

R, Python, and Ruby clients for GBIF species occurrence data

Scott Chamberlain^a, Carl Boettiger^b

^a*rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA*

^b*rOpenSci, Department of Environmental Science, Policy and Management, University of California, Berkeley, CA, USA*

Abstract

Corresponding Author:

Scott Chamberlain

rOpenSci, Museum of Paleontology, University of California, Berkeley, CA, USA

Email address: myrmecocystus@gmail.com

10 Background. The number of individuals of each species in a given location forms the basis for many
11 sub-fields of ecology and evolution. Data on individuals, including which species, and where they're found
12 can be used for a large number of research questions. Global Biodiversity Information Facility (hereafter,
13 GBIF) is the largest of these. Programmatic clients for GBIF would make research dealing with GBIF data
14 much easier and more reproducible.

15 Methods. We have developed clients to access GBIF data for each of the R, Python, and Ruby program-
16 ming languages: `rgbif`, `pygbif`, `gbifrb`.

17 Results. For all clients we describe their design and utility, and demonstrate some use cases.

Discussion. Programmatic access to GBIF will facilitate more open and reproducible science - the three
GBIF clients described herein are a significant contribution towards this goal.

Email addresses: `myrmecocystus(at)gmail.com` (Scott Chamberlain), `cboettig(at)berkeley.edu` (Carl Boettiger)

18 Introduction

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

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

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

34 Global Biodiversity Information Facility

35 There are many online services that collect and maintain specimen records. However, Global Biodiversity
36 Information Facility (hereafter, GBIF, <http://www.gbif.org>) is the largest collection of biodiversity records
37 globally, currently with 820 million records, roughly 5.9 million taxa, 36,000 datasets from 1,300 publishers
38 (as of 2016-02-09). Many large biodiversity warehouses such as iNaturalist (<http://www.inaturalist.org>),
39 VertNet (<http://vertnet.org>), and USGS's Biodiversity Information Serving Our Nation (BISON; [http://](http://bison.usgs.ornl.gov)
40 bison.usgs.ornl.gov) all feed into GBIF.

41 The most important organizational level in GBIF occurrence data is the occurrence record. The fields in
42 a record vary, but include information about taxonomy (kingdom, phylum, genus, species names) and their
43 identifiers, dataset metadata, and locality information including geospatial position. Going upstream, each
44 record is part of a dataset, where each dataset is submitted by an organization, organizations are organized
45 into nodes, datasets are published through institutions (which may be hosted at another organization), and
46 a network is a group of datasets (managed by GBIF).

Each occurrence record has some taxonomic name associated with it, which itself is linked to a lot of other taxonomic data - including a master taxonomic backbone that integrates taxonomies across many taxonomic authorities.

The organization of GBIF matters because you can navigate GBIF data through these hierarchical organizational levels - it helps to be familiar with the terminology and how each group relates to another.

The clients

Although we discuss libraries for R, Python, and Ruby here, we focus mostly on the R library `rgbif` as it has seen the most developer and user attention, and is the most mature.

rgbif

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 (<https://ropensci.org>), a developer network making R software to facilitate reproducible research.

pygbif

`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 disciplines.

```
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.

gbifrb

`gbifrb` (Chamberlain) is a library for working with GBIF data in the Ruby programming environment. Like Python, Ruby is a general purpose programming language used widely in all sectors. Unlike Python, Ruby is not used extensively in scientific disciplines. However, a Ruby client for GBIF data can be an important tool given how widely Ruby is used for web and web service development.

```
gem install gbifrb
```

```
require 'gbifrb'
```

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

Library interfaces

`rgbif`, `pygbif`, and `gbifrb` 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 (registry) information is available via this service, however. An interface to the GBIF maps API is in development for `rgbif`, but is non-existent for both `pygbif` and `gbifrb`. All three libraries have a suite of functions dealing with each of registry, taxonomic, names, and occurrences - we'll go through each in turn describing design of the user interface and example usage.

95 *GBIF headers*

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

100 *Registry*

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

- 102 • Datasets
- 103 • Installations
- 104 • Networks
- 105 • Nodes
- 106 • Organizations

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

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

114 *Datasets*

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

120 The `dataset_search()` function in `rgbif` is one way to search for datasets. Here, we search for the
121 term “oregon”, which finds any datasets that have words matching that term.

```

res <- dataset_search(query = "oregon")
res$data$datasetTitle[1:10]

#> [1] "Oregon State Ichthyology Collection"
#> [2] "Oregon State University Herpetological Collection"
#> [3] "Mygalomorph spiders from southwestern Oregon, USA, with descriptions of four new species"
#> [4] "A new species of Helobdella (Hirudinida: Glossiphoniidae) from Oregon, USA"
#> [5] "Cave millipedes of the United States. XVI. Two new species from Oregon Caves National Monument"
#> [6] "Annotated Checklist of the large branchiopod crustaceans of Idaho, Oregon and Washington, USA"
#> [7] "A new species of Chrysobothris Eschscholtz from Oregon and Washington, with notes on other Borealis"
#> [8] "Three new species of Grylloblatta Walker (Insecta: Grylloblattodea: Grylloblattidae), from southwestern Oregon"
#> [9] "A new species of Cladotanytarsus (Lenziella) from Oregon supports the systematic concept of ruficornis"
#> [10] "Two new species of Fluminicola (Caenogastropoda, Lithoglyphidae) from southwest Oregon, USA, and a new record from Oregon"

```

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

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

```

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

```

In Ruby, we can do the same. Here, search for datasets of type `OCCURRENCE`:

```

require 'gbifrb'
registry = Gbif::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 (29853).

