

taxadb: A High-Performance Local Taxonomic Database Interface

Kari E A Norman^a, Scott Chamberlain^b, Carl Boettiger^{a,1}

^a*Dept of Environmental Science, Policy, and Management, University of California Berkeley, Berkeley CA 94720-3114, USA*

^b*The rOpenSci Project, University of California Berkeley, Berkeley CA 94720-3114, USA*

Abstract

A familiar and growing challenge in ecological and evolutionary research is that of establishing consistent taxonomy when combining data from separate sources. While this problem is already well understood and numerous naming authorities have been created to address the issue, most researchers lack a fast, consistent, and intuitive way to retrieve taxonomic names. Here, we present **taxadb** R package to address this gap. In contrast to existing tools, **taxadb** provides all of the following in one place:

- 1) **taxadb** accesses established naming authorities to resolve synonyms, taxonomic identifiers, and hierarchical classification.
- 2) **taxadb** creates a local database, managed automatically from within R, that provides fast operations on millions of taxonomic names.
- 3) **taxadb** provides a consistent and intuitive data format.
- 4) **taxadb** is built on a simple, extensible, and language-agnostic design that can easily accommodate new authorities.

As ecologists and evolutionary biologists synthesize datasets across larger and larger assemblies of species, we face a continual challenge of maintaining consistent taxonomy. How many species are in the combined data? Do the studies use the same names for the same species, or do they use different synonyms for the same species? Failing to correct for such differences can lead to significant inflation of species counts and miss-aligned datasets. These challenges have become particularly acute as it becomes increasingly common for researchers to work across a larger number and diversity of species in any given analysis, which may preclude the resources or substantive taxonomic expertise for all clades needed to resolve scientific names (Patterson et al. 2010).

While these issues have long been recognized in the literature (Boyle et al. 2013; Dayrat 2005; Bortolus 2008; Maldonado et al. 2015; Remsen 2016), and a growing number of databases and tools have emerged over the past few decades (e.g. ITIS 2019; Biotechnology Information 2019; Roskov Y. 2018; Rees 2014; Alvarez and Luebert 2018; Wagner 2016; Foster, Chamberlain, and Grünwald 2018; Gries, Gilbert, and Franz 2014), it remains difficult to resolve taxonomic names to a common authority in a transparent, efficient, and automatable manner. Here, we present an R package, **taxadb**, which seeks to address this gap.

Databases of taxonomic names such as the Integrated Taxonomic Information System (ITIS; ITIS 2019), the National Center for Biological Information's (NCBI) Taxonomy database, (Biotechnology Information 2019), the Catalogue of Life (COL; Roskov Y. 2018), and over one hundred other providers have sought to address these problems by providing expert-curated lists of accepted taxonomic names, synonyms, associated taxonomic rank, hierarchical classifications, and scientific authority (e.g. author and date) establishing a scientific name. The R language (R Core Team 2019) is widely used in ecology and evolution (Lai et al. 2019) and the **taxize** package (Chamberlain and Szöcs 2013) has become a popular way for R users to interact with naming providers and name resolution services. **taxize** implements bindings to the web APIs (Application Programming Interface) hosted by many popular taxonomic name providers. Nevertheless, this means that functions in the **taxize** are impacted by several major drawbacks that are inherent in the implementation of these central API servers, such as:

- Queries require internet access at all times.

- Queries are slow and inefficient to implement and perform; frequently requiring separate API calls for each taxonomic name.
- The type of query is highly limited by the API design. For instance, it is usually impossible to make queries across the entire corpus of names, such as “which accepted name has the most known synonyms?”.
- Both query formats and responses differ substantially across different naming providers, making it difficult to apply a script designed for one provider to different provider.
- Most queries are not reproducible, as the results depend on the state of the central server (and potentially the quality of the internet connection)(Rees and Cranston 2017). Many names providers update the server data either continuously or at regular intervals, including both revising existing names (for spelling or changes in accepted name designation) and adding new names.

