, we simply fetch data from the `.txt` file. If you

185 want to dig into the metadata, citations, etc., it is easily accessible from the folder on your computer.

186 *Search API*

187 The search API follows the GBIF API and is broken down into the following functions:

188    • Get a single numeric count of occurrences - `rgbif: occ_count()` / `pygbif: occurrences.count`

189    • Search for occurrences - `rgbif: occ_search()` / `pygbif: occurrences.count`

190    • A simplified and optimized version of `rgbif: occ_search()` or `occ_data()` / `pygbif:`

191    `occurrences.count`

192    • Get occurrences by occurrence identifier - `rgbif: occ_get()` / `pygbif: occurrences.count`

193    • Get occurrence metadata - `rgbif: occ_metadata()` / `pygbif: occurrences.count`

*Search for occurrences.* **R**

195 The main search work-horse is `occ_search()`. This function allows very flexible search definitions. In
196 addition, this function does paging internally, making it such that the user does not have worry about
197 the 300 records per request limit - but of course we can't go over the 200,000 maximum limit.

198 The output of `occ_search()` presents a compact `data.frame` so that no matter how large the
199 `data.frame`, the output is easily assessed because only a few of the records (rows) are shown, only a few
200 columns are shown (with others shown in name only), and metadata is shown on top of the `data.frame`
201 to indicate data found and returned, media records found, unique taxonomic hierarchies returned, and
202 the query executed.

203 The output of these examples, except one, aren't shown, but all run correctly.

204 Search by species name, using `name_backbone()` first to get key

```
(key <- name_suggest(q = 'Helianthus annuus', rank = 'species')$key[1])
#> [1] 3119195
occ_search(taxonKey = key, limit = 2)
#> Records found [21538]
#> Records returned [2]
#> No. unique hierarchies [1]
#> No. media records [2]
#> Args [taxonKey=3119195, limit=2, offset=0, fields=all]
#> First 10 rows of data
#>
#>               name        key decimalLatitude decimalLongitude
#> 1 Helianthus annuus 1249279611        34.04810       -117.79884
#> 2 Helianthus annuus 1249286909        32.58747        -97.10081
#> Variables not shown: issues (chr), datasetKey (chr), publishingOrgKey
#>      (chr), publishingCountry (chr), protocol (chr), lastCrawled (chr),
#>      lastParsed (chr), extensions (chr), basisOfRecord (chr), taxonKey
#>      (int), kingdomKey (int), phylumKey (int), classKey (int), orderKey
#>      (int), familyKey (int), genusKey (int), speciesKey (int),
#>      scientificName (chr), kingdom (chr), phylum (chr), order (chr),
```

```
#>      family (chr), genus (chr), species (chr), genericName (chr),
#>      specificEpithet (chr), taxonRank (chr), dateIdentified (chr), year
#>      (int), month (int), day (int), eventDate (chr), modified (chr),
#>      lastInterpreted (chr), references (chr), identifiers (chr), facts
#>      (chr), relations (chr), geodeticDatum (chr), class (chr), countryCode
#>      (chr), country (chr), rightsHolder (chr), identifier (chr),
#>      verbatimEventDate (chr), datasetName (chr), gbifID (chr),
#>      verbatimLocality (chr), collectionCode (chr), occurrenceID (chr),
#>      taxonID (chr), license (chr), recordedBy (chr), catalogNumber (chr),
#>      http...unknown.org.occurrenceDetails (chr), institutionCode (chr),
#>      rights (chr), eventTime (chr), identificationID (chr),
#>      coordinateAccuracy (dbl), coordinateAccuracyInMeters (dbl),
#>      occurrenceRemarks (chr)
```

205  Instead of getting a taxon key first, you can search for a name directly

```
occ_search(scientificName = 'Ursus americanus')
```

206  Search for many species

```
splist <- c('Cyanocitta stelleri', 'Junco hyemalis', 'Aix sponsa')
keys <- sapply(splist, function(x) name_suggest(x)$key[1], USE.NAMES = FALSE)
occ_search(taxonKey = keys, limit = 5, return = 'data')
```

207  Spatial search, based on well known text format (Herring, 2011), or a bounding box set of four coordinates

```
# well known text
occ_search(geometry = 'POLYGON((30.1 10.1, 10 20, 20 40, 40 40, 30.1 10.1))', limit = 20)
# bounding box
occ_search(geometry = c(-125.0,38.4,-121.8,40.9), limit = 20)
```

