

1 Harmonizing taxon names in 2 biodiversity data: a review of tools, 3 databases, and best practices

4 Matthias Grenié^{1,2*}, Emilio Berti^{1,3}, Juan Carvajal-Quintero^{1,2}, Gala Mona Louise Dädlow^{1,2},
5 Alban Sagouis^{1,4}, Marten Winter^{1,2}

6 ¹ German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig,
7 Puschstraße 4, 04103 Leipzig, Germany

8 ² Leipzig University, Ritterstraße 26, 04109 Leipzig, Germany

9 ³ Friedrich-Schiller University Jena, Jena, Germany

10 ⁴ Department of Computer Science, Martin Luther University, Halle-Wittenberg, Halle,
11 Germany

12 * corresponding author <matthias.grenie@idiv.de>

13 **Keywords:** taxonomic tools, taxonomic harmonization, taxonomic name matching,
14 taxonomy, R packages, taxonomic databases

15 Abstract

16 1. The process of standardizing taxon names, taxonomic name harmonization, is
17 necessary to properly merge data indexed by taxon names. The large variety of

18 taxonomic databases and related tools are often not well described. It is often unclear
19 which databases are actively maintained or what is the original source of taxonomic
20 information. In addition, software to access these databases is developed following
21 non-compatible standards, which creates additional challenges for users. As a result,
22 taxonomic harmonization has become a major obstacle in ecological studies that
23 seek to combine multiple datasets.

24 2. Here, we review and categorize a set of major taxonomic databases publicly
25 available as well as a large collection of R packages to access them and to
26 harmonize lists of taxon names. We categorized available taxonomic databases
27 according to their taxonomic breadth (e.g. taxon-specific vs multi-taxa) and spatial
28 scope (e.g. regional vs global), highlighting strengths and caveats of each type of
29 database. We divided R packages according to their function, (e.g. syntax
30 standardization tools, access to online databases, etc.) and highlighted overlaps
31 among them. We present our findings (e.g. network of linkages, data and tool
32 characteristics) in a ready-to-use Shiny web application (available at:
33 <https://mgrenie.shinyapps.io/taxharmonizexplorer/>).

34 3. We also provide general guidelines and best practice principles for taxonomic name
35 harmonization. As an illustrative example, we harmonized taxon names of one of the
36 largest databases of community time series currently available. We showed how
37 different workflows can be used for different goals, highlighting their strengths and
38 weaknesses and providing practical solutions to avoid common pitfalls.

39 4. To our knowledge, our opinionated review represents the most exhaustive evaluation
40 of links among and of taxonomic databases and related R tools. Finally, based on our
41 new insights in the field, we make recommendations for users, database managers,
42 and package developers alike.

43 Introduction

44 In the era of big data, combining, harmonizing, and analyzing massive amounts of ecological
45 data has played a central role in improving our understanding of biodiversity in a changing
46 world (Hampton et al., 2013; La Salle et al., 2016; Michener & Jones, 2012; Wüest et al.,
47 2020). While promising, this new era is also challenging. As exabytes of primary biodiversity
48 data become publicly available, issues of quality control in data integration, interoperability
49 and redundancy have become pressing concerns to address (Jin & Yang, 2020; Kissling et
50 al., 2018; Lenters et al., 2021; Nelson & Ellis, 2019; Soberón & Peterson, 2004; Thomas,
51 2009; Wüest et al., 2020).

52 One of the biggest challenges in biodiversity data handling is maintaining a consistent
53 taxonomy of species names associated with different biological attributes (Jin & Yang, 2020;
54 Meyer et al., 2016; Tessarolo et al., 2017; Thomas, 2009). The dynamic nature of taxonomy,
55 reinforced by the growing availability of information and the increasing use of genetic
56 methods to identify species results in ever-changing taxon names considered accepted.
57 Taxonomists start by sampling individuals in the field and when considered as not yet
58 described, name them, based on best knowledge and defined procedures (Dayrat, 2005).
59 These names become *de facto* accepted. However, some names can become obsolete,
60 when e.g. researchers realize later on this species was named already before. Those names
61 then are used as synonyms of another now accepted name (Lepage et al., 2014). In addition
62 to the names *per se*, taxonomists refer to species through taxonomic concepts—i.e. biological
63 entities—(Lepage et al., 2014). Which taxonomic concepts researchers use, i.e. are defined
64 as legitimate and valid, can vary across research cultures (Lepage et al., 2014). For some
65 taxonomic groups general consensus on one taxonomic concept is far from being reached
66 (Chawuthai et al., 2016), generating confusion. This dynamic process results in difficulties

67 for end-users to point to single valid names referring unambiguously to single taxonomic
68 concepts. The use of taxonomic databases helps resolve the different relationships that exist
69 between names and taxonomic concepts (one-to-one, one-to-many, many-to-one, or even
70 many-to-many, see Lepage et al., 2014).

71 In an attempt to unify taxonomy across the tree of life, multiple initiatives have proposed
72 curated lists of taxon names referenced against accepted taxon names. Taxonomic
73 databases ([Box 1](#)) are usually based on extensive community and individual expert
74 knowledge. Decisions which taxon names are accepted are usually based on robust
75 scientific evidence. These decisions might also have to be based on less objective reasons,
76 like reliability of original resources in comparison to conflicting studies or on individual
77 preferences for grammar and spelling (e.g. Isoëtes vs Isoetes) (Isaac et al., 2004).

78 However, despite significant efforts in creating a single authoritative list of the world's taxa
79 (e.g., [37]), taxonomic unification has largely advanced through multiple independent efforts
80 with different aims and scopes (e.g., per taxon group or region; Costello, 2020; Garnett et
81 al., 2020). For example, some taxonomic databases, i.e. databases that primarily offer
82 reference taxonomic data, focus on specific taxonomic groups (e.g. Freiberg et al., 2020),
83 others on environmental realms (e.g. [34]), providing a reference at either global or regional
84 scale such as national databases ([Figure 1](#)). The last decade brought a lot of progress in
85 taxonomy in general to overcome the “taxonomic impediment” (Rouhan & Gaudeul, 2021),
86 the lack of comprehensive information per taxonomic group. These efforts have generated a
87 large number of taxa lists with taxonomic-curated information dispersed across very different
88 repositories (König et al., 2019). For example, we are aware of four global taxonomic
89 databases focusing on plants (Leipzig Catalogue of Vascular Plants [22]; World Flora Online
90 [30]; Plants of the World Online [23]; World Plants, Hassler, 2021). While we know that
91 different databases provide different scientific opinions on taxonomy (i.e. using different

92 taxonomic concepts), meaning that they all contribute to the scientific debate and none of
93 them is right or wrong, how should the non-taxonomy expert end user (e.g. macroecologists)
94 know which resource is most suitable for her/his purposes? Researchers in need of
95 validating taxon names are confronted with many different taxonomic databases that have
96 often overlapping spatial or taxonomic coverage without a clear way to select which
97 database to use.

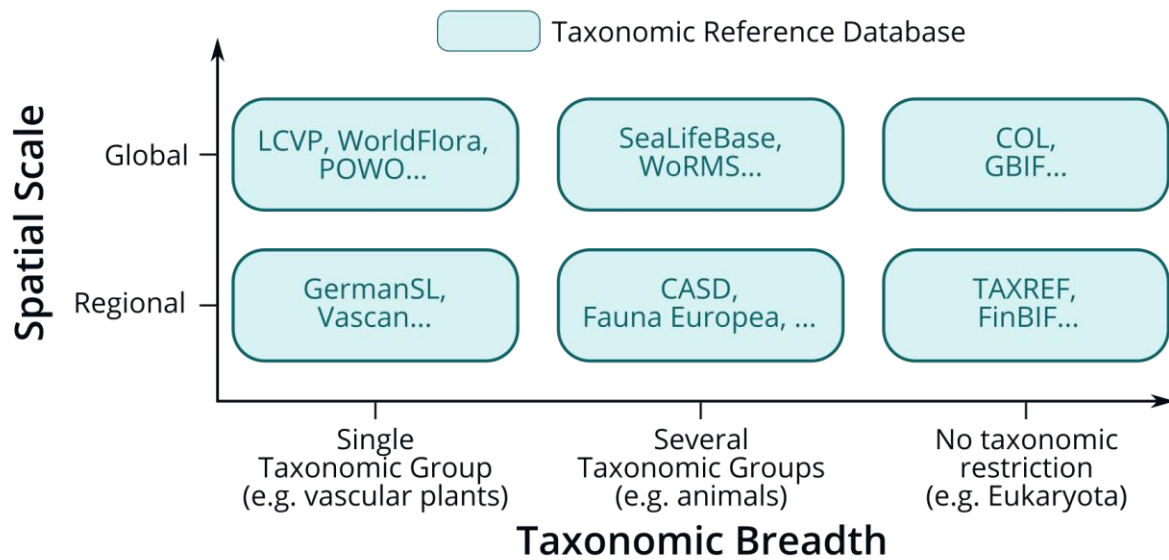
98 Taxonomic information, through taxon names ([Figure 2](#)), can serve as a common basis to
99 index and merge different biodiversity data (e.g., Dyer et al., 2017; occurrences: GBIF: The
100 Global Biodiversity Information Facility, 2020; conservation status: IUCN, 2021; traits: K. E.
101 Jones et al., 2009; Kattge et al., 2020; phylogenetic relationships: Smith & Brown, 2018;
102 Upham et al., 2019; invasion status: van Kleunen et al., 2019). Aside from the challenges
103 with maintaining updated and comprehensive taxonomic databases by themselves,
104 combining and harmonizing additional biological data can be problematic since such
105 datasets may have been created and updated at different times (sometimes spanning
106 several decades), may use different taxonomic databases to standardize taxon names, and
107 may not even be linked to any consistent taxonomic concept (Edwards et al., 2000; Farley et
108 al., 2018; König et al., 2019). Ultimately, if taxonomic name harmonization is not properly
109 executed, researchers are likely to introduce and propagate errors that can lead to
110 misquantified biodiversity components or mismatched data (Bortolus, 2008). Larger amounts
111 of data increase the issue, due to taxonomic inaccuracies introduced for increasing numbers
112 of species and taxonomic breadth (D. J. Patterson et al., 2010).

113 Driven by the needs in data harmonization, multiple tools have emerged for this task. This
114 has generated a diverse toolbox but no clear guidance on how these tools could be
115 combined into a meaningful and efficient workflow. Improving our knowledge of the
116 landscape of available taxonomic reference and tools is thus critical to developing robust

117 and comprehensive workflows to achieve high levels of data quality and accurate
118 downstream analyses.

119 Here, we fill this gap by reviewing publicly available taxonomic databases and R packages
120 for taxonomic harmonization, describing common pitfalls to avoid when using them, and
121 proposing hands-on approaches to achieve accurate and precise harmonized list of taxon
122 names. To our knowledge, our study represents the most comprehensive review and
123 assessment of tools and issues related to taxonomic name harmonization. We present and
124 discuss main steps towards robust and meaningful harmonization workflows. Specifically, we
125 review taxonomic databases, R packages, and show how they depend on and interact with
126 each other. We focus on R as it is the programming language of choice for ecologists (Lai et
127 al., 2019). We present a Shiny R application that guides users through the labyrinth of tools
128 and resources. We assess the efficiency of different possible taxonomic harmonization
129 workflows through a concrete use-case. We then formulate recommendations for end users,
130 tool developers and taxonomic data managers.