Instead of binding existing web APIs, **taxadb** is built around a set of compressed text files which are automatically downloaded, imported, and stored on a local database by **taxadb**. The largest of the taxonomic naming providers today contain under 6 million name records with uncompressed file sizes over a GB, which can be compressed to around 50 MB and downloaded in under a minute on a 1 MB/s connection. By using a local database as the backend, **taxadb** allows R users to interact with large data files without large memory (RAM) requirements. A query for a single name over the web API requires a remote server to respond, execute the query, and serialize the response, which can take several seconds. Thus it does not take many taxa before transferring the entire data set to query locally is more efficient. Moreover, this local copy can be cached on the user’s machine, requiring only the one-time setup, and enabling offline use and reproducible queries. Rather than returning data in whatever format is given by the provider, **taxadb** provides a data structure following a consistent, standardized layout or schema following Darwin Core, which provides standard terms for biodiversity data (Wieczorek et al. 2012). Table 1 summarizes the list of all naming providers currently accessed by **taxadb**. More details are provided in the Data Sources Vignette, <https://docs.ropensci.org/taxadb/articles/data-sources.html>.

Table 1: Descriptions of the providers supported by **taxadb** with their reference abbreviation and the total number of identifiers contained by each provider.

| Provider | Abbreviation | Number of Identifiers | Description |
|--|--------------|-----------------------|--|
| Integrated Taxonomic Information System (ITIS 2019) | itis | 597120 | originally formed to standardize taxonomic name usage across many agencies in the United States federal government |
| National Center for Biological Information’s Taxonomy database (NCBI) | ncbi | 188221 | nomenclature for sequences in the International Nucleotide Sequence Database Collaboration database |
| Catalogue of Life (COL) | col | 1998435 | comprehensive taxonomic effort, includes some other providers (e.g. itis) |
| Global Biodiversity Information Facility Taxonomic Backbone (GBIF 2019) | gbif | 3546672 | taxonomic backbone of the GBIF database, assembled from other sources including COL |
| FishBase (Froese and Pauly 2019) | fb | 34299 | nomenclature for global database of fishes |
| Open Tree Taxonomy (OTT) | ott | 4455820 | comprehensive tree of life based on phylogenetic trees and taxonomic data |
| International Union for Conservation of Nature and Natural Resources (IUCN 2019) | iucn | 131927 | taxonomy for classification of species status |

Package Overview

```
library(tidyverse)
library(taxadb)
```

After loading our package and the tidyverse package for ease in manipulating function output, we look up the taxonomic identifier for Atlantic Cod, *Gadus morhua*, and the compliment:

```
get_ids("Gadus morhua")
```

```
## [1] "ITIS:164712"
```

```
get_names("ITIS:164712")
```

```
## [1] "Gadus morhua"
```

Our first call to any `taxadb` functions will automatically set up a local, persistent database if one has not yet been created. This one-time setup will download, extract, and import the compressed data into persistent database storage (using the appropriate location specified by the operating system (see Ratnakumar, Mick, and Davis 2016), or configured using the environmental variable `TAXADB_HOME`). The example above searches for names in ITIS, the default provider, which can be configured using the `provider` argument. Any future function calls to this function or any other function using data from the same provider will be able to access this data rapidly without the need for processing or an internet connection.

Users can also explicitly trigger this one-time setup using `td_create()` and specifying the provider abbreviation (see Table 1), or simply using `all` to install all available providers:

```
td_create("all")
```

`taxadb` functions like `get_ids()` and `td_create()` take an optional argument, `db`, to an external database connection. `taxadb` will work with most DBI-compliant databases such as MySQL or Postgres, but will be much faster when using a column-oriented database engine such as `duckdb` or `MonetDBLite`. These latter options are also much easier for most users, since each can be installed directly as an R package. `taxadb` will default to the fastest available option. `taxadb` can also run without a database backend by setting `db=NULL`, though some functions will require a lot (2-20 GB) of free RAM for this to work with many of the larger providers.