208  Get only occurrences with lat/long data

13

```
occ_search(hasCoordinate = TRUE, limit = 5)
```

209 Get only those occurrences with spatial issues. Spatial issues are a set of issues that are returned in
210 the `issues` field. They each indicate something different about that record. For example, the issue
211 `COUNTRY_COORDINATE_MISMATCH` indicates that the interpreted occurrence coordinates fall outside of
212 the indicated country. You can see how that might be useful when it comes to cleaning your data prior
213 to analysis/visualization.

```
occ_search(hasGeospatialIssue = TRUE, limit = 5)
```

214 **Python**

215 The equivalent occurrence search workhorse in `pygibf` is `occurrences.search`.

216 Search by species name, using `name_backbone()` first to get key

```
from pygbif import species
from pygbif import occurrences as occ
key = species.name_suggest(q = 'Helianthus annuus', rank = 'species')['data'][0]['key']
occ.search(taxonKey = key, limit = 2)
```

217 Instead of getting a taxon key first, you can search for a name directly

```
occ.search(scientificName = 'Ursus americanus')
```

218 Search for many species

```
from pygbif import species
from pygbif import occurrences as occ
splist = ['Cyanocitta stelleri', 'Junco hyemalis', 'Aix sponsa']
keys = [ species.name_suggest(x)['data'][0]['key'] for x in splist ]
occ.search(taxonKey = keys, limit = 5)
```

219 Spatial search, based on well known text format (Herring, 2011), or a bounding box set of four coordinates

```
from pygbif import occurrences as occ
# well known text
occ.search(geometry = 'POLYGON((30.1 10.1, 10 20, 20 40, 40 40, 30.1 10.1))', limit = 20)
# bounding box
occ.search(geometry = '-125.0,38.4,-121.8,40.9', limit = 20)
```

220 Get only occurrences with lat/long data

```
from pygbif import occurrences as occ
occ.search(hasCoordinate = True, limit = 5)
```

221 Get only those occurrences with spatial issues. Spatial issues are a set of issues that are returned in
222 the `issues` field. They each indicate something different about that record. For example, the issue
223 `COUNTRY_COORDINATE_MISMATCH` indicates that the interpreted occurrence coordinates fall outside of
224 the indicated country. You can see how that might be useful when it comes to cleaning your data prior
225 to analysis/visualization.

```
from pygbif import occurrences as occ
occ.search(hasGeospatialIssue = True, limit = 5)
```

226 *Data cleaning.* GBIF provides optional data issues with each occurrence record. These issues fall into
227 many different pre-defined classes, covering issues with taxonomic names, geographic data, and more
228 (see `occ_issues_lookup()` to find out more information on GBIF issues; and the same data on GBIF's
229 development site).

230 `occ_issues()` provides a way to easily filter data downloaded via `occ_search()` based on GBIF issues.

```
out <- occ_search(issue = 'DEPTH_UNLIKELY', limit = 500)
NROW(out)
#> [1] 4
out %>% occ_issues(-cudc) %>% .$data %>% NROW
#> [1] 4
```
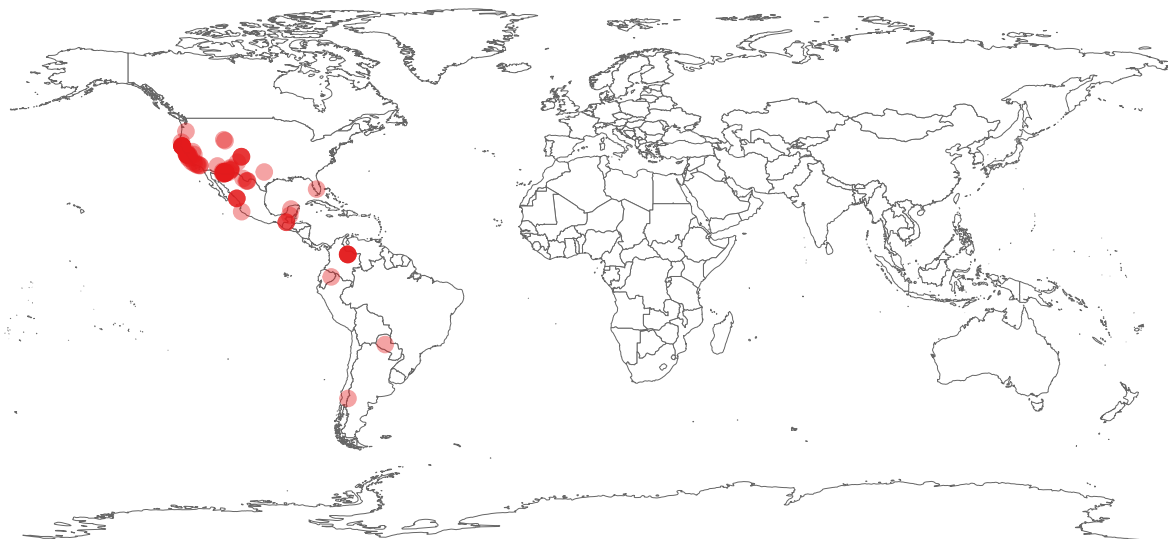
231 There's no equivalent interface in `pygbif`.

