



This is a repository copy of *A review of the heterogeneous landscape of biodiversity databases: Opportunities and challenges for a synthesized biodiversity knowledge base.*

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/187190/>

Version: Accepted Version

Article:

Feng, X, Enquist, BJ, Park, DS et al. (34 more authors) (2022) A review of the heterogeneous landscape of biodiversity databases: Opportunities and challenges for a synthesized biodiversity knowledge base. *Global Ecology and Biogeography*, 31 (7). pp. 1242-1260. ISSN 1466-822X

<https://doi.org/10.1111/geb.13497>

© 2022 John Wiley & Sons Ltd. This is the peer reviewed version of the following article: Feng, X., Enquist, B. J., Park, D. S., Boyle, B., Breshears, D. D., Gallagher, R. V., Lien, A., Newman, E. A., Burger, J. R., Maitner, B. S., Merow, C., Li, Y., Huynh, K. M., Ernst, K., Baldwin, E., Foden, W., Hannah, L., Jørgensen, P. M., Kraft, N. J. B., ... López-Hoffman, L. (2022). A review of the heterogeneous landscape of biodiversity databases: Opportunities and challenges for a synthesized biodiversity knowledge base. *Global Ecology and Biogeography*, 31(7), 1242– 1260., which has been published in final form at <https://doi.org/10.1111/geb.13497>. This article may be used for non-commercial purposes in accordance with Wiley Terms and Conditions for Use of Self-Archived Versions. This article may not be enhanced, enriched or otherwise transformed into a derivative work, without express permission from Wiley or by statutory rights under applicable legislation. Copyright notices must not be removed, obscured or modified. The article must be linked to Wiley's version of record on Wiley Online Library and any embedding, framing or otherwise making available the article or pages thereof by third parties from platforms, services and websites other than Wiley Online Library must be prohibited. Unless indicated otherwise, they may be downloaded and/or printed for private study, or other acts as permitted by national copyright laws. The publisher or other rights holders may allow further reproduction and re-use of the full text version. This is indicated by the licence information on the White Rose Research Online record for the item.

Takedown

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing eprints@whiterose.ac.uk including the URL of the record and the reason for the withdrawal request.

1 **Feng, X. Lovett, J.C. + 39 authors. 2022. A review of the heterogeneous landscape of**
2 **biodiversity databases: opportunities and challenges for a synthesized biodiversity**
3 **knowledge base. Global Ecology and Biogeography**

4
5
6 **Review**

7
8 **Title:** A review of the heterogeneous landscape of biodiversity databases: opportunities and
9 challenges for a synthesized biodiversity knowledge base

10
11 **Abstract**

12 **Aim:** Addressing global environmental challenges requires access to biodiversity data across
13 wide spatial, temporal and biological scales. Recent decades have witnessed an exponential
14 increase of biodiversity information aggregated by biodiversity databases (hereafter ‘databases’).
15 However, heterogeneous coverage, protocols, and standards of databases hampered the data
16 integration among databases. To stimulate the next stage of data integration, here we present a
17 synthesis of major databases, and investigate i) how the coverages of databases vary across
18 taxonomy, space, and record type; ii) the degree of integration among databases; iii) how
19 integration of databases can increase biodiversity knowledge; iv) the barriers to databases
20 integration.

21 **Location:** Global

22 **Time period:** Contemporary

23 **Major taxa studied:** Plants and Vertebrates

24 **Methods:** We reviewed the scope of twelve well-established databases and assessed the status
25 of their integration. We synthesized information from these databases to assess major knowledge
26 gaps and barriers to fully integration. We estimated how improved integration can increase the
27 coverage and depth of biodiversity knowledge.

28 **Results:** Each reviewed database had unique focus of data coverages. Data flows were common
29 among databases, though not always clearly documented. Functional trait databases were more
30 isolated than those pertaining to species distributions. Poor compatibility between taxonomic
31 systems used by different databases posed a major challenge to integration. We demonstrated
32 that integration of distribution databases can lead to greater taxonomic coverage that corresponds
33 to 23 years’ advancement in knowledge accumulation, and improvement in taxonomic coverage
34 could be as high as 22.4% for trait databases.

35 **Main conclusions:** Rapid increase of biodiversity knowledge can be achieved through the
36 integration of databases, providing the data necessary to address critical environmental
37 challenges. Our synthesis provides an overview of the integration status of databases. Full
38 integration across databases will require tackling the major impediments to data integration –
39 taxonomic incompatibility, lags in data exchange, barriers to effective data synchronization, and
40 isolation of individual initiatives.

41
42 **Keywords:** Big Data, Biodiversity Informatics, Biogeography, Database integration, Functional
43 trait, Taxonomic System

44

45 **1. Introduction**

46 In the face of rapid global changes, a grand challenge is how to efficiently catalogue, assess,
47 anticipate, and respond to changes in biodiversity and associated ecosystem services (Chapin *et al.*,
48 *et al.*, 2000; Ceballos *et al.*, 2015; Díaz *et al.*, 2019). Addressing this challenge requires
49 unprecedented access to biodiversity data across fine to broad spatial, temporal and biological
50 scales (Beck *et al.*, 2012). The past few decades have witnessed fast growth of biodiversity
51 information (Bisby, 2000; Hardisty *et al.*, 2013; Hobern *et al.*, 2019). Rapid digitization of
52 existing biodiversity collections and ongoing collection of new information are expanding data
53 availability worldwide (Sullivan *et al.*, 2014; Page *et al.*, 2015; Chandler *et al.*, 2017b). Indeed,
54 the Global Biodiversity Information Facility (GBIF) – the world’s leading repository of
55 biodiversity observations – recently reached 1.6 billion records (accessed March 2021).
56 However, we are still a long way from fully characterizing the taxonomy, geographic ranges and
57 functions of all species on Earth (Lomolino, 2004; Hortal *et al.*, 2015; Stork, 2018). Addressing
58 these shortfalls requires novel efforts in data synthesis to integrate the information held in the
59 world’s biodiversity projects, some 600+ of which had been created as of 2014 (Belbin, 2014)
60 and nearly half of which are essentially invisible or inaccessible to the research community due
61 to lack of cataloguing and integration (Blair *et al.*, 2020).

62
63 Data aggregation has been an ongoing goal of the biodiversity community (Nelson & Ellis,
64 2019), and a tremendous amount of work has been done by existing biodiversity data
65 aggregators, such as GBIF, iDigBio, and VertNet. However, the challenges are many: existing
66 biodiversity data aggregators often have singular objectives and consequently adhere to different
67 protocols and standards (Mesibov, 2018) (termed “data domains” in (König *et al.*, 2019)), and
68 datasets are highly heterogeneous spatially, temporally, and taxonomically (Reichman *et al.*,
69 2011; Cornwell *et al.*, 2019). The differences among biodiversity data aggregators can
70 accumulate over time; thus, biodiversity data aggregators run the risk of “speciating,” or
71 becoming isolated, which can impede data sharing and integration. In response, the community
72 has been calling for greater alignment between efforts and actively working on coordination
73 mechanisms for developing shared roadmaps for biodiversity informatics (Hobern *et al.*, 2019).
74 We therefore assert that a new synthesis is needed for the next stage of biodiversity data
75 integration, i.e., information from existing biodiversity data aggregators should be further
76 integrated to reduce shortfalls in biodiversity knowledge and achieve a more complete picture of
77 Earth’s biodiversity (Hobern *et al.*, 2019; König *et al.*, 2019; Kattge *et al.*, 2020).

78
79 To facilitate better integration among biodiversity data domains, we first need to assess the
80 current state of connectivity and integration among databases. Though biodiversity data
81 generally are well organized in individual databases, overlaps in their data coverage and the
82 extent of communication between databases remains unclear. Indeed, attention has rarely been
83 paid to the post-aggregation processes and interactions among commonly used databases (such
84 as nontransparent data flows between two databases) and synthesis studies of biodiversity data
85 from multiple databases are still scarce in the literature (Cornwell *et al.*, 2019; König *et al.*,
86 2019). To address this gap, we conducted a synthesis of existing biodiversity databases, and
87 aimed to answer four questions: **(i)** How does the coverage of a suite of major biodiversity
88 databases differ across taxon, space, and record type? **(ii)** How are existing biodiversity
89 databases integrated? **(iii)** How would the integration of databases increase biodiversity
90 knowledge? and **(iv)** What are the barriers that prevent data integration? To answer these

91 questions, we first reviewed the scope of existing major biodiversity databases and assessed the
92 status of their integration. We also demonstrated that the integration of biodiversity databases
93 could rapidly narrow major knowledge gaps. Finally, we discussed barriers that need to be
94 overcome to obtain a more complete picture of the biodiversity on Earth.