131



132

133 **Figure 1. Typology of taxonomic databases according to their taxonomic breadth and**

134 **their spatial scale.** The x-axis represents increasing taxonomic breadth from a single

135 taxonomic group to no clear taxonomic restriction (for example considering all Biota or all

136 Eukaryota). The y-axis represents spatial scale from regional to global. Each box represents

137 a specific type of taxonomic database, with examples. LCVP: Leipzig Catalogue of Vascular

138 Plants; WorldFlora: World Flora Online; POWO: Plants of the World Online; GermanSL:

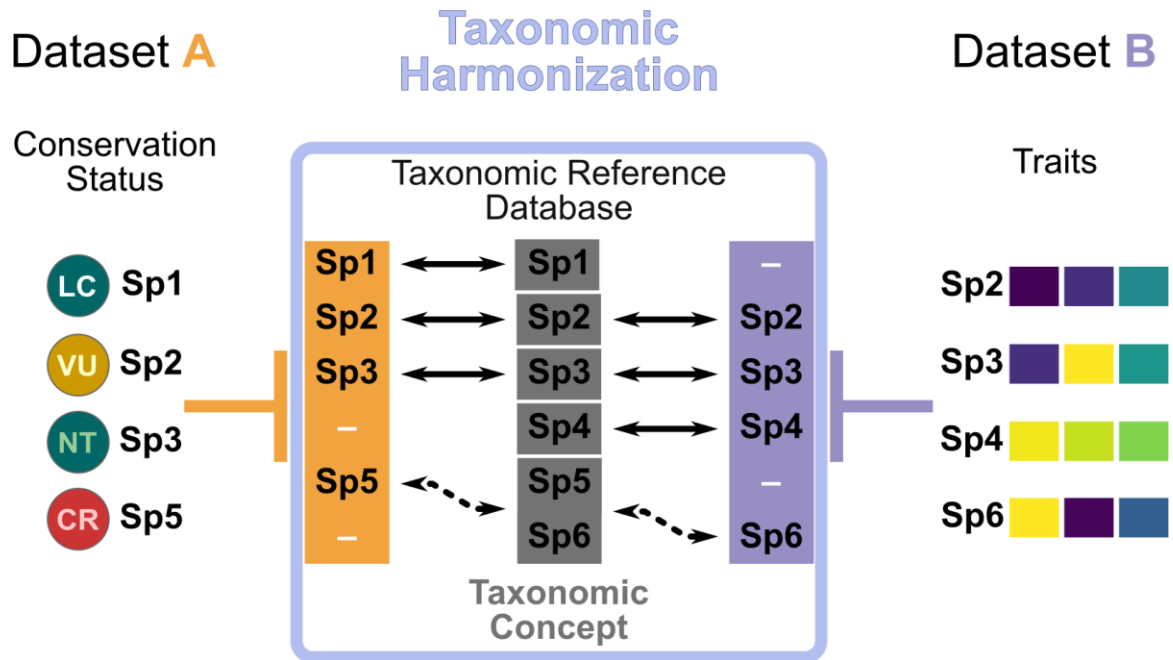
139 German Simple List; Vascan: Database of Vascular Plants of Canada; WoRMS: World

140 Register of Marine Species; CASD: Chinese Animal Scientific Database; COL: Catalogue of

141 Life; GBIF: Global Biodiversity Information Facility; TAXREF: French Taxonomic Referential;

142 FinBIF: Finnish Biodiversity Information Facility.

143



144

145 **Figure 2. Taxonomy as a unifying key for ecological datasets.** The two sides represent
 146 two exemplary datasets, with A containing conservation status of taxa (here species) and B
 147 their traits (colors show different traits). The datasets are indexed by taxon names “Sp1” to
 148 “Sp6”. The rounded rectangle in the middle depicts the taxonomic harmonization process: 1)
 149 The names are extracted from each dataset, respectively in the orange and purple
 150 rectangles; 2) Both lists are then compared to a taxonomic database which harmonizes all
 151 names. Here the names “Sp1” and “Sp6” refer to the same taxon in the taxonomic database
 152 (as indicated by the dashed lines). Without taxonomic harmonization, the exact match of
 153 names would have resulted in the loss of Sp5 and Sp6 when merging both datasets.

154

155

156 **Box 1. The taxonomic terminology diversity**

157 Across the literature, the terms **taxonomic reference (list)** (e.g. Freiberg et al., 2020),
158 **taxonomic authority (list/file)** (Vanden Berghe et al., 2015), **taxonomic databases** (Rees,
159 2014), **taxonomic backbone** (e.g. Schulman et al., 2021), or **taxonomic checklist**
160 (Costello, 2020) are used interchangeably, often without clear definitions. The terminological
161 diversity makes it difficult to understand differences between terms and potentially to find the
162 correct resources. For example, the expression “**taxonomic authority**” can be confused
163 with the authority when citing a species name, which is the citation of the author name
164 associated with a taxon. Different expressions can sometimes reflect differences in sizes of
165 provided databases, from a simple species **list** (e.g. to define the list of species names that
166 occur in a given area), to a full nomenclatural **reference** (with a taxonomy), to systems that
167 also provide synonymy resolution.

168 In this article, we use “**taxonomic databases**” as a generic expression of digital collections
169 of taxonomic information on many individual species, with processes to mitigate potential
170 conflicts between taxonomic designation.

171

172

173 The wild world of taxonomic resources

174 A typology of taxonomic databases

175 We categorized taxonomic databases (see Box 1) along two axes: taxonomic breadth and
176 covered spatial scale ([Figure 1](#)). Taxonomic breadth describes the amount of taxonomic
177 groups covered by the database. We use the term “taxonomic group” as a broad term to
178 describe a group of taxa or taxonomic ranks at which people work (e.g. birds - class *Aves*,
179 butterflies - order *Lepidoptera*). Databases have varying taxonomic and spatial breadths,
180 from narrow taxonomic breadth but global scale (e.g. eBird [17]) to broad taxonomic breadth
181 but regional/national scope (e.g. the Chinese Animal Species Database [4]). Some
182 databases even aim to provide information without any taxonomic restriction at a global
183 level, e.g. Catalogue of Life [37].

184 Because navigating the landscape of taxonomic databases can be difficult for users, we
185 provide a wide overview of available databases on as many taxonomic groups as possible at
186 varying spatial scales and taxonomic breadths ([Table 1](#)). As one covering many databases,
187 this list provides an entry point for users to get a sense of potential sources of taxonomy.
188 The immense variety of taxonomic databases, especially at regional scales, prevents our list
189 from being exhaustive but it includes most existing global databases.

190 **Table 1. A list of taxonomic databases.** We included all databases accessed by the tools
191 we referenced in the next section. Square brackets indicate [supplementary references](#).

Taxonomic Breadth → --- Spatial Scale ↓	Narrow = Single taxonomic group	Medium = Several taxonomic groups	Wide = No taxonomic restriction
---	---	---	---

<p>Regional</p>	<p>Vascular Plants GermanSL (https://germansl.infinite-nature.org/) [1], USDA (https://plants.usda.gov/home) [2], Vascan (https://data.canadensys.net/ipt/resource?r=vascan) [3]</p>	<p>Animals CASD (http://zoology.especies.cn/) [4], All plants and fungus FB2020 (http://floradobrasil.jbrj.gov.br/) [5]</p>	<p>No taxonomic restriction Dyntaxa (https://www.dyntaxa.se/) [6], EUBON (http://biodiversity.eubon.eu/web/guest/eu-bon-taxonomic-backbone) [7], FinBIF (https://laji.fi/en/) [8], NBN (https://nbn.org.uk/) [9], PESI (https://www.eu-nomen.eu/portal/index.php) [10], SP2000CN (http://sp2000.org.cn/) [11], TaiCOL (https://taibnet.sinica.edu.tw/eng) [12], TAXREF (https://inpn.mnhn.fr/programme/documentation/referentiels-especes-taxref) [13], TWN (https://twnlist.aquadesk.nl/) [14]</p>
<p>Global</p>	<p>Algae AlgaeBase (https://www.algaebase.org/) [15] Amphibians ASW (https://amphibiansoftheworld.amnh.org/) [16], Birds eBird/Clements (https://ebird.org/science/use-ebird-data/the-ebird-taxonomy) [17] Fungi Index Fungorum (http://www.indexfungorum.org/) [18] Fish FishBase</p>	<p>Marine organisms SeaLifeBase (https://sealifebase.ca/home/index.php) [33], WoRMS (https://www.marinespecies.org/) [34], Animals ZooBank † (http://zoobank.org/) [35]</p>	<p>No taxonomic restriction BOLD (http://www.barcodinglife.org/) [36], COL (https://www.catalogueoflife.org/) [37], EOL (https://eol.org/) [38], GBIF (https://www.gbif.org/) [39], GNI□□* (https://index.globalnames.org/) [40], GNR* (https://resolver.globalnames.org/) [41], GNV (https://verifier.globalnames.org/) [41]</p>

	<p>(https://www.fishbase.in) [19]</p> <p>Mammals</p> <p>MMD (https://www.mammaldiversity.org/) [20],</p> <p>Plants</p> <p>IPNI † (https://www.ipni.org/) [21], LCVP [22], POWO (http://powo.science.ke.w.org/) [23], TPL* [24], TNRS (http://tnrs.iplantcollaborative.org/) [25], [26], Tropicos (https://tropicos.org/) [27], WCSP (https://wcsp.science.ke.w.org/) [28], WCVF (https://wcvf.science.ke.w.org/) [29], World Flora Online (http://worldfloraonline.org/) [30],</p> <p>Reptiles</p> <p>ReptileDB (https://www.reptile-database.org) [31],</p> <p>Spiders</p> <p>WSC (https://wsc.nmbe.ch/) [32]</p>		<p>es.org/) [42],</p> <p>ION (http://www.organismnames.com/) [43],</p> <p>ITIS (https://www.itis.gov/) [44],</p> <p>IUCN (https://www.iucnredlist.org/) [45],</p> <p>NatServe (https://explorer.natureserve.org/) [46],</p> <p>NCBI (https://www.ncbi.nlm.nih.gov/taxonomy) [47],</p> <p>Neotoma (http://neotomadb.org/) [48],</p> <p>OTL (https://opentreeoflife.github.io/) [49],</p> <p>PBDB (https://paleobiodb.org/#/) [50],</p> <p>Wikidata (https://www.wikidata.org/)</p> <p>,</p> <p>Wikipedia (https://www.wikipedia.org/),</p> <p>Wikispecies (https://species.wikimedia.org/)</p>
--	---	--	--

192 * Databases that can be considered as outdated.

193 † Rather a nomenclatural database (collection of names) than a taxonomic reference

194 ¹ The Plant List (<https://www.theplantlist.org/>), while still widely used and easy to access, has not

195 been updated since the release of its version 1.1 in September 2013. It has been superseded notably

196 by World Flora Online and other initiatives such as POWO and LCVP.

197

198 **The wide landscape of R packages for taxonomy**

199 With the increasing amount of data used in ecological studies, taxonomic harmonization
 200 cannot rely on manual curation. Computational tools are needed to help extract, evaluate,
 201 manipulate, and visualize taxonomic information. Additionally, the use of computational tools
 202 increases the reproducibility of analyses compared to manual edits. In this section, we
 203 present the most extensive review, to our knowledge, of R packages that can be used to
 204 process taxonomic information ([Table 2](#)).

205 **Description of the landscape of tools**

206 **Table 2. Identified R packages useful for taxonomic name harmonization.** Square
 207 brackets indicate [supplementary references](#).

Category Name	Packages
Infrastructure	taxa [51], taxlist [52], taxview [53]
Database Access (Online)	algaeClassify [54], AmphiNom [55], arakno [56], dyntaxa [6], finbif [57], kewr [58], natserv [59], ncbi [60], neotoma2 [61], paleobioDB [62], plantlist [63], rcol [64], rebird [65], rentrez [66], rfishbase [67], rgbif [68], ritis [69], Rocc [70], rotl [71], rredlist [72], rreptiledb [73], rtaxref [74], SP2000 [75], taxize [76], [77], taxonomyCleanr [78], Taxonstand [79], [80], taxotools [81], [82], taxreturn [83], TNRS [84], [85], tw [86], wikipata [87], worms [88], worrms [89], zbank [90]
Database Access (Offline)	AmphiNom [55], flattax [91], flora [92], lcvplants [22], mammals [93], ncbi [60], taxadb [94], [95], taxalight [96], taxastand [97], taxizedb [98], taxonlookup [99], taxonomizr [100], tpl [101], vegdata [102], WorldFlora [103]
Data Wrangling	metacoder [104], monographaR [105], rgnparser [106], splister [107], taxastand [97], taxreturn [83], taxspell [108], traitdataform [109], vegdata [102], vegetable [110], yatah [111]

208

209 We identified some packages that provide standardized technical infrastructure for
210 taxonomic experts to develop and work with taxonomic information within R. Infrastructure
211 packages provide basic “building blocks” for other packages to build onto. **taxa** [51], used by
212 **metacoder** [104], provides R-native objects and methods to represent taxonomic data.
213 **taxlist** [52] contains objects and functions to store taxa lists, synonyms, taxonomic
214 hierarchy, and functional traits in a standardized format; it is used by **vegdata** [102].
215 **taxview** [53] provides basic visualization of taxonomic hierarchies; it is used by no other
216 packages. The fact that virtually no other packages rely on them means that several tools
217 reinvent the wheel instead of relying on standardized functions. More widespread reliance on
218 infrastructure packages and associated methods within the small community of R taxonomy
219 package developers could foster best development practices, easier interoperability, as well
220 as increased reproducibility, as it has been for example done already for spatial data through
221 the **sp** and **sf** packages (Bivand et al., 2013; Pebesma, 2018; Pebesma & Bivand, 2005).

222 We identified 47 packages providing direct access to online taxonomic databases. These
223 packages let the users search a given taxon name in one (or several) online taxonomic
224 database(s) and get back a list of potential matching names, considering both accepted
225 names and synonyms. Details about the packages, e.g. which taxonomic databases they
226 access are available in [S2](#) and our specifically for this review developed shiny app
227 **taxharmonizexplorer**. You can explore which package(s) access which database(s) as
228 well as additional useful characteristics through **taxharmonizexplorer** described in the
229 following section.

230 Accessing online databases does not come free of issues: (1) Online databases can be
231 updated continuously, potentially leading to different versions used when harmonizing at
232 different times or on- and offline, hindering reproducibility. (2) Database access is not always
233 guaranteed because of technical issues with online resources (maintenance needed, server
234 outage, Internet accessibility). (3) Some databases implement a form of request limitation,
235 enforcing a maximum number of queries that can be made in a given period of time (e.g. one
236 query every three seconds), with one query matching a single species only. (4) Online query
237 execution speed can be limited compared to local queries (of the order of several seconds
238 against tens of milliseconds, see [94, 95]) and potentially impossible if the Internet
239 connection is unstable. (5) Databases also limit the complexity of queries with no standard
240 format across databases, e.g. the user can only get a list of accepted names from an input
241 name and not ask more precise questions like “What are all names with epithet *alba*?”.