`taxadb` uses the widely known SQLite database by default, but users are encouraged to install the optional, suggested database backends by passing the option `dependencies = TRUE` to the install command. This installs a `MonetDBLite` database instance (Raasveldt and Mühleisen 2018), a columnar-oriented relational database requiring no additional installation while also providing persistent disk-based storage. This also installs `duckdb`, another local columnar database which is rapidly emerging as an alternative to `MonetDB` and `SQLite`. `taxadb` will automatically detect and use these database engines if available, and automatically handles opening, caching, and closing the database connection. For large queries, `MonetDBLite` or `duckdb` deliver impressive improvements. Our benchmark on resolving the 750 species names in the Breeding Bird Survey against over 3 million names known in the 2019 Catalogue of Life takes 8 minutes in `SQLite` but less than a second in `MonetDBLite`.

Functions in `taxadb` are organized into several families:

- queries that return vectors: `get_ids()` and it's complement, `get_names()`,
- queries that filter the underlying taxonomic data frames: `filter_name()`, `filter_rank()`, `filter_id()`, and `filter_common()`,
- database functions `td_create()`, `td_connect()` and `taxa_tbl()`,
- and helper utilities, such as `clean_names()`.

Taxonomic Identifiers

Taxonomic identifiers provide a fundamental abstraction which lies at the heart of managing taxonomic names. For instance, by resolving scientific names to identifiers, we can identify which names are synonyms –

different scientific names used to describe the same species – and which names are not recognized. Each naming authority provides its own identifiers for the names it recognizes. For example, the name *Homo sapiens* has the identifier 9606 in NCBI and 180092 in ITIS. To avoid possible confusion, `taxadb` always prefixes the naming provider, e.g. NCBI:9606. Some taxonomic naming providers include separate identifiers for synonyms, see Box 1. Unmatched names may indicate an error in data entry or otherwise warrant further investigation. Taxon identifiers are also easily resolved to the original authority (scientific publication) establishing the name. The common practice of appending an author and year to a scientific name, e.g. *Poa annua annua* (Smith 1912), serves a valuable role in disambiguating different uses of the same name but can be notoriously harder to resolve to the appropriate reference, while variation in this convention creates many distinct versions of the same name (Patterson et al. 2010).

These issues are best illustrated using a minimal example. We'll consider the task of combining data on bird extinction risk as assessed by the IUCN (International Union for Conservation of Nature and Natural Resources 2019) with data on average adult biomass, as estimated in the Elton Traits v1.0 database (Wilman et al. 2016). To keep the example concise enough for visual presentation we will focus on a subset involving just 10 species (Table 2, 3).

```
trait_data <- read_tsv(system.file("extdata", "trait_data.tsv", package="taxadb"))
status_data <- read_tsv(system.file("extdata", "status_data.tsv", package="taxadb"))
```

Table 2: The subset of the IUCN status data used for subsequent taxonomic identifier examples.

| <i>iucn_name</i> | category |
|-----------------------------------|----------|
| <i>Pipile pipile</i> | CR |
| <i>Pipile cumanensis</i> | LC |
| <i>Pipile cunjubi</i> | LC |
| <i>Pipile jacutinga</i> | EN |
| <i>Megapodius decollatus</i> | LC |
| <i>Scleroptila gutturalis</i> | LC |
| <i>Margaroperdix madagarensis</i> | LC |
| <i>Falciipennis falciipennis</i> | NT |

Table 3: The subset of the Elton trait data used for subsequent taxonomic identifier examples.

| <i>elton_name</i> | mass |
|---------------------------------------|---------|
| <i>Aburria pipile</i> | 1816.59 |
| <i>Aburria cumanensis</i> | 1239.22 |
| <i>Aburria cunjubi</i> | 1195.82 |
| <i>Aburria jacutinga</i> | 1240.96 |
| <i>Megapodius reinwardt</i> | 666.34 |
| <i>Francolinus lewalliantoides</i> | 376.69 |
| <i>Margaroperdix madagascariensis</i> | 245.00 |
| <i>Catreus wallichii</i> | 1436.88 |
| <i>Falciipennis falciipennis</i> | 685.61 |
| <i>Falciipennis canadensis</i> | 473.65 |