95

96 **2. Review of biodiversity databases**

97 Many biodiversity databases have been built over the past two decades, with varying emphases
98 on taxonomy, spatial location, and record type. To synthesize the major attributes of existing
99 biodiversity databases, we selected twelve well-established biodiversity databases: Atlas of
100 Living Australia (ALA; Belbin & Williams, 2016), Botanical Information and Ecology Network
101 (BIEN; Enquist *et al.*, 2016), Biodiversity Information Serving Our Nation (BISON; U.S.
102 Geological Survey, 2018), eBird (Sullivan *et al.*, 2014), Encyclopedia of Life (EOL; Parr *et al.*,
103 2014), Global Biodiversity Information Facility (GBIF), Global Inventory of Floras and Traits
104 (GIFT; Weigelt *et al.*, 2017), Integrated Digitized Biocollections (iDigBio, 2018a), iNaturalist
105 (iNaturalist), Map of Life (MOL; Jetz *et al.*, 2012), a global database of plant traits (TRY; Kattge
106 *et al.*, 2011), and VertNet (Constable *et al.*, 2010). Our selection can not cover every notable
107 database because of limited effort and the accessibility of database content or documentations,
108 though they were chosen to represent the breadth of the most commonly used, well-established
109 large-scale biodiversity databases (MacFadden & Guralnick, 2016; Chandler *et al.*, 2017a; James
110 *et al.*, 2018; Singer *et al.*, 2018; Cornwell *et al.*, 2019; König *et al.*, 2019) to maximize the
111 generalizability of our results and conclusions. We acknowledge that these databases are
112 typically under active development; thus our synthesis is based on a snapshot of their status on
113 the access date (March 2021; see Appendix 1).

114

115 **2.1 Varied focuses among biodiversity databases**

116 We reviewed associated metadata for biodiversity databases from project websites or
117 publications. We recorded database name, taxonomic scope, taxonomic system, record type,
118 number of records, and spatial coverage. We classified the record types into three categories:
119 geographic distribution, media type, and biological information (standardized trait databases or
120 generalized text descriptions). Within geographic distribution, we further classified the
121 information as specimen records, observations, checklists of geographic regions, or distribution
122 maps. Specimen records and observations both have information on specific occurrences of a
123 species at a georeferenced point location, but only specimen records are associated with physical
124 specimens. Checklists usually contain lists of species known to be present in defined geographic
125 regions (e.g., political divisions or protected areas). Distribution maps are those that were drawn
126 by experts or generated through models with various degrees of complexity. Media data type
127 were classified as image, audio, and video. Biological information included standardized trait
128 and generalized text descriptions.

129

130 Our review showed that each of these biodiversity databases holds unique scientific value
131 because they cover different spatial extents, taxonomic groups, and record types (Fig. 1a). The
132 databases could be grouped into different clusters based on similarities of focus and data
133 coverage. For example, EOL, iNaturalist, and eBird form a cluster of databases that indexes
134 media data and biological descriptions, while also sharing public education objective (Fig. 1b).
135 TRY and GIFT form another cluster that mainly focuses on indexing functional traits of plants.
136 GBIF, BISON, iDigBio, and VertNet form yet another cluster that emphasizes indexing species

137 occurrences. The cluster of ALA, MOL, and BIEN share the property of indexing both species
138 occurrences and geographic range maps. Here our grouping of databases considered the different
139 attributes equally, though assigning different weights on the attributes can lead to different
140 grouping outcomes. For example, many of the databases seek to document all taxa across the
141 globe (e.g., GBIF, EOL, eBird) or to index many types of data (e.g., EOL, ALA, iNaturalist).

142

143 **2.2 Data integration status among biodiversity databases**

144 To understand how existing biodiversity databases are integrated, we reviewed the data flow
145 among the databases. Biodiversity databases (e.g., GBIF) are typically data aggregators of
146 digitalized information from data providers, such as museums, herbariums, and research data
147 repositories, and detailed information about data providers are usually acknowledged on a
148 databases' website (e.g., BIEN data contributors-
149 [https://web.archive.org/web/20210511034441/https://bien.nceas.ucsb.edu/bien/data-](https://web.archive.org/web/20210511034441/https://bien.nceas.ucsb.edu/bien/data-contributors/)
150 [contributors/](https://web.archive.org/web/20210511034441/https://bien.nceas.ucsb.edu/bien/data-contributors/)). However, it is usually not straight forward to understand whether one database is
151 aggregated by another database, probably because of the concern of losing uniqueness of data
152 coverage, i.e. acknowledging to be aggregated by another aggregator can be interpreted as one
153 database becoming a subset of the other database. Regardless, understanding such relationships
154 among databases is important for users, as this immediately affects the determination of most
155 comprehensive data coverage (e.g., whether or not GBIF has the most complete occurrence set of
156 a species) or evaluation of data quality (e.g., whether or not to consider duplicated records when
157 using multiple databases). Therefore, we assessed data integration among biodiversity databases
158 based on their documentation and publications.

159

160 Overall, the data flows between biodiversity databases are not always clearly documented and at
161 times the relationships need to be inferred. Key technical details of data flow, such as time and
162 frequency of data exchange/flow, and the version or date of the imported data, are usually
163 lacking. The lack of 'snapshot' data archives hinders the reproduction of data content, as well as
164 the reproducibility of associated scientific research (Feng *et al.*, 2019). Unclear documentation of
165 data exchange may also lead to compliance issues with data licensing, and can prevent
166 assignment of proper credit to data collectors.

167

168 We found that data flow, unidirectional or bidirectional, is common among biodiversity
169 databases (Fig. 2 & Table S1). Among the network of databases, GBIF serves as a central
170 aggregator at a global scale that ingests species occurrence data from many databases, such as
171 BISON, iDigBio, and eBird. ALA and BISON have bidirectional data flows with GBIF – they
172 both i) aggregate biodiversity data collected from their focal regions (i.e., Australia and North
173 America respectively) and pass the data to GBIF, and ii) import other data collected from
174 Australia or North America from GBIF to their respective databases (Table S1). There are also
175 cases of unidirectional data flow from GBIF to specialized databases. For example, MOL
176 aggregates multiple types of information of species geographic distributions, including
177 occurrences from GBIF; as does BIEN.

178

179 We summarized the status of data integration across databases into four categories: synced,
180 lagged, impeded, and isolated (Fig. 3). Ideally, information in databases could be fully integrated
181 in either one or multiple directions in real (or near-real) time (i.e., *synced*). For example, data
182 published to iDigBio is automatically published to GBIF (iDigBio, 2018b; Singer *et al.*, 2018),

183 thus the content of iDigBio is considered synced with GBIF (Fig. 3). However, differences may
184 arise between otherwise fully integrated databases in the time between synchronization events
185 (*lagged*). For example, BIEN imports and integrates data from GBIF and other sources at annual
186 or longer intervals, which provides more stable and easily archived datasets, but the imported
187 GBIF content can be different from the most up-to-date GBIF data until the next
188 synchronization. This lag can be addressed by increasing the frequency of data exchange, shared
189 data import protocols, or developing novel database architecture designed for data integration
190 (LeBauer *et al.*, 2013). Differences between databases may also arise from obstacles that prevent
191 subsets of data from being shared (*impeded*). For example, iNaturalist only publishes data to
192 GBIF that are properly licensed (iNaturalist, 2018)). Differences in data licensing is one of the
193 major impediments to integration and is a problem that was rarely emphasized in biodiversity
194 data aggregation prior to the last decade. For example, GBIF initialized a license requirement in
195 2014 (GBIF, 2014) and excluded approximately 49 million existing records without appropriate
196 licenses. Clearly defined data licenses will make future data use and integration legally
197 straightforward, and will also provide a cornerstone for the Open Science movement (Escribano
198 *et al.*, 2018). Creative commons licenses are the most widely used mechanism to ensure proper
199 attribution while allowing others to copy and distribute data (Fitzgerald *et al.*, 2007).

200
201 Unlike the distribution databases discussed above, trait databases are characterized by isolation
202 status. These databases typically capture data within particular taxa or focus on a single trait,
203 such as GlobTherm for thermal tolerance (Bennett *et al.*, 2018) and AmphiBIO for amphibian
204 ecological traits (Oliveira *et al.*, 2017) (Fig. 3). A degree of isolation is unavoidable due to the
205 complex nature of trait data, which varies greatly in terms of data types, units, and measurement
206 methods (Deans *et al.*, 2015) and the taxon-specific nature of many traits (e.g., seed traits apply
207 only to seed plants). Such complexity is not resolved by following existing standard commonly
208 used by occurrence data such as Darwin Core (Wieczorek *et al.*, 2012). Effective synthesis and
209 integration of trait information will require trait-specific specifications such as trait ontologies
210 (Walls *et al.*, 2012), trait data standards (Schneider *et al.*, 2019) and embracing of Open Science
211 principles via initiatives like the Open Traits Network (Gallagher *et al.*, 2020).