*Mapping*

An obvious downstream use case for species occurrence data is to map the data. `rgbif` per se is largely not concerned with making this easier, although we do have a simple wrapper around `ggplot2` to make it easy to get a quick plot of occurrence data. For example, here we plot 100 occurrences for *Puma concolor.*

```
key <- name_backbone(name='Puma concolor')$speciesKey
dat <- occ_search(taxonKey = key, limit = 100, hasCoordinate = TRUE)
gbifmap(dat$data)
```



Another package, mapr, is the perfect mapping companion to `rgbif`. It has convenient functions for handling input data from `rgbif`, `spocc`, or arbitrary data.frame's, and output plots for base plots, `ggplot2`, `ggmap` (`ggplot2` with map layers underneath), and interactive maps on GitHub gists or with Leaflet.js.

There's no equivalent interface in `pygbif`.

*GBIF data in other R packages*

We discuss usage of GBIF data in other R packages throughout the manuscript, but provide a synopsis here for clarity.

*taxize*

Some of the GBIF taxonomic services are also available in taxize, an R package that focuses on getting
data from taxonomic data sources on the web. For example, with `get_gbifid()` one can get GBIF IDs
used for a set of taxonomic names - then use those IDs in other functions in `taxize` to get additional
information, like taxonomically downstream children.

*spocc*

GBIF occurrence data is available in the R package spocc via `rgbif`. `spocc` is a unified interface
for fetching species occurrence data from many sources on the web. For example, a user can collect
occurrence data from GBIF, iDigBio, and iNaturalist, and easily combine them, then use other packages
to clean and visualize the data.

**Use cases**

The following are three use cases for `rgbif`: niche modeling, spatial change in biodiversity, and
distribution mapping.

*Ecological niche modeling*

In this example, we plot actual occurrence data for *Bradypus* species against a single predictor variable,
BIO1 (annual mean temperature). This is only ont step in a species distribution modelling nworkflow.
This example can be done using BISON data as well with our rbison package.

*Load libraries*

```
library("rgbif")
library("dismo")
library("maptools")
library("plyr")
```

*Raster files*

Make a list of files that are installed with the dismo package, then create a rasterStack from these

```
files <- list.files(paste(system.file(package = "dismo"), "/ex", sep = ""),
                     "grd", full.names = TRUE)
predictors <- stack(files)
```