Here, in R 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.

```
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
PHYLUM	1
CLASS	1
ORDER	1

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

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

132 The same in Ruby:

```
require 'gbifrb'
registry = Gbif::Registry
registry.dataset_metrics(uuid: 'ec93a739-1681-4b04-b62f-3a687127a17f')
```

133 *Networks, nodes, and installations*

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


```
networks(limit = 10)$data$title
#> [1] "TrIAS"
#> [2] "Arctos"
#> [3] "Freshwater Network"
#> [4] "Ocean Biogeographic Information System (OBIS)"
```

137 And in Python:

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

138 And in Ruby:

```
require 'gbifrb'
registry = Gbif::Registry
registry.networks(limit: 10)
```

139 *Taxonomic names*

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

- 141 • Search GBIF name backbone - **name_backbone()**
- 142 • Search across all checklists - **name_lookup()**
- 143 • Quick name lookup - **name_suggest()**
- 144 • Name usage of a name according to a checklist - **name_usage()**
- 145 • GBIF name parser - **parsenames()**

146 **pygbif** and **gbifrb** have all the same functions, except the name parser goes by **name_parser()** in
147 **pygbif** and **gbifrb**.

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

151 Taxonomic names are particularly tricky. Many different organizations have their own unique codes for
152 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 6575748 unique names and 2730631 species names. The backbone taxonomy is also a dataset with key `d7dddbf4-2cf0-4f39-9b2a-bb099caae36c` (<https://www.gbif.org/dataset/d7dddbf4-2cf0-4f39-9b2a-bb099caae36c>).

We can search the backbone taxonomy with the function `name_backbone()` in all three 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')]
#> # A tibble: 1 x 2
#>   usageKey kingdom
#>   <int> <chr>
#> 1 2704173 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'] ]
```

And in Ruby

```
require 'gbifrb'
species = Gbif::Species
res = species.name_backbone(name: 'Poa', rank: 'genus', family: 'Poaceae')
res.select { |k,v| k.match(/usageKey|kingdom/) }
```

168 *Name searching*

169 One of the quickest ways to search for names is using `name_suggest()`, which does a very quick search
170 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
2142856	Althepus pum	SPECIES
8589398	Pumiliopimoidae	FAMILY
8608394	Pumiliopimoa	GENUS
4579844	Pumilibaxa	GENUS
4581745	Pumicia	GENUS
8783253	Pumililema	GENUS

171 The same in Python

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

172 And in Ruby

```
require 'gbifrb'  
species = Gbif::Species  
species.name_suggest(q: 'Pum', limit: 6)
```

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

175 *Occurrences*

176 GBIF provides two ways to get occurrence data: through the `/occurrence/search` route (see `occ_search`
177 in `rgbif`, `occurrences.search` in `pygbif`, `Occurrences.search` in `gbifrb`), or via the `/occurrence/download`
178 route (many functions, see below).

`occ_search()/occurrences.search/Occurrences.search` are the main functions for the search route, and are more appropriate when you want less data, while the download functions are more appropriate for larger data requests.

Small vs. large amounts of data of course is all relative. GBIF imposes for any given search a limit of 200,000 records in the search service, after which point you can't download any more records for that search. However, you can download more records for different searches.

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

The download service, like the occurrence search service, is rate-limited. That is, you can only have one to three downloads running simultaneously for your user credentials. However, simply check when a download job is complete, then you can start a new download request. See “Queuing Download Requests” below for help automating many download requests in R.

Occurrences Download API

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

Using the download service can consist of as few as three steps: 1) Request data via a search; 2) 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
```

203 You can check on when the download is ready using the functions `occ_download_list()` and `occ_download_meta()`.
 204 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
```

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

207 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 audiolink
#> 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), contributor (chr),
#> coverage (lgl), created (chr), creator (lgl), date (lgl), dateAccepted (lgl), dateCopyrighted (lgl),
#> dateSubmitted (lgl), description (lgl), educationLevel (lgl), extent (lgl), format (lgl),
#> hasFormat (lgl), hasPart (lgl), hasVersion (lgl), identifier (chr), instructionalMethod (lgl),
```

208 In Python

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

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

210 The Ruby library `gbifrb` does not yet have occurrence download functionality.

211 *Downloaded data format*

212 The downloaded dataset from GBIF is a Darwin Core Archive (DwC-A), an internationally recog-
213 nized biodiversity informatics standard (<http://rs.tdwg.org/dwc/>). The DwC-A downloaded is a com-
214 pressed folder with a number of files, including metadata, citations for each of the datasets included in
215 the download, and the data itself, in separate files for each dataset as well as one single `.txt` file. In
216 `rgbif::occ_download_import()`, we simply fetch data from the `.txt` file. If you want to dig into the
217 metadata, citations, etc., it is easily accessible from the folder on your computer.

218 *Search API*

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

- 220 • Get a single numeric count of occurrences - `rgbif: occ_count()` / `pygbif: occurrences.count` /
221 `gbifrb: Occurrences.count`
- 222 • Search for occurrences - `rgbif: occ_search()` / `pygbif: occurrences.search` / `gbifrb: Occurrences.search`
- 223 • A simplified and optimized version of `rgbif: occ_search()` or `occ_data()` / `none` / `none`
- 224 • Get occurrences by occurrence identifier - `rgbif: occ_get()` / `pygbif: occurrences.get` / `gbifrb:`
225 `Occurrences.get`
- 226 • Get occurrence metadata - `rgbif: occ_metadata()` / `pygbif: various` / `gbifrb: various`

227 *Search for occurrences*

228 The main search work-horse is `occ_search()`. This function allows very flexible search definitions. In
229 addition, this function does paging internally, making it such that the user does not have worry about the
230 300 records per request limit - but of course we can't go over the 200,000 maximum limit.

231 The output of `occ_search()` presents a compact `data.frame` so that no matter how large the `data.frame`,
232 the output is easily assessed because only a few of the records (rows) are shown, only a few columns are

233 shown (with others shown in name only), and metadata is shown on top of the `data.frame` to indicate data
234 found and returned, media records found, unique taxonomic hierarchies returned, and the query executed.

235 The output of these examples, except one, aren't shown.

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

237 **R**

```
library(rgbif)

(key <- name_suggest(q = 'Helianthus annuus', rank = 'species'))$key[1])
#> [1] 9206251

occ_search(taxonKey = key, limit = 2)
#> Records found [46490]
#> Records returned [2]
#> No. unique hierarchies [1]
#> No. media records [2]
#> No. facets [0]
#> Args [limit=2, offset=0, taxonKey=9206251, fields=all]
#> # A tibble: 2 x 71
#>   key      scientificName decimalLatitude decimalLongitude issues datasetKey
#>   <chr> <chr>                <dbl>                <dbl> <chr>   <chr>
#> 1 2542~ Helianthus an~          25.8                -100. cdrou~ 50c9509d~
#> 2 2550~ Helianthus an~          25.6                -100. cdrou~ 50c9509d~
#> # ... with 65 more variables: publishingOrgKey <chr>, installationKey <chr>,
#> #   publishingCountry <chr>, protocol <chr>, lastCrawled <chr>,
#> #   lastParsed <chr>, crawlId <int>, extensions <chr>, basisOfRecord <chr>,
#> #   taxonKey <int>, ...
```

238 **Python**