212
213 Poor compatibility between taxonomic systems adopted by different databases has posed a major
214 impediment for database integration (Fig. 2 & Table S2). As biodiversity information is
215 generally indexed by species' scientific names, a crucial step is to index information based on
216 one unified or multiple compatible taxonomic systems. Taxonomic systems reflect decisions of
217 database developers; some databases maintain flexibility in nomenclature, especially when the
218 taxa are in flux (e.g., vertebrate species stored in VertNet), whereas some databases impose
219 stronger rules. For example, EOL maintains multiple independent taxonomic systems to avoid
220 potential conflicts between non-compatible nomenclature; GBIF and COL have both employed a
221 comprehensive but single-backbone system designed to be compatible with different taxonomic
222 systems; MOL developed a backbone that includes Catalogue of Life (a global effort to compile
223 existing catalogued species) and manually curated taxonomic datasets for synonym issues; BIEN
224 standardizes taxon names according to external, expert-curated taxonomic reference databases
225 (Boyle *et al.*, 2013). The different approaches and strategies to accommodating taxonomic
226 systems among biodiversity databases may solve taxonomic issues locally for that specific
227 database (Jorge & Peterson, 2004), but deepen differences that prevent future data integration,

228 thus facilitating the “speciation” of databases. Still, resolving differences between existing
 229 taxonomic systems is just an initial step. Creation of a single authoritative list of names will take
 230 time; full reconciliation of synonyms and distinct taxon concepts may take decades (Berendsohn,
 231 1997; Franz & Peet, 2009; Boyle *et al.*, 2013; Wiser, 2016; Garnett *et al.*, 2020). This will
 232 require a global effort, as envisioned by the Global Taxonomy Initiative (Samper, 2004).

233

234 **3. Enhanced data coverage via database integration**

235 To quantify the improvement of combining multiple databases, we compared leading databases
 236 that focus on similar taxonomic groups and similar record types. We used terrestrial plants
 237 (Embryophyta; hereafter “plants”) and vertebrates (Vertebrata) as test cases, because these
 238 taxonomic groups are comparatively well collected and documented in biodiversity databases
 239 compared to others (Clark & May, 2002; Fazey *et al.*, 2005; Hecnar, 2009; Titley *et al.*, 2017;
 240 Cornwell *et al.*, 2019; König *et al.*, 2019; Kattge *et al.*, 2020). We did not use taxon, such as
 241 microbes or invertebrates, that account for large portions of biodiversity on Earth but face huge
 242 data gaps (Locey & Lennon, 2016). Specifically, we combined (i) the distribution of terrestrial
 243 plants from GBIF and non-GBIF sources, and (ii) one crucial and commonly measured trait for
 244 plants and vertebrates, respectively: maximum height (Moles *et al.*, 2009; Guralnick *et al.*, 2016)
 245 using the Botanical Information and Ecology Network (BIEN (Enquist *et al.*, 2016)), TRY
 246 initiative (Kattge *et al.*, 2011), and EOL (Parr *et al.*, 2014), and body length using VertNet
 247 (Constable *et al.*, 2010) and EOL (see Appendix 1). Our study goes beyond recent gap analyses
 248 of biodiversity data (Meyer *et al.*, 2016; Cornwell *et al.*, 2019; König *et al.*, 2019), by expanding
 249 the scope to multiple data aggregators with similar missions, in two major clades (i.e., plants and
 250 vertebrates), and using an ecological trait characterized by continuous values.

251

252 **3.1 Better coverage through data integration**

253 **3.1.1 Overall trend in data collection**

254 We found that the total number of distribution records (spatial coordinates) for plants has
 255 increased exponentially since the 1750s (Lomolino *et al.*, 2010) (Fig. 4a) as documented in GBIF
 256 and the combined dataset. A similar exponential increase was found when only spatially unique
 257 records were examined (Fig. 4b). This pattern is also supported by a model selection analysis
 258 among linear, exponential, and logistic functions (Table S3). This trend in the growth of
 259 biodiversity data is analogous to many accelerating processes in the Anthropocene (Steffen *et al.*,
 260 2015), such as urbanization, globalization, transportation, and telecommunications. One
 261 prominent example in Information Technology (IT) is the exponential growth in the number of
 262 transistors in a dense integrated circuit, which doubles roughly every two years (Moore, 1965).
 263 This pattern, termed “Moore’s Law”, is also evident in the accelerating development of cyber
 264 infrastructures for many disciplines in science. Based on the similar exponential curve for
 265 biodiversity data, we estimated that the total number of plant distribution records doubles every
 266 17 years and the number of spatially unique records doubles every 21 years. The high speed of
 267 biodiversity data accumulation represents the great power of data collection, digitization,
 268 processing, and publishing, which lays the basis for and presents the opportunities for
 269 biodiversity database integration.

270

271 In contrast to the number of distribution records, the number of species identified is gradually
 272 reaching saturation (Fig. 4c). Based on a fitted logistic curve (Table S3), we predicted that the
 273 number of catalogued plant species in distribution databases would be saturated at $365,519 \pm$

274 2,233 (mean \pm SD of the coefficient from the fitted logistic model), i.e. the saturation point of
275 predicted number of terrestrial plant species in the integrated biodiversity distribution databases,
276 with species names resolved using the Taxonomic Name Resolution Service (TNRS; version 5.0)
277 (Boyle *et al.*, 2013). This estimate is higher than the current catalogued number of terrestrial
278 plants in Catalogue of Life (COL; 354,327), though within the previously estimated range for the
279 total number of plant species on Earth (334,000 - 403,911) (Lughadha *et al.*, 2016). The slowing
280 trend in plant species discovery started in ~1949 (the inflection point of the logistic curve of the
281 cumulative number of species in GBIF; Table S1), and is in line with previous estimations
282 (Christenhusz & Byng, 2016). Such trends may suggest that we are gradually reaching saturation
283 and closing the *Linnean shortfall*, the lack of knowledge in describing and cataloging species
284 (Hortal *et al.*, 2015), for plants. The slowing trend could also be caused by species extinctions,
285 reduced funding for natural history studies, and increasing difficulties in detecting the remaining
286 rare species (Joppa *et al.*, 2011).

287

288 **3.1.2 Improvement in distribution data**

289 Integration of biodiversity databases would powerfully increase our knowledge of biodiversity.
290 For instance, GBIF is the world's largest biodiversity repository, but adding ~15 million records
291 from additional sources (compiled by BIEN) would improve its coverage by ~3.7 million
292 spatially unique records and ~20 thousand species (Fig. 4d-f). The number of distribution records
293 per taxon in GBIF could be increased by 4.4% – an average of 19 additional records per species.
294 The improvement of taxonomic coverage in GBIF would be equivalent to 23 years of new data
295 accumulation, based on extrapolation of the fitted logistic curve (Fig. 4c, Table S3). GBIF and
296 non-GBIF datasets together provide distribution data for ~ 307,985 species (76-92% of the
297 estimated richness of all plants (Lughadha *et al.*, 2016)), suggesting we are gradually decreasing
298 the *Wallacean shortfall*, the lack of knowledge in species distribution, for plant species, in
299 accordance with findings in Cornwell *et al.* (2019).

300

301 **3.1.3 Improvement in trait data**

302 Database integration also substantially improves the taxonomic coverage of trait information
303 (i.e., maximum height in plants; body length in vertebrates; see Methods). Under standardized
304 taxonomy, we found that individual plant and vertebrate trait databases always include unique
305 species-trait combinations and cover different portions of taxonomic diversity (Fig. 5). For
306 instance, trait knowledge increased in 69-82 plant orders and 86-124 vertebrate orders through
307 database integration, while the range of increase varied by database. The average improvement
308 of species-trait combination across these databases ranged from 2.0 to 8.7% for plant orders and
309 21.5-22.4% for vertebrate orders. The number of plant orders that were sparsely-sampled in
310 BIEN (i.e., <10% of species with trait observations), for example, decreased from 99 to 65
311 through data integration; a similar decrease was seen for sparsely-sampled vertebrate orders in
312 EOL from 53 down to nine (Fig. 5).

313

314 **3.1.4 Limitations of our assessment**

315 Data integration can effectively decrease the gaps in our knowledge, and the resulting more
316 comprehensive data can facilitate global scale studies of biodiversity and help identify and
317 reduce potential data biases (Reddy & Dávalos, 2003). We note that our assessment of the
318 possibilities for data integration does not address how different data sources (or “data
319 resolutions,” as defined in (König *et al.*, 2019)) should be best integrated for different study

objectives. These mismatches are apparent in cases, such as distribution data represented by presences vs. abundances, or a trait value measured at individual level vs. species level. However, indexing the availability of trait data for a focal species is a major step toward more rigorous data integration and scientific research. With the integrated data, one could cross-validate the values from different sources to ask questions such as: “Do trait values vary by methods of measurements?” or “Can species-level trait data well represent the range of values measured at the individual level?” Cross-validations will be especially useful if the user of one database is mainly the general public while the user of the other is the science community, so that more rigorous information is delivered from the science community to the general public. With the integrated data, one could also conduct scientific research at broader scales and study, for example, trait variation across time or across spatial or environmental gradients (Siefert *et al.*, 2015), or species-trait combinations within communities.