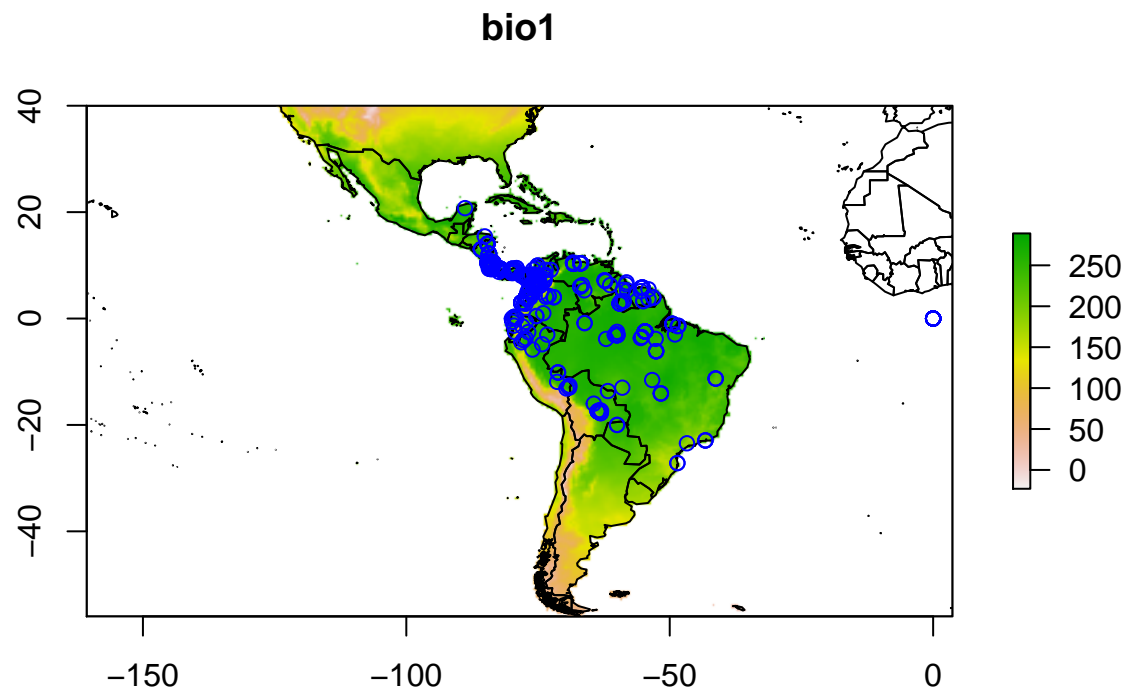266 *Get world boundaries*

```
data(wrld_simpl)
```

267 *Get GBIF data using the rOpenSci package rgbif*

```
nn <- name_lookup("bradypus*", rank = "species")
nn <- na.omit(unique(nn$data$nubKey))
df <- occ_search(taxonKey = nn, hasCoordinate = TRUE, limit = 500)
df <- df[ sapply(df, function(x) class(x$data)) %in% "data.frame" ]
df <- ldply(lapply(df, "[[", "data"))
df2 <- df[,c('decimalLongitude','decimalLatitude')]
```

268 *Plot*

269    (1) Add raster data, (2) Add political boundaries, (3) Add the points (occurrences)

```
plot(predictors, 1)
plot(wrld_simpl, add = TRUE)
points(df2, col = "blue")
```

**bio1**



270

*Biodiversity in big cities*

In this example, we collect specimen records across different cities using GBIF data from the `rgbif` package.

*Load libraries*

```
library("rgbif")
library("ggplot2")
library("plyr")
library("RCurl")
library("RColorBrewer")
```

*Get bounding boxes for some cites*

Bounding lat/long data is from here.

```
rawdat <- getURL('https://raw.githubusercontent.com/amyxzhang/boundingbox-cities/master/boundbox
dat <- read.table(text = rawdat, header = FALSE, sep="\t", col.names=c("city","minlat","maxlon",
dat <- data.frame(city=dat$city, minlon=dat$minlon, minlat=dat$minlat, maxlon=dat$maxlon, maxlat
```

```
getdata <- function(x){
  coords <- as.numeric(x[c('minlon','minlat','maxlon','maxlat')])
  num <- occ_search(geometry = coords)$meta$count
  data.frame(city=x['city'], richness=num, stringsAsFactors = FALSE)
}
```

```
out <- apply(dat, 1, getdata)
```