```
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)
```

239 **Ruby**

```
require 'gbifrb'
species = Gbif::Species
occ = Gbif::Occurrences
key = species.name_suggest(q: 'Helianthus annuus', rank: 'species')['data'][0]['key']
occ.search(taxonKey: key, limit: 2)
```

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

241 **R**

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

242 **Python**

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

243 **Ruby**

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

244 Search for many species

245 **R**

```
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')
```

246 **Python**

```
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)
```

247 **Ruby**


```

species = Gbif::Species
occ = Gbif::Occurrences
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)

```

248 Spatial search, based on well known text format (Herring 2011), or a bounding box set of four coordinates. The well known text string and the bounding box in the below example specify the same rectangular area in California, centering approximately on Sacramento. Whereas the bounding box format requires longitude SW corner, latitude SW corner, longitude NE corner, latitude NE corner, the well known text string requires an extra long/lat pair to close the polygon.

253 R

```

# well known text
wkt <- 'POLYGON((-122.6 39.9,-120.0 39.9,-120.0 37.9,-122.6 37.9,-122.6 39.9))'
occ_search(geometry = wkt, limit = 20)
# bounding box
occ_search(geometry = c(-122.6,37.9,-120.0,39.9), limit = 20)

```

254 Python

```

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)

```

255 Ruby

```

occ = Gbif::Occurrences
# 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)

```

256 Get only occurrences with lat/long data using the `hasCoordinate` parameter

257 **R**

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

258 **Python**

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

259 **Ruby**

```
occ = Gbif::Occurrences
occ.search(hasCoordinate: true, limit: 5)
```

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

265 **R**

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

266 **Python**

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

267 **Ruby**

```
occ = Gbif::Occurrences
occ.search(hasGeospatialIssue: true, limit: 5)
```

268 **Data cleaning**

269 GBIF provides optional data issues with each occurrence record. These issues fall into many different pre-
270 defined classes, covering issues with taxonomic names, geographic data, and more (see `rgbif::occ_issues_lookup()`)

271 to find out more information on GBIF issues; and the same data on GBIF's development site).

272 `rgbif::occ_issues()` provides a way to easily filter data downloaded via `rgbif::occ_search()` based
273 on GBIF issues.

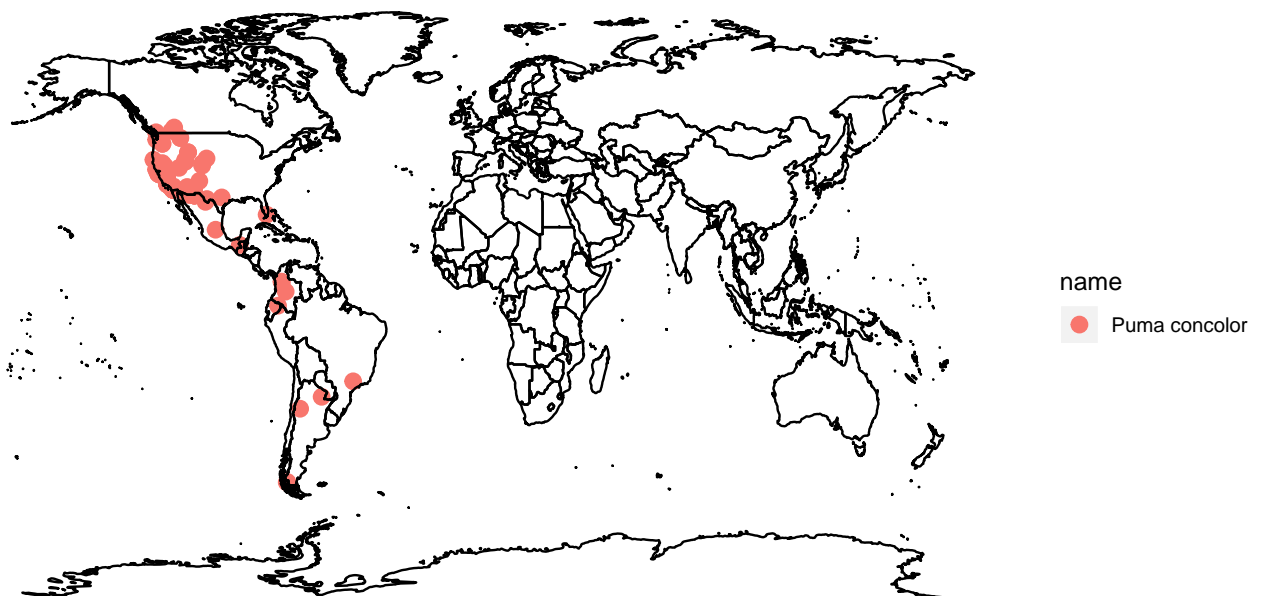
```
out <- occ_search(issue = 'DEPTH_UNLIKELY', limit = 500)
NROW(out)
#> [1] 5
out %>% occ_issues(-cudc) %>% .$data %>% NROW
#> [1] 500
```

274 There's no equivalent interface in `pygbif` or `gbifrb` yet.

275 Mapping

276 An obvious downstream use case for species occurrence data is to map the data. `rgbif` has enough to
277 do without handling mapping. Another package, `mapr`, works with `rgbif` to make this easy. For example,
278 here we plot 100 occurrences for *Puma concolor*.

```
library(mapr)
key <- name_backbone(name='Puma concolor')$speciesKey
dat <- occ_search(taxonKey = key, limit = 100, hasCoordinate = TRUE)
map_ggplot(dat, name = "species")
```



279

280 `mapr` has convenient functions for handling input data from `rgbif`, `spocc`, or arbitrary `data.frame`'s, and
281 output plots for base plots, `ggplot2`, `ggmap` (`ggplot2` with map layers underneath), and interactive maps
282 on GitHub gists or with Leaflet.js.

283 There's no equivalent interface in `pygbif` or `gbifrb`.

284 Citing GBIF data

285 All the data within GBIF was painstakingly collected by people and institutions across the globe. They
286 all, including GBIF, deserve to be cited appropriately if you use any GBIF data in your research. There
287 are clear forces making this difficult however. Publishers often limit, without reason if the publication is
288 online only, the number of citations one can include a manuscript. Even when this happens, you can always
289 include additional citations in an appendix.

290 To make citing data from GBIF easier `rgbif` has a function `gbif_citation()`, which accepts the output
291 of many different `rgbif` functions and gives citations for each dataset represented in the data therein. For
292 example, suppose you searched for occurrences for `taxonKey` 9206251 (*Helianthus annuus*):