3.2 A clearer picture of what we do not know

Importantly, database integration can provide an improved assessment of gaps in biodiversity knowledge (Meyer *et al.*, 2015; Cornwell *et al.*, 2019; König *et al.*, 2019). Following our integration of various databases (Appendix 1), approximately 58,000 plant species still lacked publicly available distribution records. This gap corresponds to approximately 15.8% of the species in Catalogue of Life – a global effort to compile existing catalogued species. The coverage of distribution records in plant orders varied from 47% (in order Hypnales) to fully covered in some orders with small number of extant species (Cornwell *et al.*, 2019) (e.g. Ceratophyllales). Further, 30.8 million km² of ice-free land surface, as assessed using Eckert IV equal area projection, currently has no valid plant geolocations (Fig. 4g). These areas are mainly in Russia (despite the considerable recent progress of data sharing by the Russian GBIF community (Shashkov & Ivanova, 2019)), central Asia, and northern Africa, and are approximately 13% of the Earth’s land area.

Trait data have considerably larger gaps: height information is absent for 333,597 plant species from 102 orders from BIEN, TRY and EOL, and body length information is absent for 38,992 vertebrate species from 127 orders from VertNet and EOL. In total, height data is unavailable for approximately 92.6% of plant species and body length for 56.8% of vertebrate species in Catalogue of Life. The data coverages were mostly below 60% for plant orders and percentages were relatively higher for vertebrate orders. Plant height and vertebrate body length are commonly used traits in ecological research that are frequently recorded in databases (Moles *et al.*, 2009; Guralnick *et al.*, 2016), suggesting other biological traits (e.g., life span, metabolic rate) or essential biodiversity variables (e.g., population abundances) (Pereira *et al.*, 2013) will likely have much larger *shortfalls* (but see analyses of plant growth form in (König *et al.*, 2019)). In the face of accelerating increases in biodiversity data availability, recognizing the remaining knowledge gaps could help guide future data compilation efforts (e.g. the gap filling activity in eBird (eBird, 2014)) and potentially turn our enhanced power of compiling information into efforts that generate critically needed knowledge (Cornwell *et al.*, 2019).

4. Challenges and Opportunities

4.1 A catalogue and synthesis of biodiversity databases

To achieve global integration of biodiversity knowledge, we would first need to know what databases are available. To facilitate this process, we need a catalogue of biodiversity databases

366 with their metadata recorded, such as spatial, temporal, taxonomic scope, as well as the types of
367 data aggregated, so that existing or new databases can be easily known, compared, and
368 effectively used. Lee Belbin has maintained the Biodiversity Information Projects of the World
369 (Belbin, 2014) – essentially containing metadata of 685 biodiversity projects. The recorded
370 metadata includes project summary, geographic, temporal, and taxonomic scope, and key
371 technique attributes (though this list is no longer accessible after 2019; but see (Blair *et al.*,
372 2020)). Similarly, GBIF has a registry system that indexes the metadata of GBIF participants,
373 institutions, and datasets; however, data associated with this registry mainly focuses on a few
374 record types, including occurrences, checklists, and sampling events
375 ([https://web.archive.org/web/20210514141441/https://www.gbif.org/article/5F1XBKbirSiq0ascK](https://web.archive.org/web/20210514141441/https://www.gbif.org/article/5F1XBKbirSiq0ascKYiA8q/gbif-infrastructure-registry)
376 [YiA8q/gbif-infrastructure-registry](https://web.archive.org/web/20210514141441/https://www.gbif.org/article/5F1XBKbirSiq0ascKYiA8q/gbif-infrastructure-registry)). Another example is Global Index of Vegetation Plot
377 Databases that indexes the metadata of vegetation-plot data that are publicly available (Dengler
378 *et al.*, 2011). In contrast, DataONE has a broader scope that indexes the metadata of large variety
379 of biological and environmental data (Michener *et al.*, 2012). Those existing efforts form a good
380 basis for a catalogue of biodiversity databases that can continuously keep track of existing data
381 aggregators and index new aggregation efforts. Still, the relationships among the biodiversity
382 databases are not always obvious. Therefore, a synthesis, ideally updated regularly, would be
383 helpful to clarify the relationships among the biodiversity databases, in particular what is the
384 unique data coverage of one database and what are the data flows among biodiversity databases.

385

386 **4.2 Overcoming the barriers to database integration**

387 After cataloguing the metadata and synthesizing the relationships among biodiversity databases,
388 many technical barriers remain. As a prerequisite to integration, the data in a database should be
389 openly available with proper data licenses to minimize impediments to data sharing (see section
390 2.2); another major barrier is the incompatible taxonomic systems. A promising effort is
391 Catalogue of Life Plus (Banki *et al.*, 2019) that builds upon existing but disconnected efforts
392 (such as the COL and GBIF backbone taxonomy) to create an open, shared and sustainable
393 consensus taxonomy, which can serve as the infrastructure for individual biodiversity databases
394 or database integration. Thirdly, existing databases adopt different mechanisms of data standards
395 and database architecture (Hardisty *et al.*, 2019), thus leading to incompatibilities for database
396 integration. For example, during the data cleaning stage, one collection of a specimen without
397 coordinates could be georeferenced differently based on different georeferencing algorithms,
398 thus likely leading to two different coordinates, and therefore appear to be two different records
399 after data integration. One solution could be creating a community-wide standard and tools for
400 data evaluation and cleaning (e.g. Belbin *et al.*, 2018; Serra-Diaz *et al.*, 2018). Community-
401 driven standards for biodiversity data, such as Darwin Core (Wieczorek *et al.*, 2012), Humboldt
402 Core (Guralnick *et al.*, 2018), and trait-data standard (Schneider *et al.*, 2019) have emerged;
403 expanding the use of those community-developed data standards by individual databases would
404 enable more effective database integration. Overall, the essential goal is to maximize
405 compatibility, and thus minimize barriers to data flow and synthesis. After solving the technical
406 barriers, the integrated content from multiple databases could be organized in multiple non-
407 exclusive ways: i) a single centralized database, ii) some decentralized but connected databases
408 (Gallagher *et al.*, 2020), or iii) multiple synced databases (LeBauer *et al.*, 2013).

409

410 **4.3 Challenges for individual aggregators after database integration**

411 It is also worth thinking the uniqueness and destiny of individual databases after integration.
412 Seemingly, integration may render individual databases irrelevant, e.g., an individual database
413 may be considered a subset of an integrated database. However, this should not be the case. While
414 data integration occurs at shared data element (e.g., taxon, place, time) and data standard, each
415 individual database could still have unique domain information. For example, while GBIF
416 aggregates species occurrence data from iNaturalist, the latter still uniquely host the media data.
417 Also, an individual database can make a unique contribution by aiming to fill data gaps (e.g.,
418 spatial or taxon gaps revealed by the integrated knowledge base).

419
420 On the other side, there has been a process of specialization of databases along the whole
421 workflow of data aggregation. Specifically, the developers of some databases have expanded
422 their scope to development of infrastructure, such as tools for data integration, data cleaning, and
423 hosting data portals. There are prominent examples among the databases that have close
424 relationships with GBIF. For example, ALA develops open-access modules for the platform that
425 can be implemented by other biodiversity initiatives (Belbin *et al.*, 2021). VertNet has been
426 actively providing data maintenance services, including data cleaning and indexing, among the
427 network of collaborative biodiversity databases (Constable *et al.*, 2010).

428
429 Besides specialized roles in data aggregation or tool development, individual databases can also
430 play unique roles for users, even when based on the same shared knowledge base. For example,
431 ALA is prominent in the education of Australian biodiversity to its Australian users, as well as in
432 facilitating scientific research by putting this biodiversity in the context of its environment.

433
434

435 **5. Concluding remarks**

436 The accelerating increase of biodiversity data offers numerous exciting prospects and challenges
437 for documenting and forecasting the location, status, function and potential fate of species on the
438 planet. However, increases in biodiversity data do not directly translate to similar increases in the
439 knowledge needed to address many fundamental and applied questions. In the face of urgent
440 environmental challenges, new approaches are urgently needed to increase biodiversity
441 knowledge and accessibility of the knowledge. We demonstrate that rapid progress can be made
442 toward better biodiversity knowledge through the integration of database infrastructures.
443 Integration can lead to large and rapid increases in knowledge of species distributions and traits
444 (see (Conde *et al.*, 2019; König *et al.*, 2019)), but the benefit goes beyond just more complete
445 knowledge: it can reduce biases and doubled efforts in biodiversity research, allow cross-
446 validations to compare conclusions drawn from different sources, and provide a clearer picture of
447 where gaps remain, thereby helping to focus future sampling and research (König *et al.*, 2019).
448 To address the shortfalls in biodiversity knowledge and achieve full integration across databases,
449 we need to fund and maintain the foundations of biodiversity information science including
450 biological surveys, taxonomic assessment (Australian Academy of Science, 2018), and
451 digitization of legacy data (Ariño, 2010), as well as tackle the major impediments to data
452 integration – taxonomic incompatibility, lags in data exchange, barriers to effective synthesis,
453 and isolation of individual initiatives.