242 To overcome these issues several packages provide or build local database copies.
243 **lcvplants** [22] accesses the LCVP database fully offline through a local copy, it also offers
244 functions to harmonize two lists of names. **ncbit** [60] provides a similar access but to the
245 NCBI database [47]. **taxadb** [94, 95] creates a unified local database from different data
246 sources as specified by the user. **taxalight** [96], which is maintained by the same
247 developers, is faster and with fewer dependencies, it will supersede **taxadb** (Boettiger,
248 personal communication). **taxizedb** [98] also downloads local copies of the database but,
249 contrary to **taxadb** and **taxalight**, it provides the data without standardizing its format
250 between sources. The user can then access the original information through SQL queries
251 tailored for each database. **taxonlookup** [99] provides a curated versioned taxonomy of
252 land plants. **taxastand** [97] lets the user load local taxonomic data in Darwin Core format.
253 **vegdata** [102] allows the download of the GermanSL database to access it offline. It also
254 offers access to any (offline) TurboVeg database available on the user’s computer within R.

255 **worldFlora** [103] lets the user access the World Flora Online database from R once it has
256 been downloaded by the user.

257 Taxonomic harmonization is not limited to accessing databases and accessing lists of
258 (un)accepted names. Several R packages offer functions to manipulate taxonomic data,
259 parse taxonomic files, or summarize taxonomic information. **monographaR** [105] uses
260 standardized tables to produce a monograph on examined specimens in a paper, with
261 associated maps and phenological diagrams. **rgnparser** [106] wraps within an R package a
262 tool built by GlobalNames in the Go language that parses scientific names into components
263 (i.e., genus, species, authority, year, etc.) efficiently. **taxlist** [52] and **vegdata** [102]
264 provide help functions to harmonize one's own taxa list, including interaction with TurboVeg.
265 **taxonomyCleanr** [78] processes and cleans taxonomic information, including a function to
266 write taxonomy in Ecological Metadata Language (EML; M. B. Jones et al., 2006).
267 **taxotools** [81, 82] contains functions to create your own taxonomic database and match it
268 with other lists, it also parses data in Darwin Core format. **yatah** [111] parses taxonomic
269 information from long strings with special characters as used in genomic data, outputs
270 summary statistics about it, and visualizes associated taxonomic hierarchy.

271 We identified several packages that deal with taxonomic assignment from genomic data but
272 considered them out of scope of this review (see [S1](#) for the inclusion criteria).

273 Tools: Lessons Learned and Future Direction

274 To avoid reinventing the wheel, whenever possible, package developers should build their
275 tools on top of existing packages and functions; however, we found little evidence for
276 package or function reuse across packages (see lack of network links in
277 [taxharmonizexplorer](#)). As an exception, **taxize** [76, 77] relies on individual packages that

278 provide functions to access specific online databases (for example it relies on **rfishbase**
279 [67] to access FishBase). The lack of dependencies between packages is inefficient from a
280 developer standpoint and unclear for end users, due to packages performing virtually
281 identical tasks but in a slightly different way, with different syntaxes, and different ways of
282 handling errors. For example, **plantlist** [63], **taxadb** [94, 95], **taxalight** [96], **taxize** [76,
283 77], **taxizedb** [98], **Taxonstand** [79, 80], and **tpl** [101] all access The Plant List data. While
284 evaluating relevant tools, we identified several packages in early development. **splister**
285 [107] and **taxastand** [97] both allow the user to match its own custom reference database,
286 which can be useful for areas or taxa where no commonly accepted taxonomy exists.
287 **taxreturn** [83] fetches data from BOLD and NCBI taxonomies for metabarcoding.
288 **taxspell** [108] checks the spelling of taxon names through dictionaries that reference the
289 most common spelling mistakes.

290 Our review was facilitated by the fact that the packages are deposited in standardized
291 central repositories such as CRAN or Bioconductor. Many packages were also accessible in
292 their last development state on open development platforms such as GitHub. Thanks to this
293 accessibility, we identified the tools in development mentioned in the paragraph above,
294 showing the trends in tools for taxonomy.

295 Of the 60 packages we included, 20 were made available through rOpenSci, many of which
296 are central in global taxonomic harmonization such as **taxize** [76, 77]. rOpenSci is a not-
297 for-profit organization that aims to “[...] *help develop R packages for the sciences via*
298 *community driven learning, review and maintenance of contributed software in the R*
299 *ecosystem*” (Boettiger et al., 2015). The fact that rOpenSci supported the development and
300 the publicity of many tools important for taxonomy underlines how rOpenSci filled quasi an
301 “ecological” package niche that was not filled by traditional scientific developers. Resolving

302 taxonomic name conflicts requires good taxonomic knowledge, which is rare outside of
303 taxonomists. While the manipulation of online databases requires a good knowledge of web
304 technologies, uncommon among scientists. The intersection of both is thus even rarer.
305 Furthermore, there are few incentives to build and maintain scientific software (Jay et al.,
306 2020; Mislán et al., 2016). The combined expertise found among rOpenSci members greatly
307 helped advance the development and maintenance of tools to interact with taxonomic data.

308 Several tools we reviewed accessed data that can be considered outdated. For example,
309 several packages access The Plant List [24], which used to be the main global taxonomic
310 database for plants, but has not been updated since 2013 and is considered outdated by its
311 authority (see <https://www.theplantlist.org/>). It refers now to the World Flora Online database
312 as the updated successor [30]. Despite this, because of its easy access, standardized
313 format, and continuous availability it is still used by packages created long after 2013. The
314 Plant List has gained ~1000 citations, since 2020, (according to Google Scholar) of which
315 very likely many used the outdated list, leading to results based on outdated knowledge.

316 Similarly, **taxize** [76, 77] accesses both Global Names Index and Global Names Resolver,
317 which are massive collections of other taxonomic databases (Mozzherin et al., 2021). Global
318 Names Index has not been updated since 2018 and it has been superseded by Global
319 Names Resolver in 2018 (Mozzherin, personal communication). Global Names Resolver has
320 in turn been superseded by Global Names Verifier (Mozzherin, 2021), with even faster
321 software and continuously updated data. While maintaining access to older databases is
322 paramount to ensure the reproducibility of taxonomic name harmonization, users should
323 check the date of last update of the resource they are accessing. The tools should explicitly
324 warn their users when they are using outdated taxonomic databases and point them to
325 alternative, more up-to-date, sources.

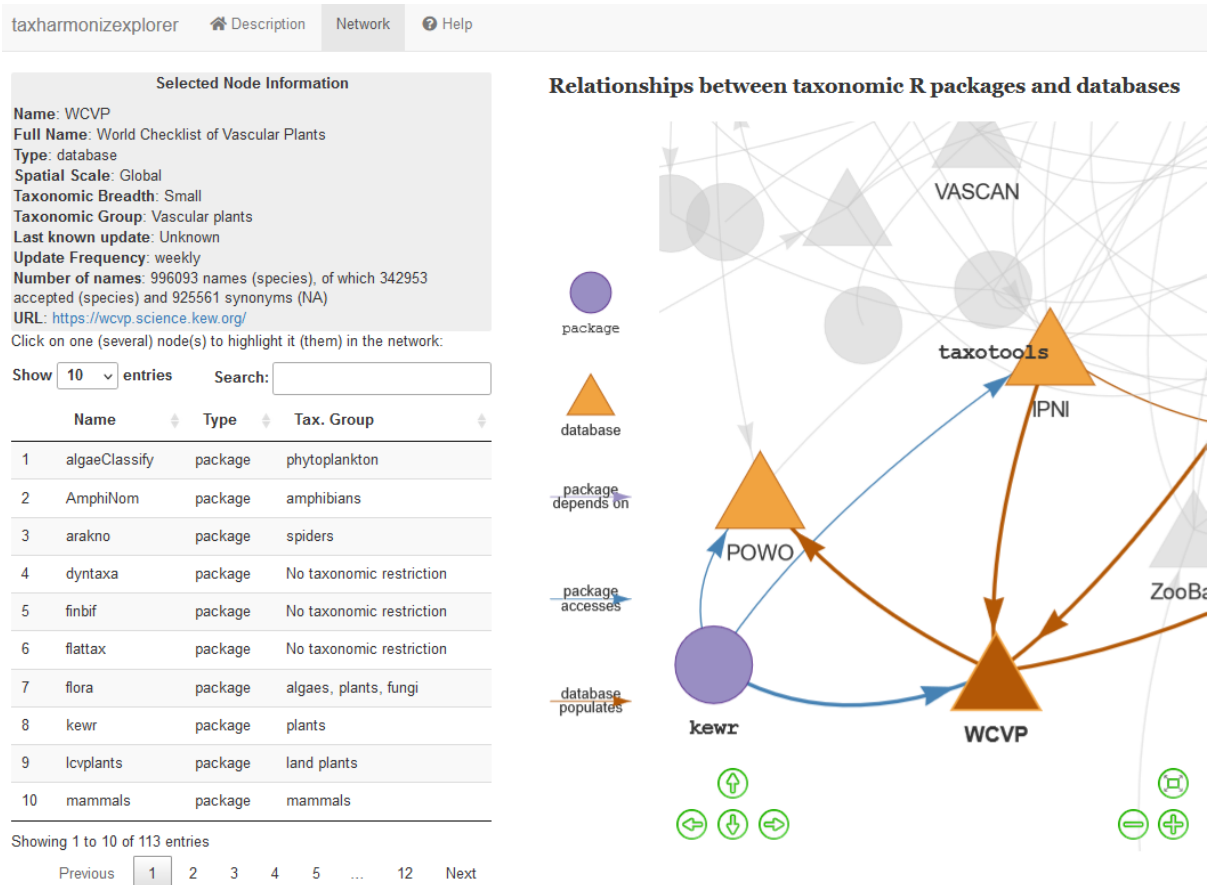
326 A tool to guide users in the network of resources

327 To help the users navigate the complex network of tools and databases, we developed a
328 shiny application that lets users explore the relationships between resources and their main
329 characteristics (date of last update, taxonomic breadth, URL, etc.). We called it
330 **taxharmonizexplorer** and it's available as a perennial archive on Zenodo (Grenié et al.,
331 2021) but also accessible online at: <https://mgrenie.shinyapps.io/taxharmonizexplorer>.

332 The application presents on the right side a network that links taxonomic databases and
333 packages ([Figure 3](#)). Global databases with a wide taxonomic breath often aggregate
334 taxonomies trying to provide a unified taxonomic backbone for all covered organisms, such
335 as Catalogue Of Life (COL) or Encyclopedia Of Life (EOL) [37, 38]. The databases are
336 connected when they rely on one another, while packages are connected when they depend
337 on each other. Finally, packages are connected to databases when they provide access to
338 the databases. The top left panel displays information about the node selected on the
339 network and includes a link to the package or database website. The bottom left of the app
340 shows a table where the user can select and search for nodes through their name, type, and
341 taxonomic group.

342 The dataset that backs the network is continuously improving as we are identifying the links
343 that connect the different databases and add new R packages. The dataset is open for
344 contributions for packages and databases that we may have missed (through GitHub or
345 email to the corresponding author).

346



347

348 **Figure 3. Screenshot showing the network view of taxharmonizexplorer.** The left
 349 section shows a table of each of the nodes in the network to let the user select manually
 350 nodes of interest, the top part presents a summary of the information on the selected node in
 351 the network. The right section displays the relationships between packages (which depends
 352 on which other), between databases (how one populates another one), and between
 353 packages and databases (which packages access which databases).

354

356 Box 2. The double-edged sword of “fuzzy matching”

357 “Fuzzy matching” is a method to match taxon names that differ by some characters.

358 **How it works.** Similarity measures are used to quantify the discrepancy between two names
359 (Meyer et al., 2016). For example, orthographic distance metrics measure similarity as the
360 reciprocal of the number of characters to be modified to obtain one string from another. The
361 obtained score indicates how close two names are to each other. The highest score name is
362 then matched to the name of interest. One common metric is measuring single-character
363 deletions, substitutions, or insertions with the Levenshtein Distance (e.g. [95]). An alternative
364 is the phonetic modified Damerau-Levenshtein distance weighting transpositions lower than
365 individual character substitutions (Taxamatch; Rees, 2014).

366 **When to use it.** Fuzzy matching is useful when orthographic and spelling errors are
367 suspected in the list of taxon names, meaning that exact matching cannot resolve them.
368 These typos can have multiple causes, e.g. transcription mistakes, wrong latin name,
369 differences in spelling style among taxonomic authorities, changes in the spelling style of
370 accepted names, etc.

371 **Risks.** When two different taxa display similar names (low orthographic distance), they can
372 be fuzzy matched to the same accepted name. If used blindly to match taxon names at
373 broad spatial and temporal scale and taxonomic coverage, there is a relatively high risk of
374 fuzzy matching a wrong name in a different part of the tree of life. The Interim Register of
375 Marine and Nonmarine Genera (Rees, 2021) provides a database of possible name colliders
376 at genus level.

377 Resort to fuzzy matching should only come at the end of the harmonization process to cast a

378 bigger net of candidate names. Use of fuzzy matching should always be explicitly stated by
379 users; tools that implement fuzzy matching by default should highlight this feature and give
380 the option to toggle it off. Tools should also mention to what extent are results based on
381 fuzzy matching. When resorting to fuzzy matching, sensitivity analyses should be performed
382 using fuzzy matching scores, e.g. by random sampling taxon names using matching scores
383 as probability weights.