```
res <- occ_search(taxonKey = 9206251, limit = 20)
```

293 We retrieved 20 occurrences from 4 different datasets.

294 To get citations for each of the four datasets, one can pass the output of `occ_search()`, here the object
295 `res`, to the `gbif_citation()` function:

```
cites <- gbif_citation(res)
cites
#> [[1]]
#> <<rgbif citation>>
#> Citation: iNaturalist.org (2020). iNaturalist Research-grade Observations.
#> Occurrence dataset https://doi.org/10.15468/ab3s5x accessed via
#> GBIF.org on 2020-02-25.. Accessed from R via rgbif
#> (https://github.com/ropensci/rgbif) on 2020-02-25
#> Rights:
#>
#> [[2]]
```

```

#> <<rgbif citation>>
#>   Citation: Shah M, Coulson S (2020). Artportalen (Swedish Species Observation
#>         System). Version 92.180. ArtDatabanken. Occurrence dataset
#>         https://doi.org/10.15468/kllkyl accessed via GBIF.org on 2020-02-25..
#>         Accessed from R via rgbif (https://github.com/ropensci/rgbif) on
#>         2020-02-25
#>   Rights:
#>
#> [[3]]
#> <<rgbif citation>>
#>   Citation: naturgucker.de. naturgucker. Occurrence dataset
#>         https://doi.org/10.15468/uc1apo accessed via GBIF.org on 2020-02-25..
#>         Accessed from R via rgbif (https://github.com/ropensci/rgbif) on
#>         2020-02-25
#>   Rights:
#>
#> [[4]]
#> <<rgbif citation>>
#>   Citation: Vanreusel W, Barendse R, Steeman R, Gielen K, Swinnen K, Desmet P,
#>         Herremans M (2020). Waarnemingen.be - Non-native plant occurrences in
#>         Flanders and the Brussels Capital Region, Belgium. Version 1.14.
#>         Natuurpunt. Occurrence dataset https://doi.org/10.15468/smdvdo accessed
#>         via GBIF.org on 2020-02-25.. Accessed from R via rgbif
#>         (https://github.com/ropensci/rgbif) on 2020-02-25
#>   Rights:

```

And we get four citations, one of each of the datasets behind the 20 occurrence records retrieved.

We don't yet have the same functionality available in `pygbif`, but it will be coming soon.

Downloads vs. searches

The distinction between “downloads” and “searches” is not especially meaningful with respect to the science done with GBIF data. However, downloads and searches are different with respect to easily being

able to cite GBIF data, which matters for GBIF and for data providers.

Downloads refers to the `rgbif` functions that start with `occ_download`, and the `pygbif` methods that start with `occurrences.download`. These methods use the GBIF “downloads API”, which you can interact with from `rgbif` and `pygbif` programatically, similar to how you would interact with the GBIF website - each interaction creating a download request, which processes in the background. The user has to ask if the request is ready to download. Downloads can request a very large amount of data.

Searches are faster and more dynamic. Search functions in `rgbif` include `occ_search`, `occ_data`, and `occ_get`, while `pygbif` methods include `occurrences.search`, `occurrences.get`. These methods can only request up to 200,000 occurrence records.

The big difference between downloads and searches with respect to citations is that each download generates a Digital Object Identifier (or DOI). This DOI can be used cite the entire set of data, that can include many datasets. There’s no good way to refer to the entire results of searches, other than the user manually creating a dataset somewhere public and linking to it from Zenodo to get a DOI. However, simply using downloads is simpler in terms of getting a single citeable DOI for the data used in a research paper/project.

In `rgbif`, the output of `occ_download_meta()` can be passed to `gbif_citation()` to get the citation for the entire download.

```
key <- "0000122-171020152545675"
res <- occ_download_meta(key)
gbif_citation(res)
#> $download
#> [1] "GBIF Occurrence Download https://doi.org/10.15468/dl.yghxj7 Accessed from R via rgbif (https://doi.org/10.15468/dl.yghxj7)"
#>
#> $datasets
#> NULL
```

If you however download the data first, then `gbif_citation()` can get citations for all the datasets contained within the download.

```
d1 <- occ_download_get(key, overwrite = TRUE)
gbif_citation(d1)
#> $download
```

```

#> [1] "GBIF Occurrence Download https://doi.org/10.15468/dl.yghxj7 Accessed from R via rgbif (https://github.com/ropensci/rgbif) on
#>
#> $datasets
#> $datasets[[1]]
#> <<rgbif citation>>
#> Citation: Grant S, Jones J (2017). Field Museum of Natural History (Zoology)
#> Invertebrate Collection. Version 18.6. Field Museum. Occurrence Dataset
#> https://doi.org/10.15468/6q5vuc accessed via GBIF.org on 2017-10-20..
#> Accessed from R via rgbif (https://github.com/ropensci/rgbif) on
#> 2020-02-24
#> Rights: To the extent possible under law, the publisher has waived all
#> rights to these data and has dedicated them to the Public Domain (CC0
#> 1.0). Users may copy, modify, distribute and use the work, including
#> for commercial purposes, without restriction.
#>
#> $datasets[[2]]
#> <<rgbif citation>>
#> Citation: Creuwels J (2017). Naturalis Biodiversity Center (NL) - Crustacea.
#> Naturalis Biodiversity Center. Occurrence Dataset
#> https://doi.org/10.15468/vjoltu accessed via GBIF.org on 2017-10-20..
#> Accessed from R via rgbif (https://github.com/ropensci/rgbif) on
#> 2020-02-24
#> Rights: To the extent possible under law, the publisher has waived all
#> rights to these data and has dedicated them to the Public Domain (CC0
#> 1.0). Users may copy, modify, distribute and use the work, including
#> for commercial purposes, without restriction.

```

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.

323 *taxize*

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

328 *spocc*

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

333 **R vs. Python vs. Ruby**

334 Both R and Python are commonly used in science, and can be used for similar tasks. Python, however,
335 is a more general programming language, and can be used in more contexts than R can be used in. Ruby
336 is used very little in science; but, like Python, Ruby is very widely used as a general purpose programming
337 language, with heavy use in web development and web services.

338 The three clients can do a lot of the same tasks. We envision *rgbif* being more common in workflows
339 of academics asking research questions, whereas *pygbif* and *gbifrb* can do that as well, but may be more
340 easily used in a website.

341 The R client *rgbif* has had much more development time than *pygbif* and *gbifrb*, but with time
342 *pygbif* and *gbifrb* will become equally mature.

343 **Use cases**

344 The following are three use cases for the R library *rgbif*: niche modeling, spatial change in biodiversity,
345 and distribution mapping.

346 *Ecological niche modeling*

347 In this example, we plot actual occurrence data for *Bradypus* species against a single predictor variable,
348 BIO1 (annual mean temperature). This is only one step in a species distribution modelling workflow.

349 This example can be done using BISON data as well with our *rbison* package.

350 *Load libraries*