454 **References**

- 455 Ariño, A.H. (2010) Approaches to estimating the universe of natural history collections data.
456 *Biodiversity Informatics*, **7**, 81-92.
- 457 Australian Academy of Science (2018) Discovering Biodiversity: a decadal plan for taxonomy
458 and biosystematics in Australia and New Zealand 2018–2027. In:
- 459 Banki, O., Hobern, D., Döring, M. & Remsen, D. (2019) Catalogue of Life Plus: A collaborative
460 project to complete the checklist of the world's species. *Biodiversity Information Science
461 and Standards*, **3**, e37652.
- 462 Beck, J., Ballesteros-Mejia, L., Buchmann, C.M., Dengler, J., Fritz, S.A., Gruber, B., Hof, C.,
463 Jansen, F., Knapp, S., Kreft, H., Schneider, A.-K., Winter, M. & Dormann, C.F. (2012)
464 What's on the horizon for macroecology? *Ecography*, **35**, 673-683.
- 465 Belbin, L. (2014) *Biodiversity Information Projects of the World*. Retrieved from:
466 <https://web.archive.org/web/20180609082447/http://www.tdwg.org/biodiv-projects/>
467 (accessed 1 May 2018).
- 468 Belbin, L. & Williams, K.J. (2016) Towards a national bio-environmental data facility:
469 experiences from the Atlas of Living Australia. *International Journal of Geographical
470 Information Science*, **30**, 108-125.
- 471 Belbin, L., Wallis, E., Hobern, D. & Zenger, A. (2021) The Atlas of Living Australia: History,
472 current state and future directions. *Biodiversity data journal*, **9**, e65023-e65023.
- 473 Belbin, L., Chapman, A., Wiczorek, J., Zermoglio, P., Thompson, A. & Morris, P. (2018) Data
474 Quality Task Group 2: Tests and Assertions. *Biodiversity Information Science and
475 Standards*, **2**, e25608.
- 476 Bennett, J.M., Calosi, P., Clusella-Trullas, S., et al. (2018) GlobTherm, a global database on
477 thermal tolerances for aquatic and terrestrial organisms. *Scientific Data*, **5**, 180022-
478 180022.
- 479 Berendsohn, W.G. (1997) A Taxonomic Information Model for Botanical Databases: The IOPI
480 Model. *Taxon*, **46**, 283-309.
- 481 Bisby, F.A. (2000) The Quiet Revolution: Biodiversity Informatics and the Internet. *Science*,
482 **289**, 2309.
- 483 Blair, J., Gwiazdowski, R., Borrelli, A., Hotchkiss, M., Park, C., Perrett, G. & Hanner, R. (2020)
484 Towards a catalogue of biodiversity databases: An ontological case study. *Biodiversity
485 Data Journal*, **8**, e32765.
- 486 Boyle, B., Hopkins, N., Lu, Z., Raygoza Garay, J.A., Mozzherin, D., Rees, T., Matasci, N.,
487 Narro, M.L., Piel, W.H., McKay, S.J., Lowry, S., Freeland, C., Peet, R.K. & Enquist, B.J.
488 (2013) The taxonomic name resolution service: an online tool for automated
489 standardization of plant names. *BMC Bioinformatics*, **14**, 16.
- 490 Catalogue of Life (2021) *Species 2000 & ITIS Catalogue of Life, 2021-04-05. Digital resource at*
491 www.catalogueoflife.org. *Species 2000: Naturalis, Leiden, the Netherlands. ISSN 2405-*
492 *8858.*
- 493 Ceballos, G., Ehrlich, P.R., Barnosky, A.D., Garcia, A., Pringle, R.M. & Palmer, T.M. (2015)
494 Accelerated modern human-induced species losses: Entering the sixth mass extinction.
495 *Science Advances*, **1**, e1400253.
- 496 Chamberlain, S.A. & Szocs, E. (2013) taxize: taxonomic search and retrieval in R. *F1000Res*, **2**,
497 191.
- 498 Chandler, M., See, L., Copas, K., Bonde, A.M.Z., López, B.C., Danielsen, F., Legind, J.K.,
499 Masinde, S., Miller-Rushing, A.J., Newman, G., Rosemartin, A. & Turak, E. (2017a)

- 500 Contribution of citizen science towards international biodiversity monitoring. *Biological*
501 *Conservation*, **213**, 280-294.
- 502 Chandler, M., See, L., Copas, K., Bonde, A.M.Z., Lopez, B.C., Danielsen, F., Legind, J.K.,
503 Masinde, S., Miller-Rushing, A.J., Newman, G., Rosemartin, A. & Turak, E. (2017b)
504 Contribution of citizen science towards international biodiversity monitoring. *Biological*
505 *Conservation*, **213**, 280-294.
- 506 Chapin, F.S., 3rd, Zavaleta, E.S., Eviner, V.T., Naylor, R.L., Vitousek, P.M., Reynolds, H.L.,
507 Hooper, D.U., Lavorel, S., Sala, O.E., Hobbie, S.E., Mack, M.C. & Díaz, S. (2000)
508 Consequences of changing biodiversity. *Nature*, **405**, 234-242.
- 509 Christenhusz, M.J.M. & Byng, J.W. (2016) The number of known plants species in the world
510 and its annual increase. *Phytotaxa*, **261**, 201-217.
- 511 Clark, J.A. & May, R.M. (2002) Taxonomic Bias in Conservation Research. *Science*, **297**, 191.
- 512 Conde, D.A., Staerk, J., Colchero, F., et al. (2019) Data gaps and opportunities for comparative
513 and conservation biology. *Proceedings of the National Academy of Sciences*, **116**, 9658.
- 514 Constable, H., Guralnick, R., Wieczorek, J., Spencer, C., Peterson, A.T. & VertNet Steering, C.
515 (2010) VertNet: a new model for biodiversity data sharing. *PLoS Biology*, **8**, e1000309.
- 516 Cornwell, W.K., Pearse, W.D., Dalrymple, R.L. & Zanne, A.E. (2019) What we (don't) know
517 about global plant diversity. *Ecography*, **0**
- 518 Deans, A.R., Lewis, S.E., Huala, E., et al. (2015) Finding Our Way through Phenotypes. *PLoS*
519 *Biology*, **13**, e1002033.
- 520 Dengler, J., Jansen, F., Glöckler, F., Peet, R.K., De Cáceres, M., Chytrý, M., Ewald, J., Oldeland,
521 J., Lopez-Gonzalez, G., Finckh, M., Mucina, L., Rodwell, J.S., Schaminée, J.H.J. &
522 Spencer, N. (2011) The Global Index of Vegetation-Plot Databases (GIVD): a new
523 resource for vegetation science. *Journal of Vegetation Science*, **22**, 582-597.
- 524 Díaz, S., Settele, J., Brondízio, E., et al. (2019) Summary for policymakers of the global
525 assessment report on biodiversity and ecosystem services of the Intergovernmental
526 Science-Policy Platform on Biodiversity and Ecosystem Services. In:
527 eBird (2014) *eBird's missing species*. Retrieved from: [https://ebird.org/news/ebirds-missing-](https://ebird.org/news/ebirds-missing-species/)
528 [species/](https://ebird.org/news/ebirds-missing-species/) (accessed 1 January 2020).
- 529 Enquist, B.J., Condit, R., Peet, R.K., Schildhauer, M. & Thiers, B.M. (2016) Cyberinfrastructure
530 for an integrated botanical information network to investigate the ecological impacts of
531 global climate change on plant biodiversity. *PeerJ Preprints*, **4**, e2615v2.
- 532 Enquist, B.J., Feng, X., Donoghue, J.C.I., et al. The commonness of rarity: global distribution
533 across the land plants. In prep.
- 534 Escribano, N., Galicia, D. & Ariño, A.H. (2018) The tragedy of the biodiversity data commons: a
535 data impediment creeping nigher? *Database*, **2018**
- 536 Fazey, I., Fischer, J. & Lindenmayer, D.B. (2005) What do conservation biologists publish?
537 *Biological Conservation*, **124**, 63-73.
- 538 Feng, X., Park, D.S., Walker, C., Peterson, A.T., Merow, C. & Papeş, M. (2019) A checklist for
539 maximizing reproducibility of ecological niche models. *Nature Ecology & Evolution*, **3**,
540 1382-1395.
- 541 Fitzgerald, B.F., Coates, J.M. & Lewis, S.M. (2007) *Open Content Licensing: Cultivating the*
542 *Creative Commons*. Sydney University Press, Sydney, Australia.
- 543 Franz, N.M. & Peet, R.K. (2009) Towards a language for mapping relationships among
544 taxonomic concepts. *Systematics and Biodiversity*, **7**, 5-20.