384

385

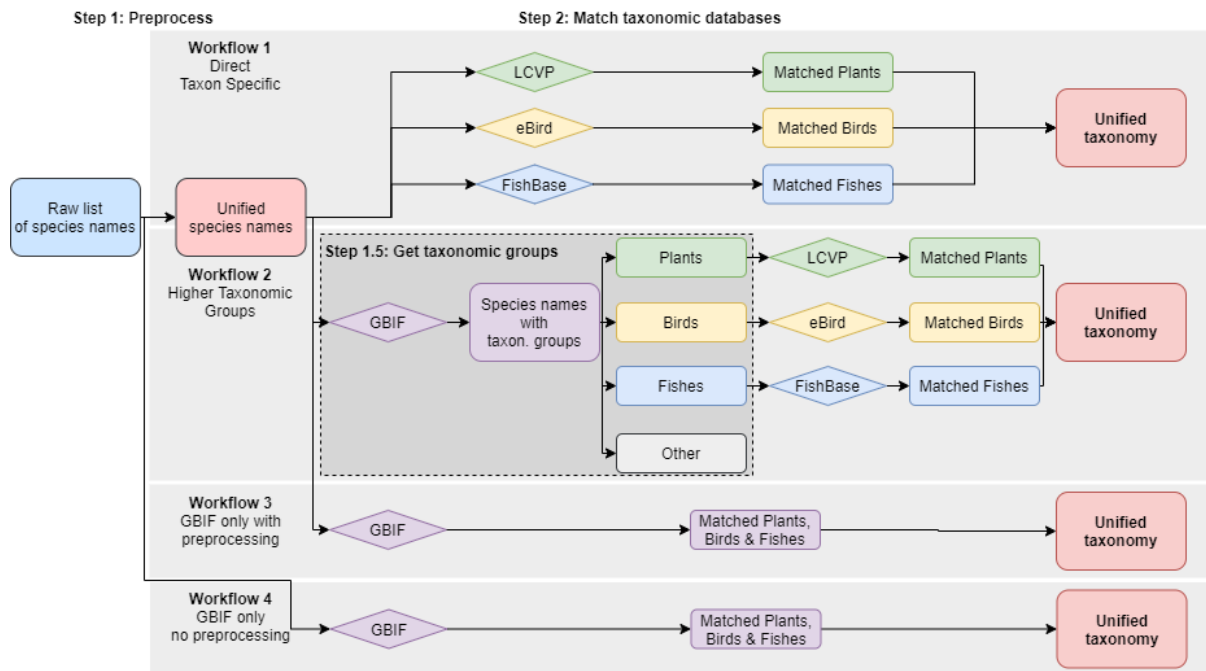
386

387 Stepping out of the taxonomic harmonization 388 labyrinth: recommendations and a comparison of 389 example workflows

390 In this section, we provide general guidelines and best practices to harmonize taxonomy in
391 large biodiversity datasets to avoid common pitfalls. As an illustrative example, we
392 harmonize taxon names from BioTIME (v. 02_04_2018, BioTIME Consortium, 2018;
393 Dornelas et al., 2018), the largest global compilation of time series assemblages, which
394 includes 44,440 taxa spanning multiple taxonomic groups at broad spatial and temporal
395 scales. BioTIME is often used (~145 citations) and is particularly interesting as it gathers
396 information from different data sources (361 studies), which potentially leads to taxonomic
397 inconsistencies between them. For the sake of simplicity we only focus here on birds, fishes,
398 and vascular plants in BioTIME. We detailed the process and tools used for our taxonomic
399 harmonization (packages, including versions, specific functions, and parameter values
400 used). To achieve full reproducibility we encourage others to detail their workflow in a similar
401 fashion, as taxonomic harmonization workflows can be highly sensitive to the exact version
402 of the tools or data used..

403 We applied four different workflows (WF, [Figure 4](#)), to harmonize the taxonomy of BioTIME.
404 WF1 and WF2 use taxon-specific databases whenever available. WF1 matches all species
405 names against all chosen taxon-specific databases and conflicts are resolved afterwards,
406 whereas in WF2 taxa are first assigned to higher taxonomic groups (birds, fish or vascular
407 plants) and only then matched against relevant taxon-specific databases. WF1 and WF2 can
408 be summarized as follows: Step 1, taxon names are preprocessed to unify writing style. Step

409 1.5 (only in WF2) taxa are assigned to high taxonomic groups using a multi-taxa global
410 database. Step 2, taxon names are matched against taxon-specific databases. The other
411 two workflows, WF3 and WF4, only use GBIF to harmonize all names. In WF3 names are
412 pre-processed (Step 1 as in WF1 and WF2), while in WF4 taxon names are passed directly
413 from BioTIME to GBIF. We included these two workflows because they are intuitive and
414 easy to implement and, as such, appeal particularly to non-taxonomists. We compared the
415 performance of the different workflows by the number of identified names in the different
416 taxonomic groups (birds, fishes, and vascular plants).



417

418 **Figure 4. Diagram of different taxonomic harmonization workflows.** The workflows differ
 419 in the number of steps they consider and the databases they leverage on. Rounded
 420 rectangles are lists of taxon names while diamonds represent taxonomic databases against
 421 which the names are matched. The different colors used at step 2 represent different
 422 taxonomic groups.

423 Step 1: Preprocess names (a.k.a. clean/unify writing style)

424 Taxon names writing style can vary between sources, complicating harmonization (D.
425 Patterson et al., 2016; D. J. Patterson et al., 2010) and becoming a source for errors. These
426 differences arise because of the disparate use of upper and lower case, abbreviations,
427 annotations, depictions of hybrids, authorships, etc. Removing these syntactic issues and
428 standardizing taxon names is thus the starting point of taxonomic harmonization. To match
429 all possible variations of a scientific name, these need to be divided into their stable (e.g.
430 genus, species epithet, and authorships) and prone-to-change elements (e.g. annotations)
431 and then combined into only stable elements (Mozzherin et al., 2017). The result is a
432 syntactically normalized list of names. We recommend keeping authorship, whenever
433 possible, along the taxon names because it decreases errors. Using taxa authorship
434 information also disambiguates between accepted and synonyms names (e.g., the IRMNG
435 referencing binomial homonyms, Rees, 2021).

436 To standardize the writing style of taxon names across BioTIME, we used the function
437 `gn_parse_tidy()` from package `rgnparser` v.0.2.0 [106]. After parsing taxon names, we
438 only kept the two first words of each parsed name, which ideally represent the scientific
439 binomial name of species (*Genus species*). We did not keep authorship as most names in
440 BioTIME did not have it. We applied this step for all workflows except WF4. We found that of
441 the 44,326 names reported in the original file, 4,734 taxa (11%) had spelling style
442 differences, i.e. species with the same binomial name after parsing. Of the remaining 39,592
443 unique taxon names, 6,692 were composed of only one word. We removed these taxa as
444 our aim was to match only binomial names. Importantly, the remaining 32,900 names also
445 contained common names and undetermined taxa with taxonomic abbreviation and
446 keywords, e.g. "*Family fam*". As our aim was to programmatically harmonize taxonomy using

447 available R packages, we kept such binomial entries as they were returned from `rgnparser`
448 [106]; such inaccuracies will be solved in the next steps. GBIF offers an alternative name
449 parser, which can be used through `rgbif` with the `parsenames()` function [68].

450 Step 1.5: (if needed) Divide taxa in higher taxonomic groups

451 In WF2, taxon names are passed only to the relevant taxon-specific databases, e.g. plants
452 are matched only against a plant-specific database. Multi-taxa global databases (e.g. GBIF
453 [39]) can provide classification to divide taxa into taxonomic groups. The potential errors
454 should be fairly limited for higher taxonomic groups as multi-taxa databases generally offer
455 reliable higher taxonomy (regna, phylum, class, etc.), even though some binomial names
456 could match across different phyla (e.g., the *Aotus* genus is present in both plants and
457 monkeys). These cases are referenced in the Interim Register of Marine and Non-Marine
458 Genera (IRMNG, Rees, 2021).

459 BioTIME originally assigns taxonomic groups, but these are at the study level rather than for
460 each species. For example, the species *Abalistes stellatus* was correctly assigned to the fish
461 group except in one study, where it was assigned to the benthos group (to which most of the
462 species in this study belong). To achieve maximal taxonomic accuracy, we reclassified
463 species names into higher taxonomic groups using GBIF. We queried all names against
464 GBIF and, based on higher clades (mostly taxonomic classes, e.g. *Sarcopterygii*, and
465 unranked clades, e.g. *Tracheophyta*), we grouped names into three groups that could be
466 referred to by taxon-specific databases: birds, fishes, and vascular plants.

467 Step 2: Match taxonomic databases

468 The selection of databases and packages for harmonization depends on the taxonomic
469 breadth and the spatial coverage of the species list under study ([Figure 1](#)). In general, we
470 recommend using the most updated and taxa-specific databases. For example, if this
471 contains species names for one taxonomic group (e.g. fishes) from a specific region (e.g.
472 France), the most appropriate approach should be to use a taxon-specific global database
473 (e.g. FishBase [19]) or a regional database (e.g. TAXREF [13]). For instance, if the aim is to
474 merge the list of species names with other global datasets, then FishBase would be
475 preferred, whereas if the goal is to provide a comprehensive list of species in France, then
476 TAXREF can be used instead. This approach can present some caveats in specific cases.
477 For example, if the regional studied dataset comprises non-native or aquatic species that
478 may not be present in the regional or terrestrial focused database respectively, but would
479 likely be present in a global database. Another example would be using fuzzy matching ([Box](#)
480 [2](#)) on a database of large taxonomic scope which could end up matching names in the
481 wrong part of the tree of life (e.g. *Fucus* to *Ficus*).

482 The type of search, exact matching vs. fuzzy matching (see details in [Box 2](#)), performed
483 during taxon name matching can strongly affect the results. While fuzzy matching can
484 correct misspellings, it increases the chances of mismatching errors. A way to safeguard
485 against potential mismatches is to perform a first harmonization without fuzzy matching and
486 then a second process (Step 3 below) including fuzzy matching algorithms only if many
487 species names are left without matches. The use of higher taxonomic ranks can also help
488 control that fuzzy matched names correspond to the appropriate part of the tree of life.

489 Finally, we strongly recommend tracking package versions and version or date of access of
490 the taxonomic database(s) used. Tracking versions increases replicability, as different

491 versions of packages and databases can give different results. For example, **taxadb** [94, 95]
492 uses yearly snapshots of taxonomic databases, provided by the developers, to create a local
493 database. On the other hand, **taxize** [76, 77] uses the last available version accessing
494 databases online APIs.

495 As BioTIME has global scope, we used only global databases. The choice of taxonomic
496 references and R packages to use was informed by our Shiny app, providing a direct
497 example of its utility. The databases and R packages used were: eBird v.2021 and **rebird**
498 v.1.2.0 for birds, FishBase v.21.04 and **rfishbase** v.3.1.8 for fishes, **lcvplants** v.1.1.1 and
499 **LCVP** v. 1.0.4 for plants, and GBIF (accessed August 2021) and **rgbif** v.3.6.0 for assigning
500 taxonomic groups in WF2 and for WF3 and WF4. We only used exact matching. Of the
501 32,900 parsed names, WF1 matched, as unique names, 878 birds, 5123 fishes, and 4435
502 plants ([Table 3](#)). WF2 matched slightly less ($n = 25$) species names, caused by
503 misclassification of higher taxonomic groups, mostly plants ($n = 23$), by GBIF (Step 1.5).
504 WF3 and WF4 matched the highest number of species, with 795 and 803 more species than
505 WF1, respectively. The higher number of species matched was, however, due for a large
506 proportion to species names that were considered synonyms in WF1 and WF2 and that were
507 thus assigned to the same accepted name by taxon-specific databases. For instance, 734
508 synonyms were identified in WF2, while there were only 484 in WF3. Because of this, WF3
509 and WF4 should be generally avoided when suitable taxon-specific databases are available.

510 In summary, the workflows using taxon-specific databases performed relatively similar in the
511 number of matched names, with WF1 matching slightly more species than WF2, but
512 requiring three times the queries needed for WF2. WF3 and WF4 were faster, easier, and
513 matched the most species names, but this was at the expense of not resolving many
514 synonyms. Which of these workflows is best depends ultimately on the goal of the taxonomic

515 harmonization process and users must choose what suits most the task at hand. Yet, using
516 taxon-specific databases (WF2) to match species names already divided into high taxonomic
517 groups seems an optimal trade-off between computational speed, programmatic complexity,
518 accuracy and robustness of the harmonization process.

519 Step 3: (do at your own risk) Resolve unmatched names with 520 fuzzy matching

521 If not satisfied with the number of matches achieved through Steps 1-2, further steps can be
522 implemented to maximize the number of matched names, looking for misspellings not
523 corrected in Steps 1-2. These spelling errors correspond to errors associated with the wrong
524 spelling of latin names (e.g. the use *Breviraja caerulia* instead *Breviraja caerulea*), either due
525 to typos or caused by using different databases (Costello et al., 2013; D. Patterson et al.,
526 2016; D. J. Patterson et al., 2010). Some misspellings may have been corrected during Step
527 2 if species names were matched using fuzzy matching.