```
library("sp")
library("rgbif")
library("dismo")
library("maptools")
library("dplyr")
```

351 *Raster files*

352 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)
```

353 *Get world boundaries*

```
data(wrld_simpl)
```

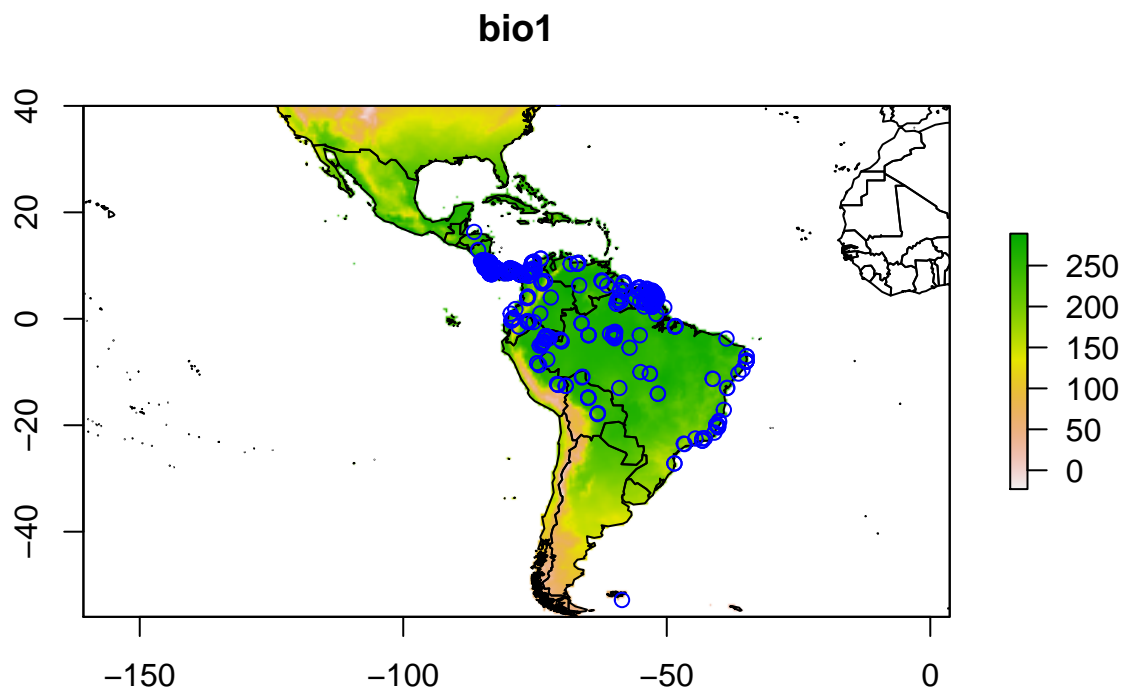
354 *Get GBIF data using the rOpenSci package rgbif*

```
nn <- name_lookup("bradypus*", rank = "species")
nn <- na.omit(unique(nn$data$nubKey))
df <- occ_search(taxonKey = nn, hasCoordinate = TRUE, limit = 500)
df_data <- df[ sapply(df, function(x) any(class(x$data) %in% "tbl_df")) ]
df_data <- dplyr::bind_rows(lapply(df_data, "[[", "data"))
df2 <- df_data %>% dplyr::select(decimalLongitude, decimalLatitude)
```

355 *Plot*

356 (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")
```



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("httr")
library("RColorBrewer")
library("wicket")
```

Get bounding boxes for some cities

Bounding lat/long data is from <https://raw.githubusercontent.com/amyxzhang/boundingbox-cities/master/boundbox.txt>.

```
url <- 'https://raw.githubusercontent.com/amyxzhang/
boundingbox-cities/master/boundbox.txt'
rawdat <- content(GET(sub("\n", "", url)), as = "text")
```

```

dat <- read.table(
  text = rawdat, header = FALSE,
  sep = "\t", col.names = c("city", "minlat", "maxlon", "maxlat", "minlon"),
  stringsAsFactors = FALSE)
dat <- data.frame(
  city = dat$city, minlon = dat$minlon,
  minlat = dat$minlat, maxlon = dat$maxlon,
  maxlat = dat$maxlat,
  stringsAsFactors = FALSE
)

```

365 A helper function to get count data. GBIF has a count API, but we can't use that with a geometry
 366 search as that API doesn't support geospatial search. We can however use the search API via `occ_search()`
 367 and set `limit = 1` so that we

```

getdata <- function(x){
  coords <- as.numeric(x[c('minlon', 'minlat', 'maxlon', 'maxlat')])
  wkt <- wicket::wkt_correct(wicket::bounding_wkt(values = coords))
  num <- occ_search(geometry = wkt, limit = 1)$meta$count
  data.frame(
    city = x['city'],
    richness = num,
    stringsAsFactors = FALSE
  )
}

```

```

out <- apply(dat, 1, getdata)