- 545 Gallagher, R.V., Falster, D.S., Maitner, B.S., et al. (2020) Open Science principles for
 546 accelerating trait-based science across the Tree of Life. *Nature Ecology & Evolution*, **4**,
 547 294-303.
- 548 Garnett, S.T., Christidis, L., Conix, S., et al. (2020) Principles for creating a single authoritative
 549 list of the world's species. *PLOS Biology*, **18**, e3000736.
- 550 GBIF (2014) *New approaches to data licensing and endorsement*. Retrieved from:
 551 <https://www.gbif.org/news/82363/new-approaches-to-data-licensing-and-endorsement>
 552 (accessed 1 May 2018).
- 553 Guralnick, R., Walls, R. & Jetz, W. (2018) Humboldt Core – toward a standardized capture of
 554 biological inventories for biodiversity monitoring, modeling and assessment. *Ecography*,
 555 **41**, 713-725.
- 556 Guralnick, R.P., Zermoglio, P.F., Wiczorek, J., LaFrance, R., Bloom, D. & Russell, L. (2016)
 557 The importance of digitized biocollections as a source of trait data and a new VertNet
 558 resource. *Database*, **2016**, baw158-baw158.
- 559 Hardisty, A., Roberts, D. & The Biodiversity Informatics, C. (2013) A decadal view of
 560 biodiversity informatics: challenges and priorities. *BMC Ecology*, **13**, 16.
- 561 Hardisty, A.R., Belbin, L., Hobern, D., McGeoch, M.A., Pirzl, R., Williams, K.J. & Kissling,
 562 W.D. (2019) Research infrastructure challenges in preparing essential biodiversity
 563 variables data products for alien invasive species. *Environmental Research Letters*, **14**,
 564 025005.
- 565 Hecnar, S.J. (2009) Human bias and the biodiversity knowledge base: An examination of the
 566 published literature on vertebrates. *Biodiversity*, **10**, 18-24.
- 567 Hobern, D., Baptiste, B., Copas, K., et al. (2019) Connecting data and expertise: a new alliance
 568 for biodiversity knowledge. *Biodiversity data journal*, **7**, e33679-e33679.
- 569 Hortal, J., de Bello, F., Diniz-Filho, J.A.F., Lewinsohn, T.M., Lobo, J.M. & Ladle, R.J. (2015)
 570 Seven shortfalls that beset large-scale knowledge of biodiversity. *Annual Review of*
 571 *Ecology, Evolution, and Systematics*, **46**, 523-549.
- 572 iDigBio (2018a) *Integrated Digitized Biocollections (iDigBio)*. Retrieved from:
 573 <https://www.idigbio.org> (accessed 1 May 2018).
- 574 iDigBio (2018b) *Data Ingestion Guidance*. Retrieved from:
 575 https://www.idigbio.org/wiki/index.php/Data_Ingestion_Guidance (accessed).
- 576 iNaturalist Retrieved from: <https://www.inaturalist.org/> (accessed 1 May 2018).
- 577 iNaturalist (2018) *Research Grade Observations*. Retrieved from:
 578 <https://www.inaturalist.org/posts/16429-research-grade-observations> (accessed 20 October
 579 2018).
- 580 James, S.A., Soltis, P.S., Belbin, L., Chapman, A.D., Nelson, G., Paul, D.L. & Collins, M.
 581 (2018) Herbarium data: Global biodiversity and societal botanical needs for novel
 582 research. *Applications in Plant Sciences*, **6**, e1024.
- 583 Jetz, W., McPherson, J.M. & Guralnick, R.P. (2012) Integrating biodiversity distribution
 584 knowledge: toward a global map of life. *Trends in Ecology and Evolution*, **27**, 151-159.
- 585 Joppa, L.N., Roberts, D.L. & Pimm, S.L. (2011) How many species of flowering plants are
 586 there? *Proceedings of the Royal Society B: Biological Sciences*, **278**, 554-559.
- 587 Jorge, S. & Peterson, A.T. (2004) Biodiversity Informatics: Managing and Applying Primary
 588 Biodiversity Data. *Philosophical Transactions: Biological Sciences*, **359**, 689-698.
- 589 Kattge, J., Díaz, S., Lavorel, S., et al. (2011) TRY - a global database of plant traits. *Global*
 590 *Change Biology*, **17**, 2905-2935.

- 591 Kattge, J., Bönisch, G., Díaz, S., et al. (2020) TRY plant trait database – enhanced coverage and
592 open access. *Global Change Biology*, **26**, 119-188.
- 593 König, C., Weigelt, P., Schrader, J., Taylor, A., Kattge, J. & Kreft, H. (2019) Biodiversity data
594 integration—the significance of data resolution and domain. *PLoS Biology*, **17**,
595 e3000183.
- 596 LeBauer, D.S., Wang, D., Richter, K.T., Davidson, C.C. & Dietze, M.C. (2013) Facilitating
597 feedbacks between field measurements and ecosystem models. *Ecological Monographs*,
598 **83**, 133-154.
- 599 Locey, K.J. & Lennon, J.T. (2016) Scaling laws predict global microbial diversity. *Proceedings*
600 *of the National Academy of Sciences*, **113**, 5970.
- 601 Lomolino, M.V. (2004) Conservation biogeography. *Frontiers of biogeography: new directions*
602 *in the geography of nature*, 293-296.
- 603 Lomolino, M.V., Riddle, B.R., Whittaker, R.J. & Brown, J.H. (2010) *Biogeography*, 4th edn.
604 Sinauer Associates, Sunderland, Massachusetts.
- 605 Lughadha, E.N., Govaerts, R., Belyaeva, I., Black, N., Lindon, H., Allkin, R., Magill, R.E. &
606 Nicolson, N. (2016) Counting counts: revised estimates of numbers of accepted species
607 of flowering plants, seed plants, vascular plants and land plants with a review of other
608 recent estimates. *Phytotaxa*, **272**, 82-88.
- 609 MacFadden, B.J. & Guralnick, R.P. (2016) Horses in the Cloud: big data exploration and mining
610 of fossil and extant Equus (Mammalia: Equidae). *Paleobiology*, **43**, 1-14.
- 611 Mesibov, R. (2018) An audit of some processing effects in aggregated occurrence records.
612 *ZooKeys*, **751**, 129-146.
- 613 Meyer, C., Weigelt, P. & Kreft, H. (2016) Multidimensional biases, gaps and uncertainties in
614 global plant occurrence information. *Ecology Letters*, **19**, 992-1006.
- 615 Meyer, C., Kreft, H., Guralnick, R. & Jetz, W. (2015) Global priorities for an effective
616 information basis of biodiversity distributions. *Nature Communications*, **6**, 8221-8221.
- 617 Michener, W.K., Allard, S., Budden, A., Cook, R.B., Douglass, K., Frame, M., Kelling, S.,
618 Koskela, R., Tenopir, C. & Vieglais, D.A. (2012) Participatory design of DataONE—
619 Enabling cyberinfrastructure for the biological and environmental sciences. *Ecological*
620 *Informatics*, **11**, 5-15.
- 621 Moles, A.T., Warton, D.I., Warman, L., Swenson, N.G., Laffan, S.W., Zanne, A.E., Pitman, A.,
622 Hemmings, F.A. & Leishman, M.R. (2009) Global patterns in plant height. *Journal of*
623 *Ecology*, **97**, 923-932.
- 624 Moore, G.E. (1965) Cramming more components onto integrated circuits. *Electronics*, **38**, 114-
625 117.
- 626 Nelson, G. & Ellis, S. (2019) The history and impact of digitization and digital data mobilization
627 on biodiversity research. *Philosophical Transactions of the Royal Society B: Biological*
628 *Sciences*, **374**, 20170391.
- 629 Oliveira-Filho, A.T. (2017) *NeoTropTree, Flora arbórea da Região Neotropical: Um banco de*
630 *dados envolvendo biogeografia, diversidade e conservação*. Retrieved from:
631 <http://www.neotropree.info> (accessed 7 May 2019).
- 632 Oliveira, B.F., São-Pedro, V.A., Santos-Barrera, G., Penone, C. & Costa, G.C. (2017)
633 AmphiBIO, a global database for amphibian ecological traits. *Sci Data*, **4**, 170123.
- 634 Page, L.M., MacFadden, B.J., Fortes, J.A., Soltis, P.S. & Riccardi, G. (2015) Digitization of
635 biodiversity collections reveals biggest data on biodiversity. *BioScience*, **65**, 841-842.