If we attempted to join these data directly on the species names provided by each table, we would find very little overlap, with only one species name having both a body mass and an IUCN threat status resolved

114 (Table 4).

```
joined <- full_join(trait_data, status_data, by = c("elton_name" = "iucn_name"))
```

Table 4: Example IUCN and trait data joined directly on scientific name showing only one match. While common, joining on scientific name does not account for nomenclatural and taxonomic inconsistencies between databases and therefore results in seemingly very little overlap in species representation between the two.

| <i>elton_name</i> | mass | category |
|---------------------------------------|---------|----------|
| <i>Aburria pipile</i> | 1816.59 | - |
| <i>Aburria cumanensis</i> | 1239.22 | - |
| <i>Aburria cujubi</i> | 1195.82 | - |
| <i>Aburria jacutinga</i> | 1240.96 | - |
| <i>Megapodius reinwardt</i> | 666.34 | - |
| <i>Francolinus lewalliantoides</i> | 376.69 | - |
| <i>Margaroperdix madagascariensis</i> | 245.00 | - |
| <i>Catreus wallichii</i> | 1436.88 | - |
| <i>Falciipennis falciipennis</i> | 685.61 | NT |
| <i>Falciipennis canadensis</i> | 473.65 | - |
| <i>Pipile pipile</i> | - | CR |
| <i>Pipile cumanensis</i> | - | LC |
| <i>Pipile cujubi</i> | - | LC |
| <i>Pipile jacutinga</i> | - | EN |
| <i>Megapodius decollatus</i> | - | LC |
| <i>Scleroptila gutturalis</i> | - | LC |
| <i>Margaroperdix madagarensis</i> | - | LC |

115 If we first resolve names used in each data set into shared identifiers, (for instance, using the Catalogue of
 116 Life), we discover that there is far more overlap in the species coverage than we might have initially realized.
 117 First, we just add an ID column to each table by looking up the Catalog of Life identifier for the name
 118 provided:

```
traits <- trait_data %>% mutate(id = get_ids(elton_name, "col"))
status <- status_data %>% mutate(id = get_ids(iucn_name, "col"))
```

119 We can now join on the id column instead of names directly:

```
joined <- full_join(traits, status, by = "id")
```

120 This results in many more matches (Table 5), as different scientific names are recognized by the naming
 121 provider (Catalog of Life 2018 in this case), as *synonyms* for the same species, and thus resolve to the
 122 same taxonomic identifier. While we have focused on a small example for visual clarity here, the `get_ids()`
 123 function in `taxadb` can quickly resolve hundreds of thousands of species names to unique identifiers, thanks
 124 to the performance of fast joins in a local MonetDBLite database.

Table 5: Example IUCN and trait data joined on taxonomic ID. Multiple species have a different scientific name in the Elton and IUCN Redlist databases but can be match based on their COL taxonomic ID.

| <i>elton_name</i> | <i>iucn_name</i> | mass | category | id |
|---------------------------------------|-----------------------------------|---------|----------|--------------|
| <i>Aburria pipile</i> | <i>Pipile pipile</i> | 1816.59 | CR | COL:35517887 |
| <i>Aburria cumanensis</i> | <i>Pipile cumanensis</i> | 1239.22 | LC | COL:35537158 |
| <i>Aburria cujubi</i> | <i>Pipile cujubi</i> | 1195.82 | LC | COL:35537159 |
| <i>Aburria jacutinga</i> | <i>Pipile jacutinga</i> | 1240.96 | EN | COL:35517886 |
| <i>Megapodius reinwardt</i> | - | 666.34 | - | COL:35521309 |
| <i>Francolinus levalliantoides</i> | - | 376.69 | - | COL:35518087 |
| <i>Margaroperdix madagascariensis</i> | <i>Margaroperdix madagarensis</i> | 245.00 | LC | COL:35521355 |
| <i>Catreus wallichii</i> | - | 1436.88 | - | COL:35518185 |
| <i>Falci pennis falci pennis</i> | <i>Falci pennis falci pennis</i> | 685.61 | NT | COL:35521380 |
| <i>Falci pennis canadensis</i> | - | 473.65 | - | COL:35521381 |
| - | <i>Megapodius decollatus</i> | - | LC | COL:35537166 |
| - | <i>Scleroptila gutturalis</i> | - | LC | - |