```

368 *Merge to original table*

```

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

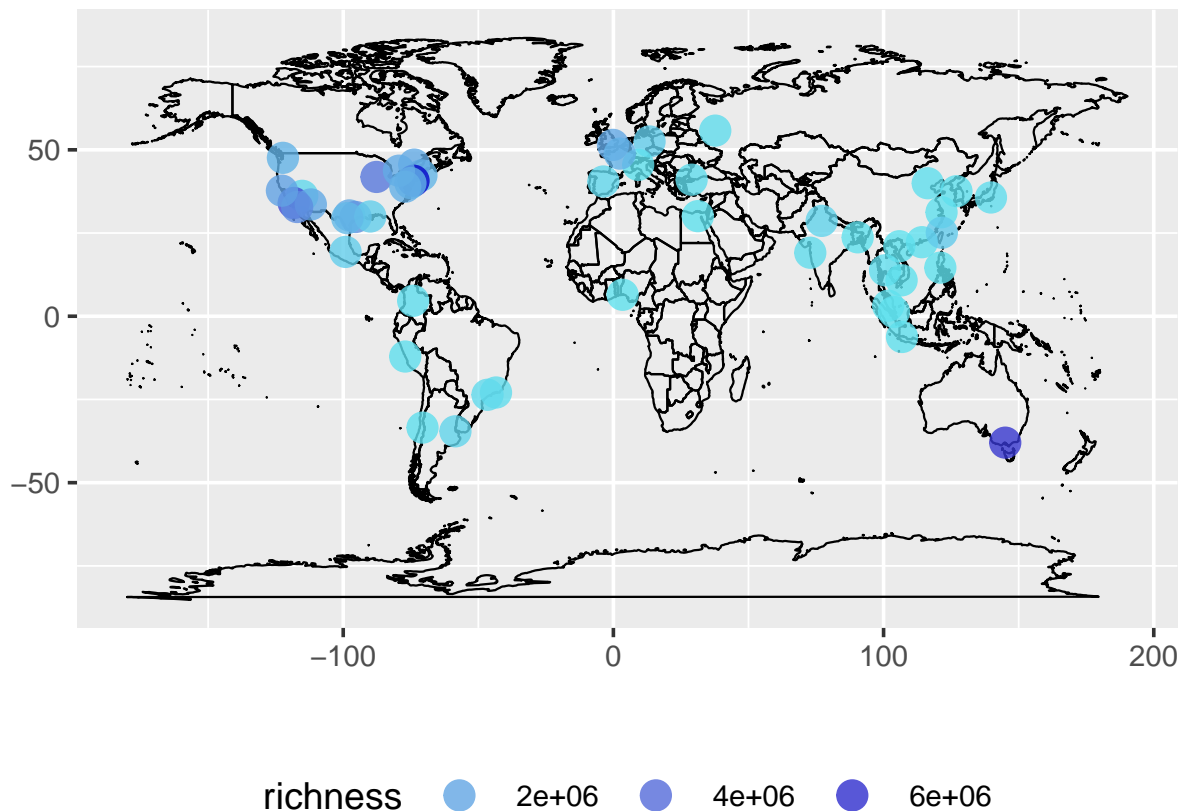
```

369 *Add centroids from bounding boxes*

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

370 *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))
```



371

372 *Valley oak occurrence data comparison*

373 This example is inspired by a tweet from Antonio J. Perez-Luque who shared his plot on Twitter. Antonio
 374 compared the occurrences of Valley Oak (*Quercus lobata*) from GBIF to the distribution of the same species

375 from the Atlas of US Trees.

376 The data in question from the example above is no longer available, so below we use a different species.

377 *Load libraries*

```
library('rgbif')
library('raster')
library('sp')
library('sf')
library('rgeos')
library('scales')
library('rnaturalearth')
```

378 *Get GBIF Data for Fraxinus excelsior*

```
keyFe <- name_backbone(name = 'Fraxinus excelsior', kingdom = 'plants')$speciesKey
dat.Fe <- occ_search(taxonKey = keyFe, return = 'data', limit = 10000L)
```

379 *Get Distribution map of F. excelsior European Forest Genetic Resources Programme*

380 From <http://www.euforgen.org/species/fraxinus-excelsior/>. And save shapefile in same directory

```
url <- 'http://www.euforgen.org/fileadmin/templates/euforgen.org/upload/Documents/Maps/Shapefile/Fraxi
tmp <- tempdir()
download.file(url, destfile = "fraxinus_excelsior.zip")
unzip("fraxinus_excelsior.zip", exdir = tmp)
fe <- sf::read_sf(file.path(tmp, "Fraxinus_excelsior_EUFORGEN.shp"))
```

381 *Get Elevation data of US*

```
eur <- rnaturalearth::ne_countries(continent = "europe", type = "map_units")
eur1 <- eur[eur$sovereignty != "Russia", ]
```

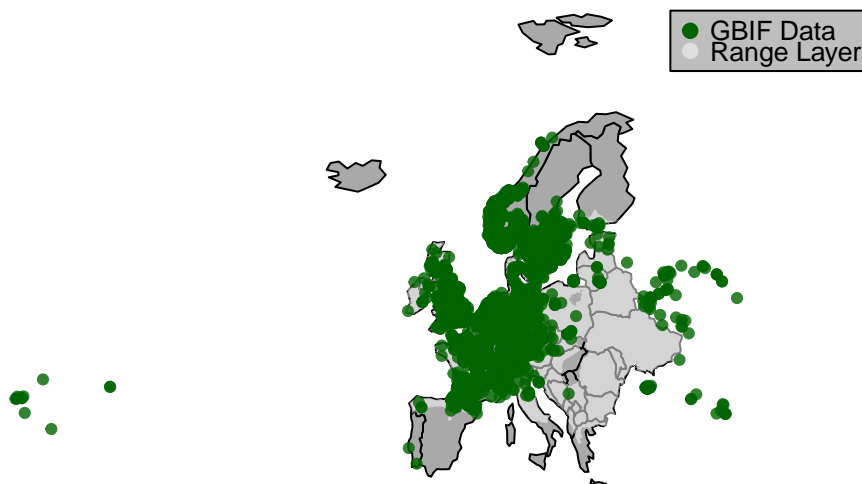
382 *Plot map*

