

1 rgbif: R client for working with GBIF species occurrence data

2 Scott Chamberlain^{*,a}

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

4 **Abstract**

- 5 1. xxx
- 6 2. xxx
- 7 3. xxx
- 8 4. xxxx

9 **Introduction**

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

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

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

25 There are many online services that collect and maintain specimen records. However, Global Biodiversity
26 Information Facility (hereafter, GBIF, <http://www.gbif.org/>) is the largest collection of biodiversity

*Corresponding author

Email address: `scott(at)ropensci.org` (Scott Chamberlain)

27 records globally, currently with 580 million records, 1.6 million taxa, 15,000 datasets from 770 publishers
28 (figures collected on 2015-10-04). Many large biodiversity warehouses such as iNaturalist (<http://www.inaturalist.org/>),
29 VertNet (<http://vertnet.org/>), and USGS's Biodiversity Information Serving Our
30 Nation (BISON; <http://bison.usgs.ornl.gov/>) all feed into GBIF.

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

36 **The `rgbif` package**

37 The `rgbif` package is completely written in R, uses an [MIT license](#) to maximize use everywhere. `rgbif`
38 is developed publicly on GitHub at <https://github.com/ropensci/rgbif>, where development versions of
39 the package can be installed, and bugs and feature requests reported. Stable versions of `rgbif` can be
40 installed from [CRAN](#), the distribution network for R packages. `rgbif` is part of the rOpenSci project,
41 a developer network making R software to facilitate reproducible research.

42 *Package interface*

43 `rgbif` is designed following the [GBIF Application Programming Interface](#), or API. The GBIF API has
44 four major components: registry, species names, occurrence data, and maps. We ignore maps in `rgbif`
45 as it is concerned with generating maps for web applications. `rgbif` has a suite of functions dealing
46 with each of registry, species names, and occurrence data - we'll go through each in turn describing
47 design and example usage.

48 *Registry*

49 The GBIF registry API services are spread across four sets of functions:

- 50 • Datasets
- 51 • Installations
- 52 • Networks

- 53 • Nodes
- 54 • Organizations

55 *Datasets*

56 Search for datasets

```
res <- dataset_search(query = "oregon")
res$data$datasetTitle[1:10]
#> [1] "SDNHM Birds Collection"
#> [2] "CM Birds Collection"
#> [3] "condoncollection"
#> [4] "Taxonomy in Flux Checklist"
#> [5] "Wool carder bees of the genus Anthidium in the Western Hemisphere"
#> [6] "Bryophyte Collection - University of Washington Herbarium (WTU)"
#> [7] "University of British Columbia Herbarium (UBC) - Bryophytes Collection"
#> [8] "UWFC Ichthyology Collection"
#> [9] "Lichen Collection - University of Washington Herbarium (WTU)"
#> [10] "UWBM Mammalogy Collection"
```

57 Get dataset metrics

```
res <- dataset_metrics(uuid='66dd0960-2d7d-46ee-a491-87b9adcfe7b1')
df <- data.frame(rank = names(res$countByRank),
                 count = unname(unlist(res$countByRank)))
knitr::kable(df)
```

rank	count
SPECIES	52452
GENUS	12930
VARIETY	4806
SUBSPECIES	4440
SERIES	1079

rank	count
TRIBE	844
FAMILY	509
SUBTRIBE	327
SUBFAMILY	303
SUBGENUS	241
FORM	239
SECTION	82
SUBVARIETY	4
KINGDOM	1

58 *Networks, nodes, and installations*

59 Here, we search for the first give GBIF networks, returning just the key and title fields.

```
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)"
```

60 *Species*

61 *Occurrences*

62 GBIF provides two ways to get occurrence data: through the `/occurrence/search` route (see
63 `occ_search()`), or via the `/occurrence/download` route (many functions, see below). `occ_search()` is

64 the main function for the search route, and is more appropriate for smaller data, while `occ_download*()`
65 functions are more appropriate for larger data requests.

66 Large is of course a subjective term. When you hit a “large dataset” will depend primarily on the size
67 of the your data request. GBIF imposes for any given search a limit of 200,000 records in the search
68 service, after which point you can’t download any more records for that search. However, you can
69 download more records for different searches.

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

74 *Download API*

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

79 Using the download service can be as few as three steps: 1) Request data via a search; 2) Download
80 data; 3) Import data into R.

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

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

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

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

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

88 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 decimalLatitude decimalLongitude
#> 1  657590544             NA              NA
#> 2  657679551             NA              NA
#> 3  657791316      37.70805      -118.4162
#> 4  658180562             NA              NA
#> 5  441881672             NA              NA
#> 6  911596181             NA              NA
#> 7   56454601             NA              NA
#> 8  657848913             NA              NA
#> 9  658187373             NA              NA
#> 10 658279212      38.95917      -106.9892
#> ..      ...      ...
```

89 *Downloaded data format.* The downloaded dataset from GBIF is actually a Darwin Core Archive
90 (DwC-A), an internationally recognized biodiversity informatics standard (<http://rs.tdwg.org/dwc/>).
91 The DwC-A downloaded is a compressed folder with a number of files, including metadata, citations for
92 each of the datasets included in the download, and the data itself, in separate files for each dataset as
93 well as one single `.txt` file. In `occ_download_import()`, we simply fetch data from the `.txt` file. If you
94 want to dig into the metadata, citations, etc., it is easily accessible from the folder on your computer.

95 *Search API*

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

- 97 • `occ_count()`
- 98 • `occ_search()`
- 99 • `occ_get()`
- 100 • `occ_metadata()`

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

104 ...

105 *Cleaning data.* GBIF provides optional data issues with each occurrence record. These issues fall into
106 many different pre-defined classes, covering issues with taxonomic names, geographic data, and more
107 (see `occ_issues_lookup()` to find out more information on GBIF issues; and the same data on [GBIF's](#)
108 [development site](#)).

109 `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] 2
```

110 *Use cases*

111 *A*

112 *B*

113 *C*

114 **Conclusions and future directions**

- 115 • pt 1

- pt 2
- pt 3
- pt 4

Acknowledgements

This project was supported in part by the Alfred P Sloan Foundation (Grant 2013-6-22).

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 github.com/ropensci/rgbif, under the open and permissive 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. *Rgbif: Interface to the global 'biodiversity' information facility 'aPI'*.
- Faulkner KT., Robertson MP., Rouget M., Wilson JR. 2014. A simple, rapid methodology for developing invasive species watch lists. *Biological Conservation* 179:25–32.

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

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

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

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

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

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