- 636 Parr, C.S., Wilson, N., Leary, P., Schulz, K.S., Lans, K., Walley, L., Hammock, J.A., Goddard,
637 A., Rice, J., Studer, M., Holmes, J.T.G. & Corrigan, R.J., Jr. (2014) The Encyclopedia of
638 Life v2: Providing Global Access to Knowledge About Life on Earth. *Biodivers Data J*,
639 e1079.
- 640 Pereira, H.M., Ferrier, S., Walters, M., et al. (2013) Essential Biodiversity Variables. *Science*,
641 **339**, 277.
- 642 Reddy, S. & Dávalos, L.M. (2003) Geographical sampling bias and its implications for
643 conservation priorities in Africa. *Journal of Biogeography*, **30**, 1719-1727.
- 644 Reichman, O.J., Jones, M.B. & Schildhauer, M.P. (2011) Challenges and opportunities of open
645 data in ecology. *Science*, **331**, 703-705.
- 646 Samper, C. (2004) Taxonomy and environmental policy. *Philosophical Transactions of the*
647 *Royal Society of London. Series B: Biological Sciences*, **359**, 721-728.
- 648 Schneider, F.D., Fichtmueller, D., Gossner, M.M., Güntsch, A., Jochum, M., König-Ries, B., Le
649 Provost, G., Manning, P., Ostrowski, A., Penone, C. & Simons, N.K. (2019) Towards an
650 ecological trait-data standard. *Methods in Ecology and Evolution*, **10**, 2006-2019.
- 651 Serra-Diaz, J.M., Enquist, B.J., Maitner, B., Merow, C. & Svenning, J.-C. (2018) Big data of tree
652 species distributions: how big and how good? *Forest Ecosystems*, **4**, 30.
- 653 Shashkov, M. & Ivanova, N. (2019) Considerable Progress in Russian GBIF Community.
654 *Biodiversity Information Science and Standards*, **3**, e37015.
- 655 Siefert, A., Violle, C., Chalmandrier, L., et al. (2015) A global meta-analysis of the relative
656 extent of intraspecific trait variation in plant communities. *Ecology Letters*, **18**, 1406-
657 1419.
- 658 Singer, R.A., Love, K.J. & Page, L.M. (2018) A survey of digitized data from U.S. fish
659 collections in the iDigBio data aggregator. *PLOS ONE*, **13**, e0207636.
- 660 Steffen, W., Broadgate, W., Deutsch, L., Gaffney, O. & Ludwig, C. (2015) The trajectory of the
661 Anthropocene: The Great Acceleration. *The Anthropocene Review*, **2**, 81-98.
- 662 Stork, N.E. (2018) How many species of insects and other terrestrial arthropods are there on
663 Earth? *Annual Review of Entomology*, **63**, 31-45.
- 664 Sullivan, B.L., Aycrigg, J.L., Barry, J.H., et al. (2014) The eBird enterprise: An integrated
665 approach to development and application of citizen science. *Biological Conservation*,
666 **169**, 31-40.
- 667 Titley, M.A., Snaddon, J.L. & Turner, E.C. (2017) Scientific research on animal biodiversity is
668 systematically biased towards vertebrates and temperate regions. *PLOS ONE*, **12**,
669 e0189577.
- 670 U.S. Department of Agriculture Forest Service *Forest Inventory and Analysis Database*.
671 Retrieved from: <https://www.fia.fs.fed.us/> (accessed 1 May 2018).
- 672 U.S. Geological Survey (2018) *Biodiversity Information Serving Our Nation (BISON)*. Retrieved
673 from: <https://bison.usgs.gov> (accessed 1 May 2018).
- 674 Walls, R.L., Athreya, B., Cooper, L., Elser, J., Gandolfo, M.A., Jaiswal, P., Mungall, C.J.,
675 Preece, J., Rensing, S., Smith, B. & Stevenson, D.W. (2012) Ontologies as integrative
676 tools for plant science. *American journal of botany*, **99**, 1263-1275.
- 677 Weigelt, P., König, C. & Kreft, H. (2017) GIFT - a global inventory of Floras and traits for
678 macroecology and biogeography. *bioRxiv*, 535005.
- 679 Wieczorek, J., Bloom, D., Guralnick, R., Blum, S., Döring, M., Giovanni, R., Robertson, T. &
680 Vieglais, D. (2012) Darwin Core: An Evolving Community-Developed Biodiversity Data
681 Standard. *PLoS ONE*, **7**, e29715.

- 682 Wisser, S.K. (2016) Achievements and challenges in the integration, reuse and synthesis of
683 vegetation plot data. *Journal of Vegetation Science*, **27**, 868-879.
- 684 Zermoglio, P.F., Guralnick, R.P. & Wieczorek, J.R. (2016) A Standardized Reference Data Set
685 for Vertebrate Taxon Name Resolution. *PLOS ONE*, **11**, e0146894.

686

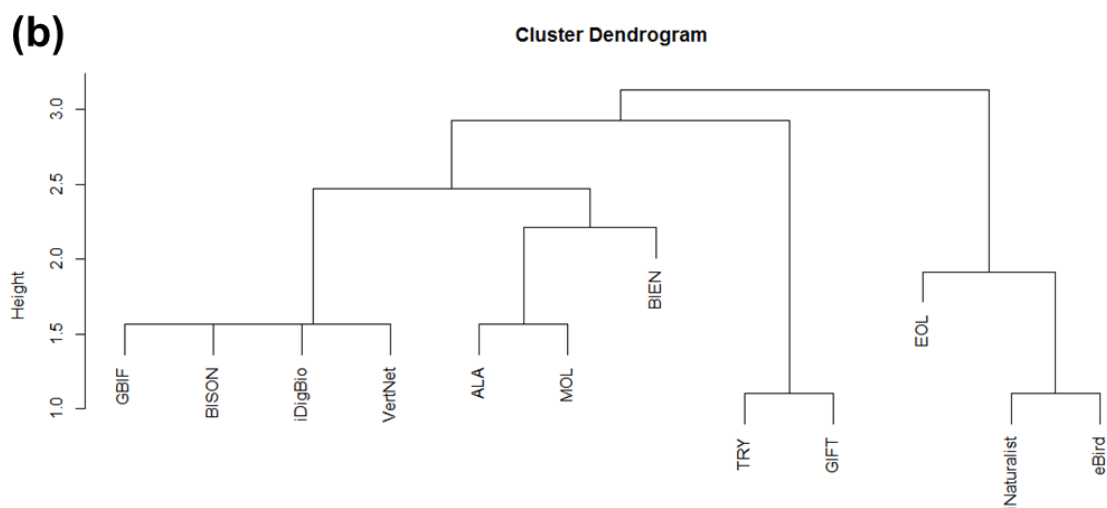
687

688 **Data and materials availability:** The plant distribution data from Global Biodiversity
689 Information Facility are accessible from <https://doi.org/10.15468/dl.87zyez>. Trait data from
690 Encyclopedia of Life are accessible from <https://eol.org/docs/what-is-eol/traitbank>. Trait data
691 from VertNet are accessible from <http://portal.vertnet.org/search>. Plant distribution and trait data
692 from Botanical Information and Ecology Network are accessible from RBIEN package. Trait
693 data from TRY are accessible from <https://try-db.org/TryWeb/dp.php>. The data from Catalogue
694 of Life are accessible from [https://download.catalogueoflife.org/col/monthly/2021-04-](https://download.catalogueoflife.org/col/monthly/2021-04-05_dwca.zip)
695 [05_dwca.zip](https://download.catalogueoflife.org/col/monthly/2021-04-05_dwca.zip). The administrative boundary dataset is accessible from
696 https://biogeo.ucdavis.edu/data/gadm3.6/gadm36_shp.zip.

