

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]
```

```
52 [1] "SDNHM Birds Collection"
53 [2] "CM Birds Collection"
54 [3] "condoncollection"
55 [4] "Taxonomy in Flux Checklist"
56 [5] "Wool carder bees of the genus Anthidium in the Western Hemisphere"
57 [6] "Bryophyte Collection - University of Washington Herbarium (WTU)"
58 [7] "University of British Columbia Herbarium (UBC) - Bryophytes Collection"
59 [8] "UWFC Ichthyology Collection"
60 [9] "Lichen Collection - University of Washington Herbarium (WTU)"
61 [10] "UWBM Mammalogy Collection"
```

```
62 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

rank	count
FORM	239
SECTION	82
SUBVARIETY	4
KINGDOM	1

63 *Networks, nodes, and installations*

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

```
networks(limit=10)$data$title
```

```
65 [1] "GBIF Backbone Sources"
66 [2] "Canadensys"
67 [3] "Southwest Collections of Arthropods Network (SCAN)"
68 [4] "VertNet"
69 [5] "Dryad"
70 [6] "GBIF Network"
71 [7] "The Knowledge Network for Biocomplexity (KNB) "
72 [8] "Online Zoological Collections of Australian Museums (OZCAM)"
73 [9] "Catalogue of Life"
74 [10] "Ocean Biogeographic Information System (OBIS)"
```

75 *Species*

76 *Occurrences*

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

81 Large is of course a subjective term. When you hit a “large dataset” will depend primarily on the size
82 of the your data request. GBIF imposes for any given search a limit of 200,000 records in the search

83 service, after which point you can't download any more records for that search. However, you can
84 download more records for different searches.

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

89 *Download API*

90 The download API syntax is similar to the occurrence search API in that the same parameters are used,
91 but the way in which the query is defined is different. For example, in the download API you can do
92 greater than searches, whereas you can not do that in the occurrence search API. Thus, we can't make
93 the query interface exactly the same for both search and download functions.

94 Using the download service can be as few as three steps.

95 Request data download given a query. Here, we xxxx

```
"xxx"
```

```
[1] "xxx"
```

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

```
(res <- occ_download_get("0000066-140928181241064", overwrite = TRUE))
```

```
100 <<gbif downloaded get>>
```

```
101   Path: ./0000066-140928181241064.zip
```

```
102   File size: 0.14 MB
```

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

105 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
...

Search API

The search API is very similar to the download API, but is meant for smaller data acquisition jobs.

Conclusions and future directions

Acknowledgements

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

Data Accessibility

All software, 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.

127 References

- 128 Beck, J., Ballesteros-Mejia, L., Buchmann, C.M., Dengler, J., Fritz, S.A., Gruber, B., Hof, C., Jansen,
129 F., Knapp, S., Kreft, H., Schneider, A.-K., Winter, M. & Dormann, C.F. (2012). Whats on the
130 horizon for macroecology? *Ecography*, **35**, 673–683. Retrieved from [http://dx.doi.org/10.1111/j.1600-](http://dx.doi.org/10.1111/j.1600-0587.2012.07364.x)
131 [0587.2012.07364.x](http://dx.doi.org/10.1111/j.1600-0587.2012.07364.x)
- 132 Brown, J.H. (1995). *Macroecology*. University of Chicago Press.
- 133 Ceballos, G., Ehrlich, P.R., Barnosky, A.D., Garcia, A., Pringle, R.M. & Palmer, T.M. (2015). Acceler-
134 ated modern human-induced species losses: Entering the sixth mass extinction. *Science Advances*, **1**,
135 e1400253–e1400253. Retrieved from <http://dx.doi.org/10.1126/sciadv.1400253>
- 136 Chamberlain, S., Ram, K., Barve, V. & Mcglinn, D. *Rgbif: Interface to the global 'biodiversity'*
137 *information facility 'aPI'*. Retrieved from <https://github.com/ropensci/rgbif>
- 138 Pimm, S.L., Jenkins, C.N., Abell, R., Brooks, T.M., Gittleman, J.L., Joppa, L.N., Raven, P.H., Roberts,
139 C.M. & Sexton, J.O. (2014). The biodiversity of species and their rates of extinction, distribution, and
140 protection. *Science*, **344**, 1246752–1246752. Retrieved from <http://dx.doi.org/10.1126/science.1246752>
- 141 R Core Team. (2014). *R: A language and environment for statistical computing*. R Foundation for
142 Statistical Computing, Vienna, Austria. Retrieved from <http://www.R-project.org/>