Box 1: Taxonomic Identifiers and Synonyms

`get_ids()` returns the `acceptedNameUsageID`, the identifier associated with the *accepted* name. Some naming providers, such as ITIS and NCBI, provide taxonomic identifiers to both synonyms and accepted names. Other providers, such as COL and GBIF, only provide identifiers for accepted names. Common practice in Darwin Core archives is to provide an `acceptedNameUsageID` only for names which are synonyms, and otherwise to provide a `taxonID`. For accepted names, the `acceptedNameUsageID` is then given as missing (NA), while for synonyms, the `taxonID` may be missing (NA). In contrast, `taxadb` lists the `acceptedNameUsageID` for accepted names (where it matches the `taxonID`), as well as known synonyms. This is semantically identical, but also more convenient for database interfaces, since it allows a name to mapped to its accepted identifier (or an identifier to map to it's accepted name usage) without the additional logic. For consistency, we will use the term “identifier” to mean the `acceptedNameUsageID` rather than the more ambiguous `taxonID` (which is undefined for synonyms listed by many providers), unless explicitly stated otherwise.

Unresolved names

`get_ids` offers a first pass at matching scientific names to id, but names may remain unresolved for a number of reasons. First, a name may match to multiple accepted names, as in the case of a species that has been split. By design, these cases are left to be resolved by the researcher using the `filter_` functions to filter underlying taxonomic tables for additional information. A name may also be unresolved due to typos or improper formatting. `clean_names` addresses common formatting issues such as the inclusion of missing species epithets (e.g. `Accipiter sp.`) that prevent matches to the Genus, or intraspecific epithets such as `Colaptes auratus cafer` that prevent matches to the binomial name. These modifications are not appropriate in all settings and should be used with care. Spell check of input names is outside the scope of `taxadb`, however existing tools such as those developed by the Global Names Architecture (<http://globalnames.org/apps/>) could be incorporated into a `taxadb` workflow.

Names may also have an ambiguous resolution wherein a name may be resolved by a different provider than the one specified, either as an accepted name or a synonym. Mapping between providers represent a meaningful scientific statement requiring an understanding of the underlying taxonomic concepts of each provider (Franz and Peet 2009; Franz and Sterner 2018; Lepage, Vaidya, and Guralnick 2014). The spirit

of `taxadb` is not to automate steps that require expert knowledge, but provide access to multiple potential “taxonomic theories”.

filter_ functions for access to underlying tables

Underlying data tables can be accessed through the family of `filter_` functions, which filter by certain attributes such as scientific name, id, common name, and rank. These functions allow us to ask general questions such as, how many bird species are there?

```
filter_rank("Aves", rank="class", provider = "col") %>%
  filter(taxonomicStatus == "accepted", taxonRank == "species") %>%
  pull(taxonID) %>%
  n_distinct()
```

```
## [1] 10354
```

We can also use this to gain a detailed look at specific species or ids. For example, we can explore why `get_ids` fails to resolve a seemingly common species:

```
multi_match <- filter_name("Abies menziesii", provider = "col")
```

Table 6: Some names may not resolve to an identifier using `get_ids()` because they match to more than one accepted ID. In such cases `filter_` functions give further detail, as in the example of *Zosterops rendovae** below which has two accepted ID matches.

| sort | taxonID | scientificName | acceptedNameUsageID | taxonomicStatus | acceptedScientificName |
|------|--------------|------------------------|---------------------|-----------------|------------------------------|
| 1 | COL:18159104 | <i>Abies menziesii</i> | COL:18157974 | synonym | <i>Pseudotsuga menziesii</i> |
| 1 | COL:18160542 | <i>Abies menziesii</i> | COL:18158639 | synonym | <i>Picea pungens</i> |
| 1 | COL:18161226 | <i>Abies menziesii</i> | COL:18158652 | synonym | <i>Picea sitchensis</i> |

We see that *Abies menziesii* is a synonym for three accepted names which the user will have to choose between (Table 6). This is an example of how `taxadb` seeks to provide users with information from existing authorities and names providers, rather than make a potentially arbitrary decision. Because they return `data.frames`, `filter_` functions provide both potential matches. Note that the simpler `get_` functions (`get_ids()`) consider multiple name matches as `NA` for the `id`, making them suitable for automated pipelines where manual resolution of duplicates is not an option.