```

plot(eur1, col = "darkgrey", legend = FALSE,
     main = 'Distribution of Fraxinus excelsior')
# add distribution range layer
plot(fe, add = TRUE, col = alpha("white", 0.5), border = FALSE)
# add Gbif presence points
points(dat.Fe$decimalLongitude, dat.Fe$decimalLatitude,
       cex = .7, pch = 19, col = alpha("darkgreen", 0.8))
legend(x = 38, y = 81, c("GBIF Data", "Range Layer"), pch = 19, bg = "grey",
      col = c('darkgreen', alpha("white", 0.5)), pt.cex = 1, cex = .8)

```

Distribution of Fraxinus excelsior



Conclusions and future directions

The `rgbif`, `pygbif`, and `gbifrb` libraries provide programmatic interfaces to GBIF's application programming interface (API) - a powerful tool for working with species occurrence data, and facilitating reproducible research. In fact, the `rgbif` package has already been used in more than 20 scholarly publications (as of 2020-02-25), including (Drozd & Šipoš 2013; Bartomeus *et al.* 2013; Barve 2014; Richardson *et al.* 2015; Feitosa *et al.* 2015; Collins *et al.* 2015; Malhado *et al.* 2015; Kong *et al.* 2015; Werner *et al.* 2015; Bone *et al.* 2015; Turner *et al.* 2015; Davison *et al.* 2015; Verheijen *et al.* 2015; Dellinger *et al.* 2015; Zizka & Antonelli 2015; Janssens *et al.* 2016; Robertson *et al.* 2016; Amano *et al.* 2016; Butterfield *et al.* 2016).

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

changes in the GBIF API. The `pygbif` and `gbifrb` libraries are in early development, and will greatly benefit from any feedback and use cases.

One area of focus in the future is to attempt to solve many use cases that have been brought up with respect to GBIF data. For example, some specimens are included in GBIF that are located in botanical gardens. For many research questions, researchers are interested in “wild” type occurrences, not those in human curated scenarios. Making removal of these occurrences easy would be very useful, but is actually quite a hard problem. There are many other problems like this, for which these three libraries will help in making more efficient and reproducible.

Acknowledgments

This project was supported in part by the Alfred P Sloan Foundation (Grant No. G-2014-13485), and in part by the Helmsley Foundation (Grant No. 2016PG-BRI004).

Data Accessibility

All scripts and data used in this paper can be found in the permanent data archive Zenodo under the digital object identifier (<https://doi.org/10.5281/zenodo.997554>). This DOI corresponds to a snapshot of the GitHub repository at <https://github.com/sckott/gbifms> that matches this preprint. Software can be found at <https://github.com/ropensci/rgbif>, <https://github.com/sckott/pygbif>, and <https://github.com/sckott/gbifrb>, all under MIT licenses. We thank all the users that have used `rgbif`, `pygbif`, and `gbifrb` and have given feedback and reported bugs. In addition, we greatly appreciate all the contributors to the three libraries, found at <https://github.com/ropensci/rgbif/graphs/contributors>, <https://github.com/sckott/pygbif/graphs/contributors>, and <https://github.com/sckott/gbifrb/graphs/contributors>.

References

- Amano, T., Lamming, J.D.L. & Sutherland, W.J. (2016). Spatial gaps in global biodiversity information and the role of citizen science. *BioScience*, **66**, 393–400. Retrieved from <http://dx.doi.org/10.1093/biosci/biw022>
- Bartomeus, I., Park, M.G., Gibbs, J., Danforth, B.N., Lakso, A.N. & Winfree, R. (2013). Biodiversity ensures plant-pollinator phenological synchrony against climate change (M. Eubanks, Ed.). *Ecology Letters*, **16**, 1331–1338. Retrieved from <http://dx.doi.org/10.1111/ele.12170>

- Barve, V. (2014). Discovering and developing primary biodiversity data from social networking sites: A novel approach. *Ecological Informatics*, **24**, 194–199. Retrieved from <http://dx.doi.org/10.1016/j.ecoinf.2014.08.008>
- Beck, J., Ballesteros-Mejia, L., Buchmann, C.M., Dengler, J., Fritz, S.A., Gruber, B., Hof, C., Jansen, F., Knapp, S., Kreft, H., Schneider, A.-K., Winter, M. & Dormann, C.F. (2012). Whats on the horizon for macroecology? *Ecography*, **35**, 673–683. Retrieved from <http://dx.doi.org/10.1111/j.1600-0587.2012.07364.x>
- Bone, R.E., Smith, J.A.C., Arrigo, N. & Buerki, S. (2015). A macro-ecological perspective on crassulacean acid metabolism (CAM) photosynthesis evolution in afro-madagascan drylands: Eulophiinae orchids as a case study. *New Phytologist*, **208**, 469–481. Retrieved from <http://dx.doi.org/10.1111/nph.13572>
- Brown, J.H. (1995). *Macroecology*. University of Chicago Press.
- Brown, K.A., Parks, K.E., Bethell, C.A., Johnson, S.E. & Mulligan, M. (2015). Predicting plant diversity patterns in madagascar: Understanding the effects of climate and land cover change in a biodiversity hotspot (L. Kumar, Ed.). *PLOS ONE*, **10**, e0122721. Retrieved from <http://dx.doi.org/10.1371/journal.pone.0122721>
- Butterfield, B.J., Copeland, S.M., Munson, S.M., Roybal, C.M. & Wood, T.E. (2016). Prestoration: Using species in restoration that will persist now and into the future. *Restor Ecol*. Retrieved from <http://dx.doi.org/10.1111/rec.12381>
- Ceballos, G., Ehrlich, P.R., Barnosky, A.D., Garcia, A., Pringle, R.M. & Palmer, T.M. (2015). Accelerated modern human-induced species losses: Entering the sixth mass extinction. *Science Advances*, **1**, e1400253–e1400253. Retrieved from <http://dx.doi.org/10.1126/sciadv.1400253>
- Chamberlain, S. *gbifrb: A ruby interface to the global biodiversity information facility API*. Retrieved from <https://github.com/sckott/gbifrb>
- Chamberlain, S. *pygbif: A python interface to the global biodiversity information facility API*. Retrieved from <https://github.com/sckott/pygbif>
- Chamberlain, S., Ram, K., Barve, V. & Mcglinn, D. *rgbif: An r interface to the global 'biodiversity' information facility API*. Retrieved from <https://github.com/ropensci/rgbif>
- Collins, R., Ribeiro, E.D., Machado, V.N., Hrbek, T. & Farias, I. (2015). A preliminary inventory of the catfishes of the lower rio nhamundá, brazil (ostariophysi, siluriformes). *BDJ*, **3**, e4162. Retrieved from <http://dx.doi.org/10.3897/bdj.3.e4162>