528 To correct spelling errors, algorithms are available to calculate the probability of
529 correspondence between an input taxon name and long lists of names. Although these fuzzy
530 searches have some risks ([Box 2](#)), functions like `gnr_resolve()` from package `taxize`
531 have arguments that reduce the probability of mismatching. Its argument `with_context`
532 restricts the search to a narrower taxonomical context, reducing the probability of matching
533 homonyms from different taxonomic groups (Costello et al., 2013; Shipunov, 2011). The
534 IRMNG database, that references colliding genera names across the tree of life, can also be
535 used to check potential typos (Rees, 2021). As fuzzy algorithms programmatically match
536 names based on their orthographic similarity, often without considering additional taxonomic

537 information, extra care should be taken if step 3 is implemented, including sensitivity
538 analyses and manual checking of matched names.

539 We applied this step only to WF2. We looked for misspellings across the 777 names
540 belonging to birds, fishes, and plants (from Step 1.5) that were not matched in WF2. We
541 used the function `gnr_resolve()` from `taxize` v.0.9.99 and selected only the best matching
542 names. We thus corrected spelling errors for 293 names and matched an additional 218
543 unique species applying again Step 2: 22 of 267 bird names, 130 of 253 fish names, and 66
544 of 257 plant names. Despite the improvement in the number of matches, these may be
545 wrong due to fuzzy matching and orthographic corrections. Therefore, we recommend
546 flagging matches obtained during this step and analyzing their influence on downstream
547 analyses to account for such potential issues ([Box 2](#)), e.g. .by randomizing the accepted
548 fuzzy matched names based on their score

549

550

551 **Table 3. Number of species matched using each workflow.** Numbers of species matched
 552 were calculated after performing Step 2 but before performing Step 3.

Workflow name	WF1	WF2	WF3	WF4
→	(Direct taxon-	(Pre-assign	(GBIF with	(GBIF without
--	specific)	taxonomic	preprocessing)	preprocessing)
Taxonomic		groups)		
Group ↓				
Birds	878	877	1092	1093
Fishes	5123	5122	5491	5496
Vascular Plants	4435	4412	4647	4649
Other	-	-	19458	19466

553

555 **Box 3. Recommendations and best practices for robust taxonomic harmonization.**

Target group	Recommendations
Users	<ol style="list-style-type: none"> 1. Learn common principles of taxonomy to be able to develop a meaningful workflow and to understand potential outputs of the used tools. 2. Use single-taxon-group databases to get the most reliable resources of taxonomic authorities. 3. Use the most recently updated databases to get the most up-to-date taxonomic knowledge. 4. Parse taxonomic names with specific tools to standardize their writing style (e.g. rgnparser). 5. If some data is already matched against one taxonomic database, use this database as a basis to harmonize the rest of the data to avoid mixing different taxonomic concepts and potential spelling styles. 6. Flag potentially inaccurate matches (fuzzy matching, orthographic corrections) for sensitivity analyses. 7. Describe your taxonomic harmonization workflow in detail, for both credit and reproducibility (e.g. Which databases and packages were used?; mention the used software and database versions; Which functions and steps were taken and why?).
Package developers	<ol style="list-style-type: none"> 1. Use updated and at best regularly maintained taxonomic databases. 2. Use infrastructure packages to enforce standard methods. 3. Check if other packages already provide the functionality to avoid duplication of tools e.g. start checking with taxharmonizexplorer https://mgrenie.shinyapps.io/taxharmonizexplorer/. 4. Put your package in a standardized repository (CRAN, Bioconductor) or at least in a long-term archive (Zenodo, OpenScienceFramework). 5. Contribute to other tools that provide similar functionality rather than create your own. 6. Use multi-language tags (keywords), and at best short abstracts in several UN languages to make them better discoverable. 7. If your tool accesses a database, always report the date of access and version of the database; if you know the database has been

	<p>superseded, issue a warning to the users.</p> <p>8. Publish widely (targeting all end user research communities) release notes about a new tool and new major updates.</p>
Database managers	<ol style="list-style-type: none"> 1. Provide detailed information on how the database was compiled: cite original publications. 2. Use harmonized explicit grammar and spelling styles rules of the taxon names and communicate them clearly. 3. Develop new databases and tools as much as possible consistent with what is already out there: don't force users to adopt a new workflow. 4. Detail publicly the links between your database and other existing databases (which backbone is it using, etc.). 5. Give clear version numbers and dates to the different versions of your database and communicate it clearly to your users (what is the update frequency and how to identify it?). 6. Give clear citation guidelines of the database as a structured file such as a BibTeX file. 7. Publish widely (targeting all end user research communities) release notes about a new database and major updates.

556

557 Conclusion

558 The correct treatment of taxon names is a prerequisite for robust biodiversity research. We
559 proposed a typology of widely used taxonomic databases and extensively reviewed R
560 packages that work with taxonomic data. Throughout our review we identified several areas
561 to be improved aiming for more integrated and user-friendly resources and processes to
562 harmonize taxon names ([Box 3](#)). Many issues we came across could have been prevented
563 by a more open and inclusive communication across research communities (e.g. ecologists,
564 data scientists, taxonomists). For instance, rigorous and widely spread communication on
565 important new or updated taxonomic resources or relevant tools would help prevent using
566 outdated data or developing redundant tools either as end user or developer. We suggest
567 publishing short release notes of taxonomic databases and tools (and major updates of

568 them) also in target journals of the respective user communities (often possible additionally
569 to data papers).

570 On a technical side, we specifically see the design and documentation of taxonomic
571 databases and tools as a major field to improve. We urge any researcher and potential tool
572 developer starting with taxonomic name harmonization to do a thorough search for the most
573 suitable (i.e. most reliable, most up-to date) databases and existing related tools. Users
574 should also document fully their harmonization workflow (software versions, functions,
575 parameters, database versions) for the sake of reproducibility. Vice versa, database
576 managers and tool developers need to make their resources discoverable for all researchers
577 globally and describe them with all necessary meta-data ([Box 3](#)). From our review, it is clear
578 that joint efforts between taxonomists and ecologists are strongly needed to understand how
579 these two related fields can inform each other better, improving taxonomic harmonization on
580 one side and making use of and improving existing tools and functions on the other.

581 Teaching and workshops focused on taxonomic name harmonization could foster knowledge
582 and best practices while helping connect both disciplines.

583 What can the broad research community do to support these services for many of us? We
584 can start by acknowledging more this type of community service, e.g. in similar ways as for
585 reviewing papers. Developing and especially maintaining databases and tools, used by
586 many, should be more visible and valuable than just counting citations. Scientific evaluation
587 should fully comprise these aspects. And developers and data managers should mention
588 these services prominently in their CVs. Funding agencies should also fund these types of
589 projects and specifically their long term maintenance or should support, at least, relevant
590 existing structures, which could serve as home for these resources.

591 Ultimately we are convinced that joint synthesis efforts across research communities
592 towards a comprehensive resource overviewing taxonomic databases and useful tools,
593 including meta-data and dependencies, will help any user to discover and work with the most
594 suitable and robust information. This resource could be hosted, for example, on platforms
595 already offering global cross-taxa information such as COL [37]. The research community
596 will always need taxonomic experts and initiatives working on these individual resources, but
597 we, as users, also need more guidance on where to find them and how to use them best.
598 Our review and the shiny app can only be a start, even hopefully a very useful one.

599 Acknowledgements

600 The authors acknowledge the support of iDiv funded by the German Research Foundation
601 (DFG–FZT 118, 202548816). The authors would like to thank all the participants of the iDiv
602 taxonomic harmonization workshop that led to the ideas developed in this paper; Sulochana
603 Swathi Kannan for her work on the characteristics of databases; Martin Freiberg, David
604 Schellenberger-Costa, Alexander Zizka, and Markus Döring for fruitful discussions about
605 taxonomic harmonization in practice; Erik F. Y. Hom for a friendly review; Brian Maitner as
606 well as one anonymous reviewer for their inputs on our work. We thank all the database
607 managers who answered our data collection emails: Darrel Frost, Michael D. Guiry, Jennifer
608 Hammock, Chantal Huijbers, Congtian Lin, Johan Liljeblad, Paul Kirk, and Chris Raper.

609 We are grateful to all the many colleagues, often not acknowledged enough, who built,
610 curated, and maintained the mentioned data and tools and continue to do so.

611 Authors' contributions

612 MW initiated the project. All authors conceived the ideas of the manuscript. MG led the
613 writing of the manuscript with substantial contributions from JCQ and MW. AS and MG led
614 the development of the companion shiny app. GMLD acquired data on databases. JCQ and
615 EB developed the example workflows. All authors contributed critically to the drafts and gave
616 final approval for publication.

617 Data Availability

618 Code and data available on GitHub (https://github.com/Rekyt/taxo_harmonization) with a
619 perennial archive on Zenodo (<https://doi.org/10.5281/zenodo.5121244>). The repository
620 contains the table of included packages and network links. It contains the code to run the
621 shiny app **taxharmonizexplorer**. The online shiny app is available at
622 <https://mgrenie.shinyapps.io/taxharmonizexplorer/>

623 References

- 624 Bivand, R. S., Pebesma, E. J., & Gomez-Rubio, V. (2013). *Applied spatial data analysis with*
625 *R, Second edition*. Springer, NY. <https://asdar-book.org/>
- 626 Boettiger, C., Chamberlain, S., Hart, E., & Ram, K. (2015). Building Software, Building
627 Community: Lessons from the rOpenSci Project. *Journal of Open Research*
628 *Software*, 3(1). <https://doi.org/10.5334/jors.bu>
- 629 Bortolus, A. (2008). Error Cascades in the Biological Sciences: The Unwanted
630 Consequences of Using Bad Taxonomy in Ecology. *AMBIO: A Journal of the Human*
631 *Environment*, 37(2), 114–118. <https://doi.org/10.1579/0044->

632 7447(2008)37[114:ECITBS]2.0.CO;2

633 Chawuthai, R., Takeda, H., Wuwongse, V., & Jinbo, U. (2016). Presenting and preserving
634 the change in taxonomic knowledge for linked data. *Semantic Web*, 7(6), 589–616.

635 Consortium, B. (2018). *BioTIME* [Data set]. Zenodo. <https://doi.org/10.5281/zenodo.3265871>

636 Costello, M. J. (2020). Taxonomy as the key to life. *Megataxa*, 1(2), 105–113.
637 <https://doi.org/10.11646/megataxa.1.2.1>

638 Costello, M. J., Bouchet, P., Boxshall, G., Fauchald, K., Gordon, D., Hoeksema, B. W.,
639 Poore, G. C. B., Soest, R. W. M. van, Stöhr, S., Walter, T. C., Vanhoorne, B.,
640 Decock, W., & Appeltans, W. (2013). Global Coordination and Standardisation in
641 Marine Biodiversity through the World Register of Marine Species (WoRMS) and
642 Related Databases. *PLOS ONE*, 8(1), e51629.
643 <https://doi.org/10.1371/journal.pone.0051629>

644 Dayrat, B. (2005). Towards integrative taxonomy. *Biological Journal of the Linnean Society*,
645 85(3), 407–417. <https://doi.org/10.1111/j.1095-8312.2005.00503.x>

646 Dornelas, M., Antão, L. H., Moyes, F., Bates, A. E., Magurran, A. E., Adam, D.,
647 Akhmetzhanova, A. A., Appeltans, W., Arcos, J. M., Arnold, H., Ayyappan, N., Badihi,
648 G., Baird, A. H., Barbosa, M., Barreto, T. E., Bässler, C., Bellgrove, A., Belmaker, J.,
649 Benedetti-Cecchi, L., ... Zettler, M. L. (2018). BioTIME: A database of biodiversity
650 time series for the Anthropocene. *Global Ecology and Biogeography*, 27(7), 760–786.
651 <https://doi.org/10.1111/geb.12729>

652 Dyer, E. E., Redding, D. W., & Blackburn, T. M. (2017). The global avian invasions atlas, a
653 database of alien bird distributions worldwide. *Scientific Data*, 4(1), 170041.
654 <https://doi.org/10.1038/sdata.2017.41>

655 Edwards, J. L., Lane, M. A., & Nielsen, E. S. (2000). Interoperability of Biodiversity
656 Databases: Biodiversity Information on Every Desktop. *Science*, 289(5488), 2312–

657 2314. <https://doi.org/10.1126/science.289.5488.2312>

658 Farley, S. S., Dawson, A., Goring, S. J., & Williams, J. W. (2018). Situating Ecology as a Big-
659 Data Science: Current Advances, Challenges, and Solutions. *BioScience*, 68(8),
660 563–576. <https://doi.org/10.1093/biosci/biy068>

661 Freiberg, M., Winter, M., Gentile, A., Zizka, A., Muellner-Riehl, A. N., Weigelt, A., & Wirth, C.
662 (2020). LCVP, The Leipzig catalogue of vascular plants, a new taxonomic reference
663 list for all known vascular plants. *Scientific Data*, 7(1), 416.
664 <https://doi.org/10.1038/s41597-020-00702-z>

