

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
16 of research questions. In addition, the pool of questions we can answer becomes much larger with more
17 and better data. In addition to wide utility, this data is important for conservation. Biodiversity loss is
18 one of the greatest challenges of our time (Pimm et al., 2014). Some have called this the sixth great
19 mass extinction (Ceballos et al., 2015). Given this challenge there is a great need for data on specimen
20 records, whether collected from live sightings in the field or specimens in museums.

21 There are many online services that collect and maintain specimen records. However, Global Biodiversity
22 Information Facility (hereafter, GBIF, <http://www.gbif.org/>) is the largest collection of biodiversity
23 records globally, currently with 580 million records, 1.6 million taxa, 15,000 datasets from 770 publishers
24 (figures collected on 2015-10-04). Many large biodiversity warehouses such as iNaturalist, VertNet, and
25 USGS's BISON all feed into GBIF.

26 Herein, we describe a library (rgbif (Chamberlain et al.)) for working with GBIF data in the R
27 programming environment (R Core Team, 2014). R is an extremely widely used language in academia,

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

28 and in non-profit and private sectors. Importantly, R makes it easy to do all of the steps of the research
29 process, including data management, data manipulation and cleaning, statistics, and vizualization.
30 Thus, an R client for getting GBIF data is a powerful tool for reproducible research.

31 **The rgbif package**

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

37 *Package interface*

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

43 *Registry*

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

- 45 • Datasets
- 46 • Installations
- 47 • Networks
- 48 • Nodes
- 49 • Organizations

50 *Datasets*

51 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"

```

52 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
TRIBE	844
FAMILY	509
SUBTRIBE	327
SUBFAMILY	303
SUBGENUS	241
FORM	239

rank	count
SECTION	82
SUBVARIETY	4
KINGDOM	1

53 *Networks, nodes, and installations*

54 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)"
```

55 *Species*

56 *Occurrences*

57 GBIF provides two ways to get occurrence data: through the `/occurrence/search` route (see
58 `occ_search`), or via the `/occurrence/download` route (many functions, see below). `occ_search()` is
59 the main funtion for the search route, and is more appropriate for smaller data, while `occ_download*()`
60 functions are more appropriate for larger data requests.

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

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

69 *Download API*

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

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

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

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

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

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

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

83 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
#> ..      ...      ...

```

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

90 *Search API*

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

- 92 • `occ_count()`
- 93 • `occ_search()`
- 94 • `occ_get()`
- 95 • `occ_metadata()`

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

99 ...

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

104 `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
```

105 *Use cases*

106 *A*

107 *B*

108 *C*

109 **Conclusions and future directions**

- 110 • pt 1
- 111 • pt 2
- 112 • pt 3
- 113 • pt 4

114 *Acknowledgements*

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

116 *Data Accessibility*

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

121 **References**

- 122 Beck J., Ballesteros-Mejia L., Buchmann CM., Dengler J., Fritz SA., Gruber B., Hof C., Jansen
123 F., Knapp S., Kreft H., Schneider A-K., Winter M., Dormann CF. 2012. Whats on the horizon for
124 macroecology? *Ecography* 35:673–683.
- 125 Brown JH. 1995. *Macroecology*. University of Chicago Press.
- 126 Ceballos G., Ehrlich PR., Barnosky AD., Garcia A., Pringle RM., Palmer TM. 2015. Accelerated
127 modern human-induced species losses: Entering the sixth mass extinction. *Science Advances* 1:e1400253–
128 e1400253.
- 129 Chamberlain S., Ram K., Barve V., Mcglinn D. *Rgbif: Interface to the global 'biodiversity' information*
130 *facility 'aPI'*.
- 131 Pimm SL., Jenkins CN., Abell R., Brooks TM., Gittleman JL., Joppa LN., Raven PH., Roberts CM.,
132 Sexton JO. 2014. The biodiversity of species and their rates of extinction, distribution, and protection.
133 *Science* 344:1246752–1246752.
- 134 R Core Team. 2014. *R: A language and environment for statistical computing*. Vienna, Austria: R
135 Foundation for Statistical Computing.