Direct database access

The full taxonomic record in the database can also be directly accessed by `taxa_tbl()`, allowing for whole-database queries that are not possible through the API or web interface of many providers. For example, we can easily check the coverage of accepted species names in each of the classes of vertebrates within the Catalogue of Life (Table 7):

```
verts <- taxa_tbl("col") %>%
  filter(taxonomicStatus == "accepted", phylum == "Chordata", taxonRank == "species") %>%
  count(class, sort = TRUE)
```

Box 2: Common Names

`taxadb` can also resolve common names to their identifier by mapping common name to the accepted scientific name. Common names have many of the same issues as scientific names but even more frequent (e.g. matching to more than one accepted name, non-standardized formatting). Common names are accessed via `filter_common` which takes a vector of common names. The user can then resolve discrepancies.

Table 7: `taxadb` also provides direct access to the database, allowing `dplyr` or SQL queries which can compute across the entire dataset, such as counting accepted species in all vertebrate classes shown here. This kind of query is effectively impossible in most REST API-based interfaces.

| class | n |
|--------------------|-------|
| Actinopterygii | 32474 |
| Aves | 10354 |
| Reptilia | 10233 |
| Amphibia | 6439 |
| Mammalia | 5852 |
| Ascidiacea | 2925 |
| Elasmobranchii | 1223 |
| Myxini | 81 |
| Thaliacea | 78 |
| Appendicularia | 68 |
| Holocephali | 56 |
| Cephalaspidomorphi | 45 |
| Leptocardii | 30 |
| Sarcopterygii | 8 |

Discussion

Some taxonomic name providers (e.g. OTT, COL, NCBI) offer periodic releases of a static names list, while many other providers (e.g. ITIS, FB, IUCN) offer name data on a rolling basis (i.e. the data returned by a given download URL is updated continuously or at arbitrary intervals without any additional indication if and how that data has changed.) `taxadb`’s `td_create()` function downloads and stores cached snapshots from each provider, which follow an annual release model to support reproducible analyses. All `taxadb` functions that download or access data include an optional argument `version` to indicate which version of the provider data should be used. By default, `taxadb` will determine the latest version available (at the time of writing this is version 2019). Appropriate metadata is stored with each snapshot, including scripts used to access and reformat the data files, as described in the “Data Sources” vignette, <https://docs.ropensci.org/taxadb/articles/data-sources.html>.

Taxonomic identifiers are an essential first step for maintaining taxonomic consistency, a key task for a wide variety of applications. Despite multiple taxonomic standardization efforts, resolving names to taxonomic identifiers is often not a standard step in the research work flow due to difficulty in accessing providers and the time consuming API queries necessary for resolving even moderately sized data sets. `taxadb` fills an important gap between existing tools and typical research patterns by providing a fast, reproducible approach for matching names to taxonomic identifiers. It could also be used to verify that conclusions were robust to the choice of naming provider. `taxadb` is not intended as an improvement or replacement for any existing approaches to taxonomic name resolution. In particular, `taxadb` is not a replacement for the APIs or databases provided, but merely an interface to taxonomic naming information contained within that data.

Lastly, we note that local database design used in `taxadb` is not unique to taxonomic names. Despite the rapid expansion of REST API-based interfaces to ecological data (Boettiger et al. 2015), in our experience, much of the data relevant to ecologists and evolutionary biologists today would be also be amenable to the local database design. The local database approach is much easier for data providers (who can leverage static scientific database repositories instead of maintaining REST servers) and often much faster for data consumers.