665 Garnett, S. T., Christidis, L., Conix, S., Costello, M. J., Zachos, F. E., Bánki, O. S., Bao, Y.,
666 Barik, S. K., Buckeridge, J. S., Hobern, D., Lien, A., Montgomery, N., Nikolaeva, S.,
667 Pyle, R. L., Thomson, S. A., Dijk, P. P. van, Whalen, A., Zhang, Z.-Q., & Thiele, K. R.
668 (2020). Principles for creating a single authoritative list of the world's species. *PLOS*
669 *Biology*, 18(7), e3000736. <https://doi.org/10.1371/journal.pbio.3000736>

670 GBIF: The Global Biodiversity Information Facility. (2020, June 24). What is GBIF? *GBIF*.
671 <https://www.gbif.org/what-is-gbif>

672 Grenié, M., Sagouis, A., & Berti, E. (2021). *Rekyt/taxo_harmonization: First release*. Zenodo.
673 <https://doi.org/10.5281/zenodo.5121245>

674 Hampton, S. E., Strasser, C. A., Tewksbury, J. J., Gram, W. K., Budden, A. E., Batcheller, A.
675 L., Duke, C. S., & Porter, J. H. (2013). Big data and the future of ecology. *Frontiers in*
676 *Ecology and the Environment*, 11(3), 156–162. <https://doi.org/10.1890/120103>

677 Hassler, M. (2021). *World Plants. Synonymic Checklist and Distribution of the World Flora.*
678 *Version 12.4; last update August 6th, 2021.* -. www.worldplants.de

679 Isaac, N. J. B., Mallet, J., & Mace, G. M. (2004). Taxonomic inflation: Its influence on
680 macroecology and conservation. *Trends in Ecology & Evolution*, 19(9), 464–469.
681 <https://doi.org/10.1016/j.tree.2004.06.004>

682 IUCN. (2021). *The IUCN Red List of Threatened Species. Version 2021-1.*
683 www.iucnredlist.org

684 Jay, C., Haines, R., & Katz, D. S. (2020). Software must be recognised as an important
685 output of scholarly research. *ArXiv:2011.07571 [Cs]*. <http://arxiv.org/abs/2011.07571>

686 Jin, J., & Yang, J. (2020). BDCleaner: A workflow for cleaning taxonomic and geographic
687 errors in occurrence data archived in biodiversity databases. *Global Ecology and*
688 *Conservation*, 21, e00852. <https://doi.org/10.1016/j.gecco.2019.e00852>

689 Jones, K. E., Bielby, J., Cardillo, M., Fritz, S. A., O'Dell, J., Orme, C. D. L., Safi, K., Sechrest,
690 W., Boakes, E. H., Carbone, C., Connolly, C., Cutts, M. J., Foster, J. K., Grenyer, R.,
691 Habib, M., Plaster, C. A., Price, S. A., Rigby, E. A., Rist, J., ... Purvis, A. (2009).
692 PanTHERIA: A species-level database of life history, ecology, and geography of
693 extant and recently extinct mammals. *Ecology*, 90(9), 2648–2648.
694 <https://doi.org/10.1890/08-1494.1>

695 Jones, M. B., Schildhauer, M. P., Reichman, O. J., & Bowers, S. (2006). The New
696 Bioinformatics: Integrating Ecological Data from the Gene to the Biosphere. *Annual*
697 *Review of Ecology, Evolution, and Systematics*, 37(1), 519–544.
698 <https://doi.org/10.1146/annurev.ecolsys.37.091305.110031>

699 Kattge, J., Bönisch, G., Díaz, S., Lavorel, S., Prentice, I. C., Leadley, P., Tautenhahn, S.,
700 Werner, G. D. A., Aakala, T., Abedi, M., Acosta, A. T. R., Adamidis, G. C., Adamson,
701 K., Aiba, M., Albert, C. H., Alcántara, J. M., Alcázar, C. C., Aleixo, I., Ali, H., ... Wirth,
702 C. (2020). TRY plant trait database – enhanced coverage and open access. *Global*
703 *Change Biology*, 26(1), 119–188. <https://doi.org/10.1111/gcb.14904>

704 Kissling, W. D., Walls, R., Bowser, A., Jones, M. O., Kattge, J., Agosti, D., Amengual, J.,
705 Basset, A., van Bodegom, P. M., Cornelissen, J. H. C., Denny, E. G., Deudero, S.,
706 Egloff, W., Elmendorf, S. C., Alonso García, E., Jones, K. D., Jones, O. R., Lavorel,

707 S., Lear, D., ... Guralnick, R. P. (2018). Towards global data products of Essential
708 Biodiversity Variables on species traits. *Nature Ecology & Evolution*, 2(10), 1531–
709 1540. <https://doi.org/10.1038/s41559-018-0667-3>

710 König, C., Weigelt, P., Schrader, J., Taylor, A., Kattge, J., & Kreft, H. (2019). Biodiversity
711 data integration—The significance of data resolution and domain. *PLOS Biology*,
712 17(3), e3000183. <https://doi.org/10.1371/journal.pbio.3000183>

713 La Salle, J., Williams, K. J., & Moritz, C. (2016). Biodiversity analysis in the digital era.
714 *Philosophical Transactions of the Royal Society B: Biological Sciences*, 371(1702),
715 20150337. <https://doi.org/10.1098/rstb.2015.0337>

716 Lai, J., Lortie, C. J., Muenchen, R. A., Yang, J., & Ma, K. (2019). Evaluating the popularity of
717 R in ecology. *Ecosphere*, 10(1), e02567.

718 Lenters, T. P., Henderson, A., Dracxler, C. M., Elias, G. A., Kamga, S. M., Couvreur, T. L.
719 P., & Kissling, W. D. (2021). Integration and harmonization of trait data from plant
720 individuals across heterogeneous sources. *Ecological Informatics*, 62, 101206.
721 <https://doi.org/10.1016/j.ecoinf.2020.101206>

722 Lepage, D., Vaidya, G., & Guralnick, R. (2014). Avibase – a database system for managing
723 and organizing taxonomic concepts. *ZooKeys*, 420, 117–135.
724 <https://doi.org/10.3897/zookeys.420.7089>

725 Meyer, C., Weigelt, P., & Kreft, H. (2016). Multidimensional biases, gaps and uncertainties in
726 global plant occurrence information. *Ecology Letters*, 19(8), 992–1006.
727 <https://doi.org/10.1111/ele.12624>

728 Michener, W. K., & Jones, M. B. (2012). Ecoinformatics: Supporting ecology as a data-
729 intensive science. *Trends in Ecology & Evolution*, 27(2), 85–93.
730 <https://doi.org/10.1016/j.tree.2011.11.016>

731 Mislán, K. A. S., Heer, J. M., & White, E. P. (2016). Elevating The Status of Code in Ecology.

732 *Trends in Ecology & Evolution*, 31(1), 4–7. <https://doi.org/10.1016/j.tree.2015.11.006>

733 Mozzherin, D. (2021). *gnames/gnverifier: V0.3.3*. Zenodo.

734 <https://doi.org/10.5281/zenodo.5111543>

735 Mozzherin, D., Myltsev, A. A., & Patterson, D. J. (2017). “gnparser”: A powerful parser for

736 scientific names based on Parsing Expression Grammar. *BMC Bioinformatics*, 18(1),

737 279. <https://doi.org/10.1186/s12859-017-1663-3>

738 Mozzherin, D., Shorthouse, D., ashipunova, & pdevries. (2021).

739 *GlobalNamesArchitecture/gni: V0.9.40 Global Names Index (no fuzzy matching)*.

740 Zenodo. <https://doi.org/10.5281/zenodo.5121908>

741 Nelson, G., & Ellis, S. (2019). The history and impact of digitization and digital data

742 mobilization on biodiversity research. *Philosophical Transactions of the Royal Society*

743 *B: Biological Sciences*, 374(1763), 20170391. <https://doi.org/10.1098/rstb.2017.0391>

744 Patterson, D. J., Cooper, J., Kirk, P. M., Pyle, R. L., & Remsen, D. P. (2010). Names are key

745 to the big new biology. *Trends in Ecology & Evolution*, 25(12), 686–691.

746 <https://doi.org/10.1016/j.tree.2010.09.004>

747 Patterson, D., Mozzherin, D., Shorthouse, D., & Thessen, A. (2016). Challenges with using

748 names to link digital biodiversity information. *Biodiversity Data Journal*, 4, e8080.

749 <https://doi.org/10.3897/BDJ.4.e8080>

750 Pebesma, E. J. (2018). Simple features for r: Standardized support for spatial vector data.

751 *The R Journal*, 10(1), 439–446. <https://doi.org/10.32614/RJ-2018-009>

752 Pebesma, E. J., & Bivand, R. (2005). Classes and methods for spatial data in R. *R News*,

753 5(2), 9–13.

754 Rees, T. (2014). Taxamatch, an Algorithm for Near (‘Fuzzy’) Matching of Scientific Names in

755 Taxonomic Databases. *PLOS ONE*, 9(9), e107510.

756 <https://doi.org/10.1371/journal.pone.0107510>

757 Rees, T. (2021). *The interim register of marine and nonmarine genera*. Available from
758 <https://www.irmng.org> at VLIZ. Accessed 2021-11-15. <https://www.irmng.org>

759 Rouhan, G., & Gaudeul, M. (2021). Plant Taxonomy: A Historical Perspective, Current
760 Challenges, and Perspectives. In P. Besse (Ed.), *Molecular Plant Taxonomy:
761 Methods and Protocols* (pp. 1–38). Springer US. [https://doi.org/10.1007/978-1-0716-
762 0997-2_1](https://doi.org/10.1007/978-1-0716-0997-2_1)

763 Schulman, L., Lahti, K., Piirainen, E., Heikkinen, M., Raitio, O., & Juslén, A. (2021). The
764 Finnish Biodiversity Information Facility as a best-practice model for biodiversity data
765 infrastructures. *Scientific Data*, *8*(1), 137. [https://doi.org/10.1038/s41597-021-00919-
766 6](https://doi.org/10.1038/s41597-021-00919-6)

767 Shipunov, A. (2011). The problem of hemihomonyms and the on-line hemihomonyms
768 database (HHDB). *Bionomina*, *4*(1), 65-72-65–72.
769 <https://doi.org/10.11646/bionomina.4.1.3>

770 Smith, S. A., & Brown, J. W. (2018). Constructing a broadly inclusive seed plant phylogeny.
771 *American Journal of Botany*, *105*(3), 302–314. <https://doi.org/10.1002/ajb2.1019>

772 Soberón, J., & Peterson, T. (2004). Biodiversity informatics: Managing and applying primary
773 biodiversity data. *Philosophical Transactions of the Royal Society of London. Series
774 B: Biological Sciences*, *359*(1444), 689–698. <https://doi.org/10.1098/rstb.2003.1439>

775 Tessarolo, G., Ladle, R., Rangel, T., & Hortal, J. (2017). Temporal degradation of data limits
776 biodiversity research. *Ecology and Evolution*, *7*(17), 6863–6870.
777 <https://doi.org/10.1002/ece3.3259>

778 Thomas, C. (2009). Biodiversity Databases Spread, Prompting Unification Call. *Science*,
779 *324*(5935), 1632–1633. https://doi.org/10.1126/science.324_1632

780 Upham, N. S., Esselstyn, J. A., & Jetz, W. (2019). Inferring the mammal tree: Species-level
781 sets of phylogenies for questions in ecology, evolution, and conservation. *PLOS*

782 *Biology*, 17(12), e3000494. <https://doi.org/10.1371/journal.pbio.3000494>

783 van Kleunen, M., Pyšek, P., Dawson, W., Essl, F., Kreft, H., Pergl, J., Weigelt, P., Stein, A.,
784 Dullinger, S., König, C., Lenzner, B., Maurel, N., Moser, D., Seebens, H., Kartesz, J.,
785 Nishino, M., Aleksanyan, A., Ansong, M., Antonova, L. A., ... Winter, M. (2019). The
786 Global Naturalized Alien Flora (GloNAF) database. *Ecology*, 100(1), e02542.
787 <https://doi.org/10.1002/ecy.2542>

788 Vanden Berghe, E., Coro, G., Bailly, N., Fiorellato, F., Aldemita, C., Ellenbroek, A., &
789 Pagano, P. (2015). Retrieving taxa names from large biodiversity data collections
790 using a flexible matching workflow. *Ecological Informatics*, 28, 29–41.
791 <https://doi.org/10.1016/j.ecoinf.2015.05.004>

792 Wüest, R. O., Zimmermann, N. E., Zurell, D., Alexander, J. M., Fritz, S. A., Hof, C., Kreft, H.,
793 Normand, S., Cabral, J. S., Szekely, E., Thuiller, W., Wikelski, M., & Karger, D. N.
794 (2020). Macroecology in the age of Big Data – Where to go from here? *Journal of*
795 *Biogeography*, 47(1), 1–12. <https://doi.org/10.1111/jbi.13633>

796 Supplementary References

797 [1] F. Jansen and J. Dengler, "Germansl-eine universelle taxonomische referenzliste fuer
798 vegetationsdatenbanken," *Tuexenia*, vol. 28, pp. 239–253, 2008.