277  *Merge to original table*

```
out <- merge(dat, ldply(out), by="city")
```

278  *Add centroids from bounding boxes*

```
out <- transform(out, lat = (minlat+maxlat)/2, lon = (minlon+maxlon)/2)
```
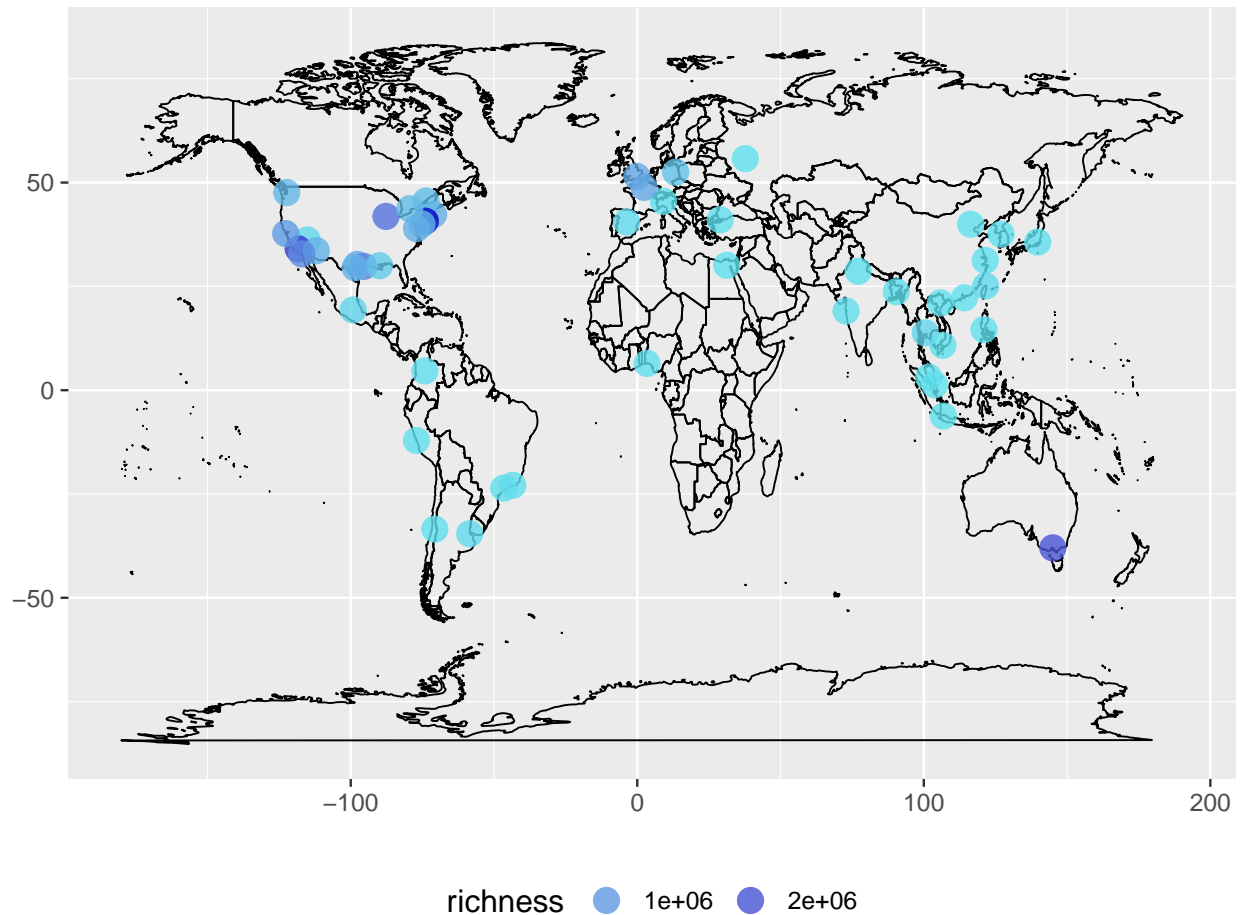
279  *Plot data*

```
mapp <- map_data('world')
ggplot(mapp, aes(long, lat)) +
  geom_polygon(aes(group=group), fill="white", alpha=0, color="black", size=0.4) +
  geom_point(data=out, aes(lon, lat, color=richness), size=5, alpha=0.8) +
  scale_color_continuous(low = "#60E1EE", high = "#0404C8") +
  labs(x="", y="") +
  theme_grey(base_size=14) +
  theme(legend.position = "bottom", legend.key = element_blank()) +
  guides(color = guide_legend(keywidth = 2))
```

richness ● 1e+06 ● 2e+06

280

## Valley oak occurrence data comparison

281 *Valley oak occurrence data comparison*

282 This example comes from Antonio J. Perez-Luque who shared his plot on Twitter. Antonio compared

283 the occurrences of Valley Oak (*Quercus lobata*) from GBIF to the distribution of the same species from

284 the Atlas of US Trees.