Acknowledgments

We thank the many researchers who contributed to the data and infrastructure of the various taxonomic providers we access through our package. Support for the development of this package was provided by United States Department of Energy through the Computational Sciences Graduate Fellowship (DOE CSGF) under grant number DE-FG02-97ER25308 awarded to K.E.A.N..

Data Availability

Code for the R package can be found on GitHub at <https://github.com/ropensci/taxadb>. The taxonomic database is also stored on Github at <https://github.com/boettiger-lab/taxadb-cache>. The original taxonomic data is stored by the individual provider and can be found through the included citations.

Authors' Contributions

K.E.A.N., S.C., and C.B. contributed to conceptual development of the package. K.E.A.N. and C.B. developed the package and contributed to the manuscript.

References

- Alvarez, Miguel, and Federico Luebert. 2018. "The Taxlist Package: Managing Plant Taxonomic Lists in R." *Biodiversity Data Journal*, no. 6 (May). <https://doi.org/10.3897/BDJ.6.e23635>.
- Biotechnology Information, National Center for. 2019. "NCBI Taxonomy." 8600 Rockville Pike, Bethesda MD, 20894 USA: U.S. National Library of Medicine. <https://www.ncbi.nlm.nih.gov/taxonomy>.
- Boettiger, Carl, Scott Chamberlain, Edmund Hart, and Karthik Ram. 2015. "Building Software, Building Community: Lessons from the rOpenSci Project." *Journal of Open Research Software* 3 (1): e8. <https://doi.org/10.5334/jors.bu>.
- Bortolus, Alejandro. 2008. "Error Cascades in the Biological Sciences: The Unwanted Consequences of Using Bad Taxonomy in Ecology." *AMBIO: A Journal of the Human Environment* 37 (2): 114–18. [https://doi.org/10.1579/0044-7447\(2008\)37\[114:ECITBS\]2.0.CO;2](https://doi.org/10.1579/0044-7447(2008)37[114:ECITBS]2.0.CO;2).
- Boyle, Brad, Nicole Hopkins, Zhenyuan Lu, Juan Antonio Raygoza Garay, Dmitry Mozzherin, Tony Rees, Naim Matasci, et al. 2013. "The Taxonomic Name Resolution Service: An Online Tool for Automated Standardization of Plant Names." *BMC Bioinformatics* 14 (1): 16. <https://doi.org/10.1186/1471-2105-14-16>.
- Chamberlain, Scott A., and Eduard Szöcs. 2013. "Taxize: Taxonomic Search and Retrieval in R." *F1000Research* 2 (October): 191. <https://doi.org/10.12688/f1000research.2-191.v2>.
- Dayrat, Benoît. 2005. "Towards Integrative Taxonomy." *Biological Journal of the Linnean Society* 85 (3): 407–17. <https://doi.org/10.1111/j.1095-8312.2005.00503.x>.
- Foster, Zachary S.L., Scott Chamberlain, and Niklaus J. Grünwald. 2018. "Taxa: An R Package Implementing Data Standards and Methods for Taxonomic Data." *F1000Research* 7 (September). <https://doi.org/10.12688/f1000research.14013.2>.
- Franz, Nico M, and Beckett W Sterner. 2018. "To Increase Trust, Change the Social Design Behind Aggregated Biodiversity Data." *Database* 2018 (January). <https://doi.org/10.1093/database/bax100>.
- Franz, N. M., and R. K. Peet. 2009. "Perspectives: Towards a Language for Mapping Relationships Among Taxonomic Concepts." *Systematics and Biodiversity* 7 (1): 5–20. <https://doi.org/10.1017/S147720000800282X>.
- Gries, Corinna, Edward Gilbert, and Nico Franz. 2014. "Symbiota – A Virtual Platform for Creating Voucher-Based Biodiversity Information Communities." *Biodiversity Data Journal* 2 (June): e1114. <https://doi.org/10.3897/BDJ.2.e1114>.
- International Union for Conservation of Nature and Natural Resources. 2019. "The IUCN Red List of Threatened Species." <https://www.iucnredlist.org>.
- ITIS. 2019. "Integrated Taxonomic Information System." <https://www.itis.gov>.
- Lai, Jiangshan, Christopher J. Lortie, Robert A. Muenchen, Jian Yang, and Keping Ma. 2019. "Evaluating the Popularity of R in Ecology." *Ecosphere* 10 (1): e02567. <https://doi.org/10.1002/ecs2.2567>.
- Lepage, Denis, Gaurav Vaidya, and Robert Guralnick. 2014. "Avibase – a Database System for Managing and Organizing Taxonomic Concepts." *ZooKeys* 420 (June): 117–35. <https://doi.org/10.3897/zookeys.420.7089>.
- Maldonado, Carla, Carlos I. Molina, Alexander Zizka, Claes Persson, Charlotte M. Taylor, Joaquina Albán, Eder Chilquillo, Nina Rønsted, and Alexandre Antonelli. 2015. "Estimating Species Diversity and Distribution in the Era of Big Data: To What Extent Can We Trust Public Databases?" *Global Ecology and Biogeography* 24 (8): 973–84. <https://doi.org/10.1111/geb.12326>.
- Patterson, D.J., J. Cooper, P.M. Kirk, R.L. Pyle, and D.P. Remsen. 2010. "Names Are Key to the Big New Biology." *Trends in Ecology & Evolution* 25 (12): 686–91. <https://doi.org/10.1016/j.tree.2010.09.004>.
- Raasveldt, Mark, and Hannes Mühleisen. 2018. "MonetDBLite : An Embedded Analytical Database." In *Proceedings of Cikm 2018 International Conference on Information and Knowledge Management (Cikm'18)*. New York, New York, USA: ACM. https://doi.org/10.475/123_4.
- Ratnakumar, Sridhar, Trent Mick, and Trevor Davis. 2016. *Rappdirs: Application Directories: Determine Where to Save Data, Caches, and Logs*. <https://CRAN.R-project.org/package=rappdirs>.
- R Core Team. 2019. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.