799 [2] USDA, NRCS, "The PLANTS database", [Online]. Available: <http://plants.usda.gov>

800 [3] L. Brouillet et al., "Database of Vascular Plants of Canada (VASCAN)," Online
801 <https://datacanadensysnetvascan> <https://www.gbif.org/dataset/3f8a1297-3259-4700-91fc-Acc4170b27ce>
802 Released 2010-12-10 Version Xx GBIF Key 3f8a1297-3259-4700-91fc-Acc4170b27ce Consult. 2021-
803 12-08, 2010, doi: 10.5886/zw3aqw.

- 804 [4] J. Liqiang, C. Lin, and H. Wangjiangning, "China animal scientific database," Available
805 [Httpwwwzoologycsdbc.cn](http://www.zoologycsdbc.cn), 2021, Accessed: Jul. 27, 2021. [Online]. Available:
806 <https://datapid.cn/CSTR:30689.11.16273678284380>
- 807 [5] Brazil Flora Group, "Brazilian Flora 2020 project - Projeto Flora do Brasil 2020," 2021, doi:
808 10.15468/1mtkaw.
- 809 [6] J. Liljeblad, "Dyntaxa. Svensk taxonomisk databas," Jul. 2021, doi: 10.15468/j43wfc.
- 810 [7] L. Halada et al., "European Biodiversity Observation Network - EBONE," 2009, Accessed: Jul.
811 26, 2021. [Online]. Available: <http://dl.gi.de/handle/20.500.12116/27839>
- 812 [8] L. Schulman, K. Lahti, E. Piirainen, M. Heikkinen, O. Raitio, and A. Juslén, "The Finnish
813 Biodiversity Information Facility as a best-practice model for biodiversity data infrastructures," *Sci.*
814 *Data*, vol. 8, no. 1, p. 137, May 2021, doi: 10.1038/s41597-021-00919-6.
- 815 [9] NBN Atlas, "National Biodiversity Network," NBN Atlas Website [Httpwww Nbnatlas Org](http://www.nbnatlas.org), 2021.
- 816 [10] Y. de Jong et al., "PESI - a taxonomic backbone for Europe," *Biodivers. Data J.*, vol. 3, p.
817 e5848, Sep. 2015, doi: 10.3897/BDJ.3.e5848.
- 818 [11] L. Q. Ji and al., "China checklist of animals," in *Catalogue of life China: Annual checklist*, B.
819 Committee, Ed. Beijing, China: Chinese Academy of Sciences, 2021. [Online]. Available:
820 <http://www.sp2000.org.cn>
- 821 [12] K. T. Shao, H. Lee, and Y. C. Lin, "TaiBNET (Catalogue of Life in Taiwan)," *Cat. Life Taiwan*–
822 [http://taibnet Sin. Edu Tw Taiwan](http://taibnet.sin.edu.tw/Taiwan), 2020.
- 823 [13] O. Gargominy et al., "TAXREF v14.0, référentiel taxonomique pour la France," in *UMS*
824 *PatriNat*, Muséum national d'Histoire naturelle, Paris, 2020. [Online]. Available:
825 <https://inpn.mnhn.fr/telechargement/referentielEspece/taxref/14.0/menu>
- 826 [14] A. Veen, "Taxa Watermanagement the Netherlands (TWN)," Apr. 2020, doi: 10.15468/yz4wjs.
- 827 [15] M. D. Guiry et al., "AlgaeBase: An On-line Resource for Algae," *Cryptogam. Algol.*, vol. 35,
828 no. 2, pp. 105–115, Jun. 2014, doi: 10.7872/crya.v35.iss2.2014.105.

- 829 [16] D. R. Frost, "Amphibian species of the world: an online reference version 6.1," Electron.
830 Database Access. Am. Mus. Nat. Hist. N. Y. USA, 2021, doi: 10.5531/db.vz.0001.
- 831 [17] J. F. Clements et al., "The eBird/Clements Checklist of Birds of the World: v2019," Cornell
832 Lab Ornithol. Ithaca N. Y., 2019, [Online]. Available:
833 <https://www.birds.cornell.edu/clementschecklist/download/>
- 834 [18] P. M. Kirk, "World catalogue of 340 K fungal names on-line.," Mycol. Res., vol. 104, no. 5, pp.
835 516–517, 2000.
- 836 [19] R. Froese and D. Pauly, "FishBase. World Wide Web electronic publication, version
837 (02/2021)," 2021, [Online]. Available: www.fishbase.org
- 838 [20] C. J. Burgin, J. P. Colella, P. L. Kahn, and N. S. Upham, "How many species of mammals are
839 there?," J. Mammal., vol. 99, no. 1, pp. 1–14, Feb. 2018, doi: 10.1093/jmammal/gyx147.
- 840 [21] The International Plant Names Index Collaborators, "International Plant Names Index," Oct.
841 2019, doi: 10.15468/uhlmlw.
- 842 [22] M. Freiberg et al., "LCVP, The Leipzig catalogue of vascular plants, a new taxonomic
843 reference list for all known vascular plants," Sci. Data, vol. 7, no. 1, Art. no. 1, Nov. 2020, doi:
844 10.1038/s41597-020-00702-z.
- 845 [23] POWO, "Plants of the World Online." Facilitated by Royal Botanic Gardens Kew. Published
846 on the Internet., 2021. Accessed: Jul. 26, 2021. [Online]. Available:
847 <http://www.plantsoftheworldonline.org/>
- 848 [24] "The plant list, version 1.1," 2013, [Online]. Available: <http://www.theplantlist.org>
- 849 [25] B. Boyle et al., "The taxonomic name resolution service: an online tool for automated
850 standardization of plant names," BMC Bioinformatics, vol. 14, no. 1, p. 16, 2013, doi: 10.1186/1471-
851 2105-14-16.
- 852 [26] B. L. Boyle et al., "Taxonomic name resolution service, version 5.0," Bot. Inf. Ecol. Netw.,
853 2021, Accessed: Apr. 22, 2021. [Online]. Available: <https://tnrs.biendata.org/>
- 854 [27] Missouri Botanical Garden, "Tropicos", [Online]. Available: <http://www.tropicos.org>

- 855 [28] WCSP, "World Checklist of Selected Plant Families." Facilitated by Royal Botanic Gardens
856 Kew. Published on the Internet., 2021. Accessed: Jul. 26, 2021. [Online]. Available:
857 <http://wcsp.science.kew.org/>
- 858 [29] R. Govaerts, E. Nic Lughadha, N. Black, R. Turner, and A. Paton, "The World Checklist of
859 Vascular Plants, a continuously updated resource for exploring global plant diversity," *Sci. Data*, vol.
860 8, no. 1, p. 215, Dec. 2021, doi: 10.1038/s41597-021-00997-6.
- 861 [30] T. Borsch et al., "World Flora Online: Placing taxonomists at the heart of a definitive and
862 comprehensive global resource on the world's plants," *TAXON*, vol. 69, no. 6, pp. 1311–1341, Dec.
863 2020, doi: 10.1002/tax.12373.
- 864 [31] P. Uetz et al., "A quarter century of reptile and amphibian databases," *Herpetol. Rev.*, vol. 52,
865 no. 2, pp. 246–255, 2021.
- 866 [32] W. S. Catalog, "World Spider Catalog. Version 22.5.," *Nat. Hist. Mus. Bern Online*
867 [Httpwscnmbch](http://www.wscnmb.ch/) Accessed 2021-11-22, 2021, doi: 10.24436/2.
- 868 [33] M. L. D. Palomares and D. Pauly, "SeaLifeBase," *World Wide Web Electron. Publ.* [Httpwww](http://www.sealifebase.org)
869 [Sealifebase Org](http://www.sealifebase.org) Version 042021, 2021, [Online]. Available: <http://www.sealifebase.org>
- 870 [34] T. Horton et al., "World register of marine species (WoRMS)," Apr. 2021, [Online]. Available:
871 <https://www.marinespecies.org>
- 872 [35] R. L. Pyle and E. Michel, "ZooBank: Developing a nomenclatural tool for unifying 250 years of
873 biological information," *Zootaxa*, vol. 1950, no. 1, pp. 39–50, 2008.
- 874 [36] S. Ratnasingham and P. D. Hebert, "BOLD: The Barcode of Life Data System ([http://www.](http://www.barcodinglife.org)
875 [barcodinglife.org](http://www.barcodinglife.org))," *Mol. Ecol. Notes*, vol. 7, no. 3, pp. 355–364, 2007.
- 876 [37] O. Bánki, "New Catalogue of Life: infrastructure," *COL*, Oct. 18, 2020.
877 <https://www.catalogueoflife.org/2020/10/18/new-col-infrastructure> (accessed May 27, 2021).
- 878 [38] C. S. Parr et al., "The Encyclopedia of Life v2: Providing Global Access to Knowledge About
879 Life on Earth," *Biodivers. Data J.*, vol. 2, p. e1079, Apr. 2014, doi: 10.3897/BDJ.2.e1079.

880 [39] GBIF: The Global Biodiversity Information Facility, "What is GBIF?," GBIF, Jun. 24, 2020.
881 <https://www.gbif.org/what-is-gbif>

882 [40] D. Mozzherin, D. Shorthouse, ashipunova, and pdevries, GlobalNamesArchitecture/gni:
883 v0.9.40 Global Names Index (no fuzzy matching). Zenodo, 2021. doi: 10.5281/zenodo.5121908.

884 [41] D. Mozzherin, D. Shorthouse, ashipunova, and pdevries, GlobalNamesArchitecture/gni:
885 v2.0.0 Resolver. Zenodo, 2021. doi: 10.5281/zenodo.5137754.

886 [42] D. Mozzherin, gnames/gnverifier: v0.3.3. Zenodo, 2021. doi: 10.5281/zenodo.5111543.

887 [43] Clarivate, "Index to Organism Names." World Wide Web electronic resource., 2020. [Online].
888 Available: www.organismnames.com.

889 [44] ITIS, "Retrieved [April, 22, 2021], from the Integrated Taxonomic Information System (ITIS),
890 www.itis.gov," Apr. 2021, doi: <https://doi.org/10.5066/F7KH0KBK>.

891 [45] IUCN, The IUCN Red List of Threatened Species. Version 2021-1. 2021. [Online]. Available:
892 www.iucnredlist.org

893 [46] NatureServe, "NatureServe explorer [web application]," 2021, Accessed: Apr. 22, 2021.
894 [Online]. Available: <https://explorer.natureserve.org/>.

895 [47] C. L. Schoch et al., "NCBI Taxonomy: a comprehensive update on curation, resources and
896 tools," Database J. Biol. Databases Curation, vol. 2020, Jan. 2020, doi: 10.1093/database/baaa062.

897 [48] J. W. Williams et al., "The Neotoma Paleocology Database, a multiproxy, international,
898 community-curated data resource," Quat. Res., vol. 89, no. 1, pp. 156–177, 2018.

899 [49] C. E. Hinchliff et al., "Synthesis of phylogeny and taxonomy into a comprehensive tree of life,"
900 Proc. Natl. Acad. Sci., vol. 112, no. 41, pp. 12764–12769, Oct. 2015, doi: 10.1073/pnas.1423041112.

901 [50] Paleobiology Database, "Paleobiology database," 2021, [Online]. Available:
902 <https://paleobiodb.org/>

903 [51] F. Zachary, C. Scott, and G. Niklaus, "Taxa: An R package implementing data standards and
904 methods for taxonomic data.," manual 272, 2018. doi: 10.12688/f1000research.14013.1.

905 [52] M. Alvarez and F. Luebert, "The taxlist package: Managing plant taxonomic lists in R,"
906 Biodivers. Data J., vol. 6, p. e23635, 2018, doi: 10.3897/bdj.6.e23635.

907 [53] S. Chamberlain and C. Boettiger, "taxview: Tools for vizualizing data taxonomically," manual,
908 2021.

909 [54] V. Patil, T. Seltmann, N. Salmaso, O. Anneville, M. Lajeunesse, and D. Straile,
910 "algaeClassify: Determine phytoplankton functional groups based on functional traits," manual, 2019.
911 [Online]. Available: <https://CRAN.R-project.org/package=algaeClassify>

912 [55] H. C. Liedtke, "AmphiNom: an amphibian systematics tool," Syst. Biodivers., vol. 17, pp. 1–6,
913 2019.

914 [56] P. Cardoso, "arakno: ARachnid KNowledge online," manual, 2021. [Online]. Available:
915 <https://CRAN.R-project.org/package=arakno>

916 [57] W. K. Morris, "Introduction to the finbif package," manual, 2021. doi:
917 10.5281/zenodo.3612814.

918 [58] B. Walker, "kewr: R package to access kew data APIs," manual, 2021. [Online]. Available:
919 <https://barnabywalker.github.io/kewr/>

920 [59] S. Chamberlain, "natserv: 'NatureServe' interface," manual, 2020. [Online]. Available:
921 <https://CRAN.R-project.org/package=natserv>

922 [60] J. Eastman, "ncbit: retrieve and build NCBI taxonomic data," manual, 2013. [Online].
923 Available: <https://CRAN.R-project.org/package=ncbit>

924 [61] S. Goring, "neotoma2: Working with the neotoma paleoecology database," manual, 2021.