285 *Load libraries*

```
library('rgbif')
library('raster')
library('sp')
library('maptools')
library('rgeos')
library('scales')
```

286 *Get GBIF Data for Quercus lobata*

```r
keyQl <- name_backbone(name='Quercus lobata', kingdom='plants')$speciesKey
dat.Ql <- occ_search(taxonKey=keyQl, return='data', limit=50000)
```

287 *Get Distribution map of Q. lobata Atlas of US Trees (Little, E.)*

288 From http://esp.cr.usgs.gov/data/little/. And save shapefile in same directory

```r
url <- 'http://esp.cr.usgs.gov/data/little/querloba.zip'
tmp <- tempdir()
download.file(url, destfile = "~/querloba.zip")
unzip("~/querloba.zip", exdir = tmp)
ql <- readShapePoly(file.path(tmp, "querloba.shp"))
```

289 *Get Elevation data of US*

```r
alt.USA <- getData('alt', country = 'USA')
```
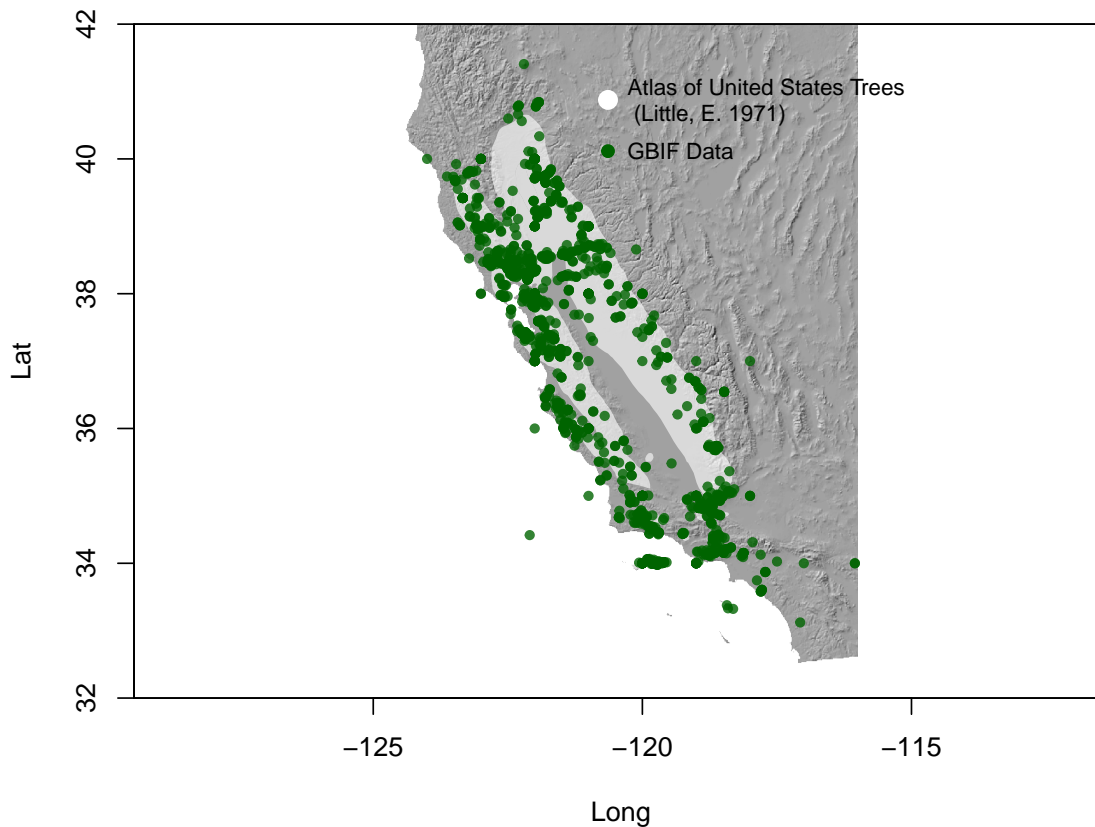
290 *Create Hillshade of US*

```r
alt.USA <- alt.USA[[1]]
slope.USA <- terrain(alt.USA, opt = 'slope')
aspect.USA <- terrain(alt.USA, opt = 'aspect')
hill.USA <- hillShade(slope.USA, aspect.USA, angle = 45, direction = 315)
```

291 *Plot map*

```r
plot(hill.USA, col = grey(0:100/100), legend = FALSE, xlim = c(-125, -116), ylim = c(32, 42), ma
# add shape from Atlas of US Trees
plot(ql, add = TRUE, col = alpha("white", 0.6), border = FALSE)
# add Gbif presence points
points(dat.Ql$decimalLongitude, dat.Ql$decimalLatitude, cex = .7, pch = 19, col = alpha("darkgre
legend(x = -121, y = 40.5, "GBIF Data", pch = 19, col = 'darkgreen', bty = 'n', pt.cex = 1, cex =
legend(x = -121, y = 41.5, "Atlas of United States Trees \n (Little, E. 1971)", pt.cex = 1.5, ce
```

22

**Distribution of Quercus lobata**



292

## Conclusions and future directions

293

294  The `rgbif` and `pygbif` libraries provide programmatic interfaces to GBIF's application programming

295  interface (API) - a powerful tool for making research using species occurrence data reproducible. In

296  fact, the `rgbif` package has already been used in 22 scholarly publications (as of 2015-11-14).