- Rees, Jonathan, and Karen Cranston. 2017. "Automated Assembly of a Reference Taxonomy for Phylogenetic Data Synthesis." *Biodiversity Data Journal* 5 (May): e12581. <https://doi.org/10.3897/BDJ.5.e12581>.
- Rees, Tony. 2014. "Taxamatch, an Algorithm for Near ('Fuzzy') Matching of Scientific Names in Taxonomic Databases." *PLOS ONE* 9 (9): e107510. <https://doi.org/10.1371/journal.pone.0107510>.
- Remsen, David. 2016. "The Use and Limits of Scientific Names in Biological Informatics." *ZooKeys* 550 (July): 207–23. <https://doi.org/10.3897/zookeys.550.9546>.
- Roskov Y., Orrell T., Abucay L. 2018. "Species 2000 & ITIS Catalogue of Life, 2018 Annual Checklist." Leiden, the Netherlands.: Species 2000: Naturalis. www.catalogueoflife.org/annual-checklist/2018.
- Wagner, Viktoria. 2016. "A Review of Software Tools for Spell-Checking Taxon Names in Vegetation Databases." *Journal of Vegetation Science* 27 (6): 1323–7. <https://doi.org/10.1111/jvs.12432>.
- Wieczorek, John, David Bloom, Robert Guralnick, Stan Blum, Markus Döring, Renato Giovanni, Tim Robertson, and David Vieglais. 2012. "Darwin Core: An Evolving Community-Developed Biodiversity Data Standard." *PLoS ONE* 7 (1). <https://doi.org/10.1371/journal.pone.0029715>.
- Wilman, Hamish, Jonathan Belmaker, Jennifer Simpson, Carolina De La Rosa, Marcelo M. Rivadeneira, and Walter Jetz. 2016. "EltonTraits 1.0: Species-Level Foraging Attributes of the World's Birds and Mammals." Figshare. <https://doi.org/10.6084/m9.figshare.c.3306933.v1>.