925 [62] Sara Varela, Javier Gonzalez-Hernandez, and Luciano Fabris Sgarbi, "paleobioDB: an R-
926 package for downloading, visualizing and processing data from the Paleobiology Database," manual,
927 2014.

928 [63] J. Zhang, B. Liu, S. Liu, Z. Feng, and K. Jiang, "plantlist: Looking up the status of plant
929 scientific names based on the plant list database, search the chinese names and making checklists of
930 plants," manual, 2021. [Online]. Available: <https://github.com/helixcn/plantlist/>

931 [64] S. Chamberlain, "rcol: Catalogue of life client," manual, 2021. [Online]. Available:
932 <https://CRAN.R-project.org/package=rcol>

933 [65] R. Maia, S. Chamberlain, A. Teucher, and S. Pardo, "rebird: R client for the eBird database of
934 bird observations," manual, 2021. [Online]. Available: <https://CRAN.R-project.org/package=rebird>

935 [66] D. J. Winter, "rentrez: an R package for the NCBI eUtils API," *R J.*, vol. 9, no. 2, pp. 520–526,
936 2017.

937 [67] C. Boettiger, D. Temple Lang, and P. Wainwright, "rfishbase: exploring, manipulating and
938 visualizing FishBase data from R," *J. Fish Biol.*, Nov. 2012, [Online]. Available:
939 <https://doi.org/10.1111/j.1095-8649.2012.03464.x>

940 [68] S. Chamberlain et al., "rgbif: Interface to the global biodiversity information facility API,"
941 manual, 2021. [Online]. Available: <https://CRAN.R-project.org/package=rgbif>

942 [69] S. Chamberlain, "ritis: Integrated taxonomic information system client," manual, 2021.
943 [Online]. Available: <https://CRAN.R-project.org/package=ritis>

944 [70] S. Mortara and A. Sánchez-Tapia, "Rocc: Workflows for biodiversity data download and
945 cleaning," manual, 2021. [Online]. Available: <https://github.com/liibre/rocc>

946 [71] F. Michonneau, J. W. Brown, and D. J. Winter, "rotl: an R package to interact with the Open
947 Tree of Life data," *Methods Ecol. Evol.*, vol. 7, no. 12, pp. 1476–1481, Dec. 2016, doi: 10.1111/2041-
948 210X.12593.

949 [72] S. Chamberlain, "rredlist: 'IUCN' red list client," manual, 2020. [Online]. Available:
950 <https://CRAN.R-project.org/package=rredlist>

951 [73] M. Lewis and O. Tallowin, "rreptiledb: Search and access data from the reptile database
952 through r," manual, 2021. [Online]. Available: <https://github.com/matthewlewis896/rreptiledb>

953 [74] M. Grenié and H. Gruson, "rtaxref: An R client for TaxRef the french taxonomical reference
954 API," manual, 2021. [Online]. Available: <https://github.com/Rekyt/rtaxref>

955 [75] L. Ding et al., "SP2000: An open-sourced R package for querying the Catalogue of Life,"
956 *Biodivers. Sci.*, vol. 29, no. 1, p. 118, Jan. 2021, doi: 10.17520/biods.2020235.

957 [76] S. Chamberlain et al., “taxize: Taxonomic information from around the web,” manual, 2020.
958 [Online]. Available: <https://github.com/ropensci/taxize>

959 [77] S. A. Chamberlain and E. Szöcs, “taxize: taxonomic search and retrieval in R,”
960 F1000Research, vol. 2, Oct. 2013, doi: 10.12688/f1000research.2-191.v2.

961 [78] C. Smith, “taxonomyCleanr: A workflow and set of functions to clean taxonomy data using R,”
962 manual, 2021. [Online]. Available: <https://github.com/EDlorg/taxonomyCleanr>

963 [79] L. Cayuela, Í. Granzow-de la Cerda, F. S. Albuquerque, and D. J. Golicher, “taxonstand: An r
964 package for species names standardisation in vegetation databases,” *Methods Ecol. Evol.*, vol. 3, no.
965 6, pp. 1078–1083, Dec. 2012, doi: 10.1111/j.2041-210X.2012.00232.x.

966 [80] L. Cayuela, I. Macarro, A. Stein, and J. Oksanen, “Taxonstand: Taxonomic standardization of
967 plant species names,” manual, 2021. [Online]. Available: [https://CRAN.R-](https://CRAN.R-project.org/package=Taxonstand)
968 [project.org/package=Taxonstand](https://CRAN.R-project.org/package=Taxonstand)

969 [81] V. Barve, “taxotools: Tools to handle taxonomic lists,” manual, 2021. [Online]. Available:
970 <https://CRAN.R-project.org/package=taxotools>

971 [82] V. Barve, Taxotools: Tools to handle taxonomic lists. R package. Zenodo, 2020. doi:
972 10.5281/zenodo.3934939.

973 [83] A. Piper, “taxreturn: An R package for retrieving and curating public DNA barcode reference
974 data,” manual, 2021. [Online]. Available: <https://github.com/alexpiper/taxreturn>

975 [84] B. Maitner and B. Boyle, “TNRS: Taxonomic name resolution service,” manual, 2021.

976 [85] B. Boyle et al., “The taxonomic name resolution service: an online tool for automated
977 standardization of plant names,” *BMC Bioinformatics*, vol. 14, no. 1, p. 16, 2013, doi: 10.1186/1471-
978 2105-14-16.

979 [86] J. van Tent, “twn: Taxa waterbeheer nederland voor r,” manual, 2021. [Online]. Available:
980 <https://CRAN.R-project.org/package=twn>

981 [87] S. Chamberlain and E. Welty, “wikitaxa: Taxonomic information from ‘wikipedia,’” manual,
982 2020. [Online]. Available: <https://CRAN.R-project.org/package=wikitaxa>

983 [88] J. Holstein, "worms: Retriving aphia information from world register of marine species,"
984 manual, 2018. [Online]. Available: <https://CRAN.R-project.org/package=worms>

985 [89] S. Chamberlain, "worrms: World register of marine species (WoRMS) client," manual, 2020.
986 [Online]. Available: <https://CRAN.R-project.org/package=worrms>

987 [90] S. Chamberlain, "zbank: 'ZooBank' API client," manual, 2018. [Online]. Available:
988 <https://CRAN.R-project.org/package=zbank>

989 [91] Z. Charlop-Powers, "flattax: NCBI flat taxonomy file," manual, 2021.

990 [92] G. Carvalho, "flora: Tools for Interacting with the Brazilian Flora 2020," manual, 2021.
991 [Online]. Available: <http://www.github.com/gustavobio/flora>

992 [93] D. Mbae, "mammals: Tracking the latest taxonomic changes to species and higher groups of
993 mammals," manual, 2021. [Online]. Available: <https://github.com/mawiramawira/mammals>

994 [94] K. E. A. Norman, S. Chamberlain, and C. Boettiger, "taxadb: A high-performance local
995 taxonomic database interface," *Methods Ecol. Evol.*, vol. 11, no. 9, pp. 1153–1159, Sep. 2020, doi:
996 10.1111/2041-210X.13440.

997 [95] C. Boettiger, K. Norman, J. Poelen, and S. Chamberlain, "taxadb: A high-performance local
998 taxonomic database interface," manual, 2021. [Online]. Available: [https://CRAN.R-](https://CRAN.R-project.org/package=taxadb)
999 [project.org/package=taxadb](https://CRAN.R-project.org/package=taxadb)

1000 [96] C. Boettiger and K. Norman, "taxalight: A lightweight and lightning-fast taxonomic naming
1001 interface," manual, 2021. [Online]. Available: <https://CRAN.R-project.org/package=taxalight>

1002 [97] J. Nitta, "taxastand: Standardize taxonomic names," manual, 2021.

1003 [98] S. Chamberlain and Z. Arendsee, "taxizedb: Tools for working with 'taxonomic' databases,"
1004 manual, 2021.

1005 [99] W. Cornwell and R. FitzJohn, *traitecoevo/taxonlookup v1.1.5*. Zenodo, 2017. doi:
1006 10.5281/zenodo.839589.

1007 [100] S. Sherrill-Mix, "taxonomizr: Functions to work with NCBI accessions and taxonomy," manual,
1008 2021. [Online]. Available: <https://CRAN.R-project.org/package=taxonomizr>

1009 [101] G. Carvalho, “tpl: interacting with the plant list,” manual, 2014.

1010 [102] F. Jansen and J. Dengler, “Plant names in vegetation databases – a neglected source of
1011 bias,” *J. Veg. Sci.*, vol. 21 (6), pp. 1179–1186, 2010, doi: 10.1111/j.1654-1103.2010.01209.x.

1012 [103] R. Kindt, “WorldFlora: An R package for exact and fuzzy matching of plant names against the
1013 World Flora Online taxonomic backbone data,” *Appl. Plant Sci.*, vol. 8, p. e11388, 2020.

1014 [104] Z. S. L. Foster, T. J. Sharpton, and N. J. Grünwald, “Metacoder: An R package for
1015 visualization and manipulation of community taxonomic diversity data,” *PLOS Comput. Biol.*, vol. 13,
1016 no. 2, p. e1005404, Feb. 2017, doi: 10.1371/journal.pcbi.1005404.

1017 [105] M. Reginato, “monographaR: An R package to facilitate the production of plant taxonomic
1018 monographs,” *Brittonia*, vol. 68, no. 2, pp. 212–216, Jun. 2016, doi: 10.1007/s12228-015-9407-z.

1019 [106] S. Chamberlain, “rgnparser: Parse scientific names,” manual, 2021. [Online]. Available:
1020 <https://CRAN.R-project.org/package=rgnparser>

1021 [107] S. Chamberlain, “splister: Match species lists against reference list,” manual, 2021. [Online].
1022 Available: <https://github.com/sckott/splister>

1023 [108] S. Chamberlain, “taxspell: Spell check taxonomic names,” manual, 2021.

1024 [109] Florian Schneider, *EcologicalTraitData/traitdataform: Conforming v0.10 of the Ecological
1025 Trait-data Standard (ETS)*. Zenodo, 2020. doi: 10.5281/zenodo.4594227.

1026 [110] Miguel Alvarez, *kamapu/vegetable: Handling Vegetation Data Sets*. Zenodo, 2020. doi:
1027 10.5281/zenodo.3776780.

1028 [111] A. Bichat, “yatah: Yet another TAXonomy handler,” manual, 2020. [Online]. Available:
1029 <https://CRAN.R-project.org/package=yatah>

1030

1031 Supplementary information

- 1032 • Supplementary Information S1: Sources and tools food taxonomic information. Word
1033 description of databases and protocol to identify relevant R packages.
- 1034 • Supplementary Information S2: Excel table describing examined R packages and
1035 databases including links between packages and databases.

1036

1037 Supplementary Information S1 - Sources and tools for

1038 taxonomic information

1039 Description of databases

1040 Please refer to [dataset S2](#). Specifically the sheet “Databases” that describe the acronym of
1041 databases, their full names, their URLs and references. The sheet “DBs metadata” describes
1042 the content of each of the columns of the “Databases” sheet.

1043 Searching for R packages

1044 To identify the sets of packages to include in our review we searched for terms “taxon”,
1045 “taxa”, “taxonomy”, “taxonomic”, and “taxonomical” on CRAN through the **pkgsearch**
1046 package, as well as GitHub using its internal search function, and Bioconductor with the
1047 rdrr.io/find/ website.

1048 Our inclusion criteria for identified tools:

- 1049 1. The tool had to be an actual R package that could be installed (exclude collection of
1050 scripts).
- 1051 2. The tool had to be functional (exclude preliminary packages that were abandoned
1052 while not entirely developed).

1053 From this initial list of packages we manually identified packages that would be relevant as
1054 to wrangle taxonomic data by reading both the titles and the description of packages. We
1055 obtained a list of 67 packages to assess. We excluded 7 packages that were focused only
1056 on genomic information and genomic data wrangling.

1057 We classified packages into wide categories:

- 1058 ● **Infrastructure**, if the package provided basic R structure that could be used to
1059 further develop other packages on taxonomy.
- 1060 ● **Database Access (Online)**, for packages accessing taxonomic databases that need
1061 to be connected at all times to provide this information.
- 1062 ● **Database Access (Offline)**, for packages that either provide taxonomic information
1063 directly offline, or may access the information offline after an initial download.
- 1064 ● **Data Wrangling**, for packages that are able to manipulate taxonomic information
1065 (provide summary statistics, modify granularity, etc.).
- 1066 ● **Data Visualization**, for packages that create plots from taxonomic information.

1067 We detailed the functionalities of the packages if they were using online or offline resources.
1068 Whether they were to be used by end users or rather used by other package developers to
1069 build upon, and if they were actively maintained.

1070 Description of packages

1071 Please refer to [dataset S2](#). Specifically the sheet “Packages” that describe in full all the
1072 examined packages, including their URLs, their last date of update, if they allow for access
1073 to taxonomic databases, etc. The sheet “Pkgs metadata” describes the content of each of
1074 the columns of the “Packages” sheet.

1075

1076 **Supplementary Information S2 - table of description of tools**
1077 **and databases**

1078 A table that contains:

- 1079 ● a list of examined packages with more data columns (URLs, last date of update, link
1080 to which database it access)
- 1081 ● a list of databases with corresponding websites an references
- 1082 ● a list of links between databases with actual sources for these links.