297  The `rgbif` package is relatively stable, and should not have many breaking changes unless necessitated

298  due to changes in the GBIF API.

299  The `pygbif` library is still in early development, and will greatly benefit from any feedback and use

300  cases.

301  One area of focus in the future is to attempt to solve many use cases that have been brought up with

302  respect to GBIF data. For example, some specimens are included in GBIF that are located in botanical

303  gardens. For many research questions, researchers are interested in "wild" type occurrences, not those

23

in human curated scenarios. Making removal of these occurrences easy would be very useful, but is actually quite a hard problem.

## Acknowledgments

## Data Accessibility

All scripts and data used in this paper can be found in the permanent data archive Zenodo under the digital object identifier (DOI). This DOI corresponds to a snapshot of the GitHub repository at github.com/sckott/msrgbif. Software can be found at https://github.com/ropensci/rgbif, under an MIT license.

## References

Beck J., Ballesteros-Mejia L., Buchmann CM., Dengler J., Fritz SA., Gruber B., Hof C., Jansen F., Knapp S., Kreft H., Schneider A-K., Winter M., Dormann CF. 2012. Whats on the horizon for macroecology? *Ecography* 35:673–683.

Brown JH. 1995. *Macroecology.* University of Chicago Press.

Brown KA., Parks KE., Bethell CA., Johnson SE., Mulligan M. 2015. Predicting plant diversity patterns in madagascar: Understanding the effects of climate and land cover change in a biodiversity hotspot. *PLOS ONE* 10:e0122721.

Ceballos G., Ehrlich PR., Barnosky AD., Garcia A., Pringle RM., Palmer TM. 2015. Accelerated modern human-induced species losses: Entering the sixth mass extinction. *Science Advances* 1:e1400253–e1400253.

Chamberlain S., Ram K., Barve V., Mcglinn D. *Pygbif: Interface to the global 'biodiversity' information facility 'aPI'.*

Chamberlain S. *Rgbif: Python client for gBIF.*

Faulkner KT., Robertson MP., Rouget M., Wilson JR. 2014. A simple, rapid methodology for developing invasive species watch lists. *Biological Conservation* 179:25–32.

Febbraro MD., Lurz PWW., Genovesi P., Maiorano L., Girardello M., Bertolino S. 2013. The use of climatic niches in screening procedures for introduced species to evaluate risk of spread: A case with the american eastern grey squirrel. *PLoS ONE* 8:e66559.

Ferretti F., Verd GM., Seret B., Šprem JS., Micheli F. 2015. Falling through the cracks: The fading history of a large iconic predator. *Fish and Fisheries*:n/a–n/a.

Ficetola GF., Rondinini C., Bonardi A., Baisero D., Padoa-Schioppa E. 2014. Habitat availability for amphibians and extinction threat: A global analysis. *Diversity and Distributions* 21:302–311.

Herring J. 2011. OpenGIS implementation standard for geographic information-simple feature access-part 1: Common architecture. *OGC Document* 4:122–127.

María Mendoza., Ospina OE., Cárdenas-Henao H., García-R JC. 2015. A likelihood inference of historical biogeography in the world's most diverse terrestrial vertebrate genus: Diversification of direct-developing frogs (craugastoridae: Pristimantis) across the neotropics. *Molecular Phylogenetics and Evolution* 85:50–58.

Pimm SL., Jenkins CN., Abell R., Brooks TM., Gittleman JL., Joppa LN., Raven PH., Roberts CM., Sexton JO. 2014. The biodiversity of species and their rates of extinction, distribution, and protection. *Science* 344:1246752–1246752.

R Core Team. 2014. *R: A language and environment for statistical computing.* Vienna, Austria: R Foundation for Statistical Computing.