- 450 Davison, J., Moora, M., Opik, M., Adholeya, A., Ainsaar, L., Ba, A., Burla, S., Diedhiou, A.G., Hiiesalu, I.,
451 Jairus, T., Johnson, N.C., Kane, A., Koorem, K., Kochar, M., Ndiaye, C., Partel, M., Reier, U., Saks, U.,
452 Singh, R., Vasar, M. & Zobel, M. (2015). Global assessment of arbuscular mycorrhizal fungus diversity
453 reveals very low endemism. *Science*, **349**, 970–973. Retrieved from [http://dx.doi.org/10.1126/science.](http://dx.doi.org/10.1126/science.aab1161)
454 [aab1161](http://dx.doi.org/10.1126/science.aab1161)
- 455 Dellinger, A.S., Essl, F., Hojsgaard, D., Kirchheimer, B., Klatt, S., Dawson, W., Pergl, J., Pyšek, P.,
456 Kleunen, M. van, Weber, E., Winter, M., Hörandl, E. & Dullinger, S. (2015). Niche dynamics of alien
457 species do not differ among sexual and apomictic flowering plants. *New Phytologist*, **209**, 1313–1323.
458 Retrieved from <http://dx.doi.org/10.1111/nph.13694>
- 459 Drozd, P. & Šipoš, J. (2013). R for all (i): Introduction to the new age of biological analyses. *Casopis*
460 *slezskeho zemskeho muzea (A)*, **62**. Retrieved from <http://dx.doi.org/10.2478/cszma-2013-0004>
- 461 Faulkner, K.T., Robertson, M.P., Rouget, M. & Wilson, J.R.U. (2014). A simple, rapid methodology
462 for developing invasive species watch lists. *Biological Conservation*, **179**, 25–32. Retrieved from <http://dx.doi.org/10.1016/j.biocon.2014.08.014>
463
- 464 Febbraro, M.D., Lurz, P.W.W., Genovesi, P., Maiorano, L., Girardello, M. & Bertolino, S. (2013). The
465 use of climatic niches in screening procedures for introduced species to evaluate risk of spread: A case
466 with the american eastern grey squirrel (H. Verbruggen, Ed.). *PLoS ONE*, **8**, e66559. Retrieved from
467 <http://dx.doi.org/10.1371/journal.pone.0066559>
- 468 Feitosa, Y.O., Absy, M.L., Latrubesse, E.M. & Stevaux, J.C. (2015). Late quaternary vegetation dynamics
469 from central parts of the madeira river in brazil. *Acta Bot. Bras.*, **29**, 120–128. Retrieved from
470 <http://dx.doi.org/10.1590/0102-33062014abb3711>
- 471 Ferretti, F., Verd, G.M., Seret, B., Šprem, J.S. & Micheli, F. (2015). Falling through the cracks: The fading
472 history of a large iconic predator. *Fish and Fisheries*, n/a–n/a. Retrieved from [http://dx.doi.org/10.](http://dx.doi.org/10.1111/faf.12108)
473 [1111/faf.12108](http://dx.doi.org/10.1111/faf.12108)
- 474 Ficetola, G.F., Rondinini, C., Bonardi, A., Baisero, D. & Padoa-Schioppa, E. (2014). Habitat availability for
475 amphibians and extinction threat: A global analysis (D. Richardson, Ed.). *Diversity and Distributions*,
476 **21**, 302–311. Retrieved from <http://dx.doi.org/10.1111/ddi.12296>
- 477 Herring, J. (2011). OpenGIS implementation standard for geographic information-simple feature access-part
478 1: Common architecture. *OGC Document*, **4**, 122–127. Retrieved from [http://www.opengeospatial.org/](http://www.opengeospatial.org/standards/sfa)
479 [standards/sfa](http://www.opengeospatial.org/standards/sfa)

- Janssens, S.B., Vandeloek, F., Langhe, E.D., Verstraete, B., Smets, E., Vandenhouwe, I. & Swennen, R. (2016). Evolutionary dynamics and biogeography of musaceae reveal a correlation between the diversification of the banana family and the geological and climatic history of southeast asia. *New Phytologist*, **210**, 1453–1465. Retrieved from <http://dx.doi.org/10.1111/nph.13856>
- Kong, X., Huang, M. & Duan, R. (2015). SDMdata: A web-based software tool for collecting species occurrence records (J.H. Badger, Ed.). *PLOS ONE*, **10**, e0128295. Retrieved from <http://dx.doi.org/10.1371/journal.pone.0128295>
- Malhado, A.C.M., Oliveira-Neto, J.A., Stropp, J., Strona, G., Dias, L.C.P., Pinto, L.B. & Ladle, R.J. (2015). Climatological correlates of seed size in amazonian forest trees (T. Nakashizuka, Ed.). *J Veg Sci*, **26**, 956–963. Retrieved from <http://dx.doi.org/10.1111/jvs.12301>
- María Mendoza, Ospina, O.E., Cárdenas-Henao, H. & García-R, J.C. (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. Retrieved from <http://dx.doi.org/10.1016/j.ympev.2015.02.001>
- Pimm, S.L., Jenkins, C.N., Abell, R., Brooks, T.M., Gittleman, J.L., Joppa, L.N., Raven, P.H., Roberts, C.M. & Sexton, J.O. (2014). The biodiversity of species and their rates of extinction, distribution, and protection. *Science*, **344**, 1246752–1246752. Retrieved from <http://dx.doi.org/10.1126/science.1246752>
- R Core Team. (2014). *R: A language and environment for statistical computing*. R Foundation for Statistical Computing, Vienna, Austria. Retrieved from <http://www.R-project.org/>
- Richardson, D.M., Roux, J.J.L. & Wilson, J.R. (2015). Australian acacias as invasive species: Lessons to be learnt from regions with long planting histories. *Southern Forests: a Journal of Forest Science*, **77**, 31–39. Retrieved from <http://dx.doi.org/10.2989/20702620.2014.999305>
- Robertson, M.P., Visser, V. & Hui, C. (2016). Biogeo: An r package for assessing and improving data quality of occurrence record datasets. *Ecography*, **39**, 394–401. Retrieved from <http://dx.doi.org/10.1111/ecog.02118>
- Turner, K.G., Fréville, H. & Rieseberg, L.H. (2015). Adaptive plasticity and niche expansion in an invasive thistle. *Ecol Evol*, **5**, 3183–3197. Retrieved from <http://dx.doi.org/10.1002/ece3.1599>
- Verheijen, L.M., Aerts, R., Bönsch, G., Kattge, J. & Bodegom, P.M.V. (2015). Variation in trait trade-offs allows differentiation among predefined plant functional types: Implications for predictive ecology. *New Phytologist*, **209**, 563–575. Retrieved from <http://dx.doi.org/10.1111/nph.13623>

510 Werner, G.D.A., Cornwell, W.K., Cornelissen, J.H.C. & Kiers, E.T. (2015). Evolutionary signals of symbiotic
511 persistence in the legumerhizobia mutualism. *Proceedings of the National Academy of Sciences*, **112**,
512 10262–10269. Retrieved from <http://dx.doi.org/10.1073/pnas.1424030112>
513 Zizka, A. & Antonelli, A. (2015). *speciesgeocodeR: An r package for linking species occurrences, user-defined*
514 *regions and phylogenetic trees for biogeography, ecology and evolution*. Cold Spring Harbor Laboratory
515 Press. Retrieved from <http://dx.doi.org/10.1101/032755>