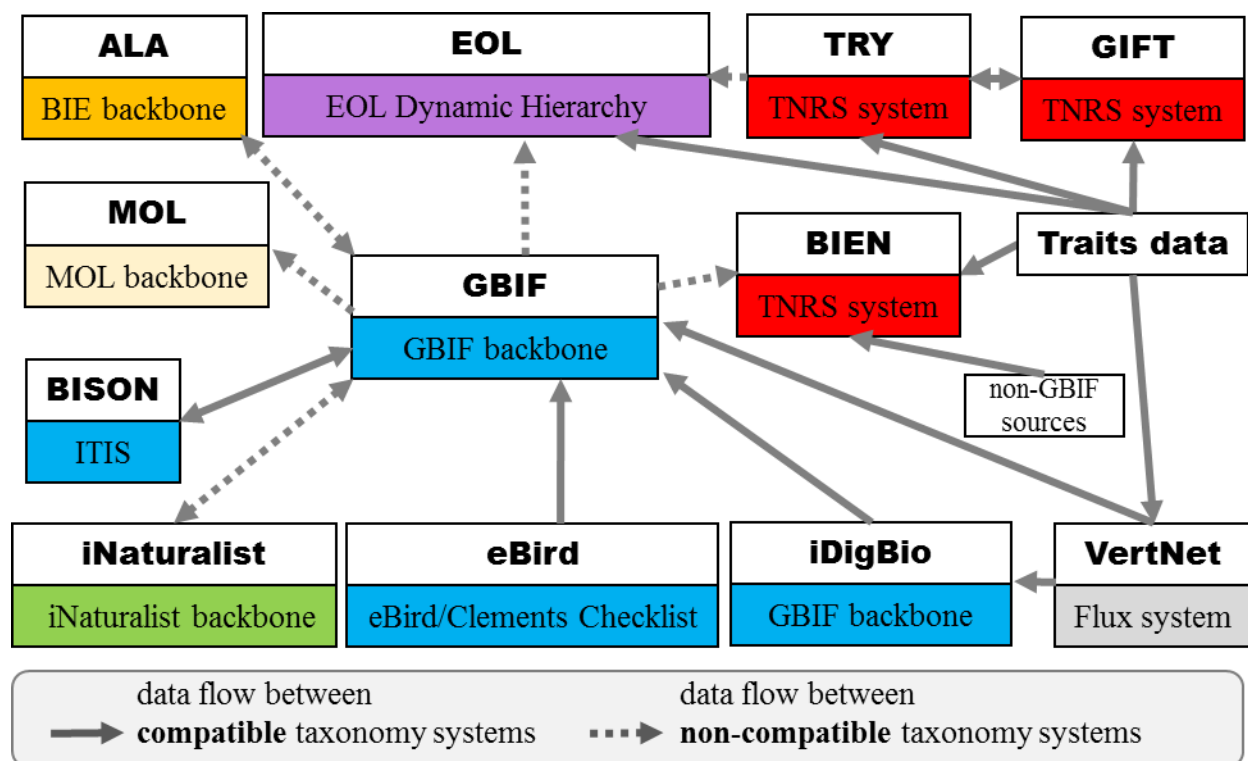
697

698

(a)		Database														
		GBIF	EOL	BISON	iDigBio	ALA	iNaturalist	MOL	BIEN	TRY	GIFT	eBird	VertNet			
Data category																
Spatial extent		Global	Global	USA & Canada	Global	Australia	Global	Global	Global	Global	Global	Global	Global	Global	Global	Global
Taxonomic group		All	All	All	All	All	All	All	All	Plants	Plants	Plants	Birds	Vertebrates		
Geographic distribution		Specimen	X		X	X		X	X						X	
		Observation	X		X		X	X	X	X			X	X		
		Checklist	X						X	X		X	X			
		Map		X			X	X	X	X				X		
Media		Images	a	X		a	X	X					X	a		
		Audio		X				X					X			
		Video		X				X					X			
Biology		Trait		X						X	Xb	Xb		X		
		Description		X			X	X	X				X			



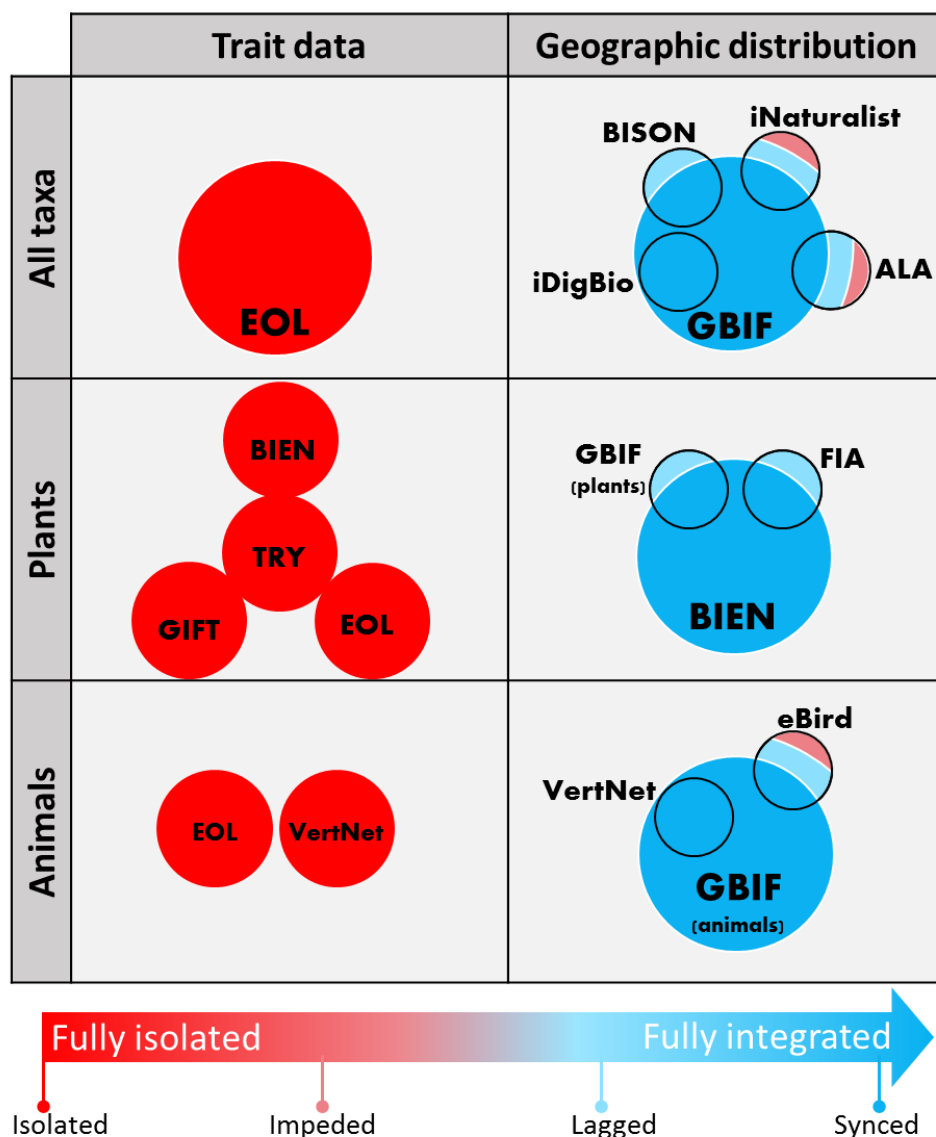
699
700 **Figure 1.** Overview of biodiversity databases reviewed in this paper. The coverages of their data
701 are shown in panel (a) indicated by “X”. Based on the data coverages, the biodiversity databases
702 are grouped into several clusters (b), where the height of the dendrogram is the relative distance
703 between clusters. Notes: a) GBIF, iDigBio, and VertNet indexes and displays images on its
704 website, while the images are mainly hosted by external institutions or facilities. b) TRY and
705 GIFT also stores geographic information about where the trait was measured.

706
707

708

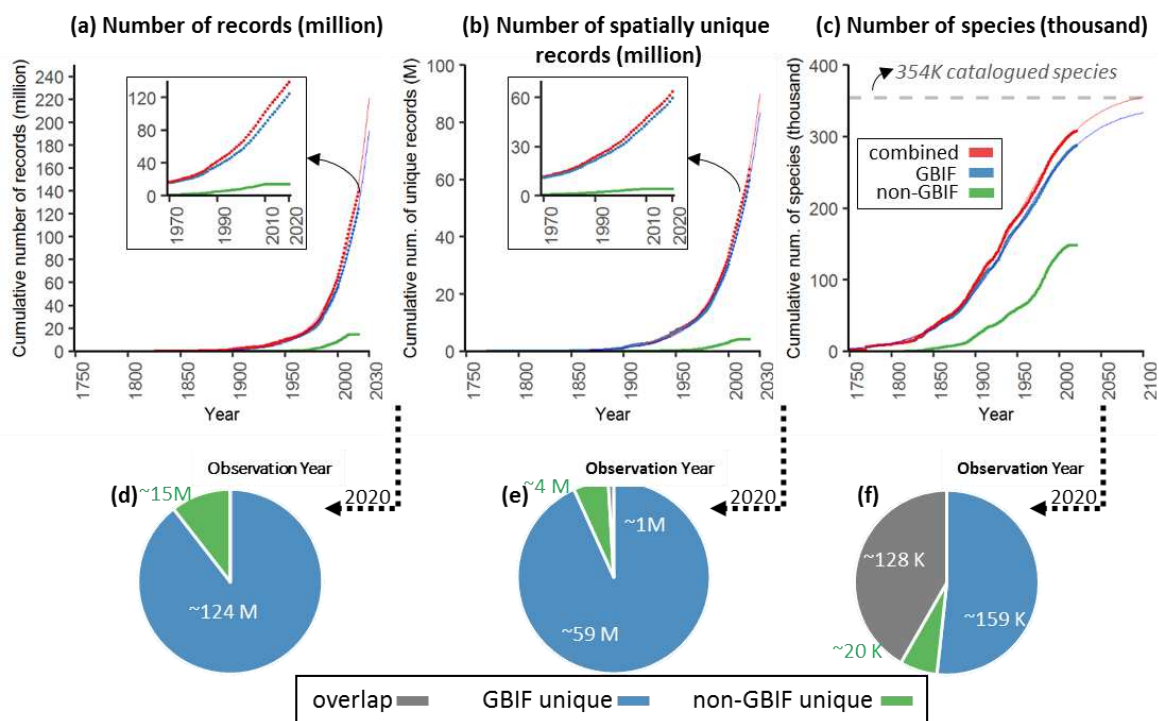
709 **Figure 2.** Data exchange between biodiversity databases with different taxonomic systems. Each
 710 box represents one database and its adopted taxonomic system (lower half). The taxonomic
 711 systems are shown in different colors, while the same color represents compatible systems. A
 712 variety of taxonomic systems exist: some databases develop backbone systems (e.g. BIE
 713 backbone, GBIF backbone, MOL backbone), some databases adopt a name scrubbing tool that
 714 standardizes names towards pre-selected taxonomic systems (e.g. BIEN, GIFT, TRY), some rely
 715 on multiple taxonomic systems (e.g. iNaturalist, EOL), and some do not implement a strong
 716 regulation on taxonomic names (e.g. VertNet). The one-way or two-way arrow represents
 717 unidirectional or bidirectional data flow between databases. ALA: Atlas of Living Australia;
 718 BIE: Biodiversity Information Explorer; BIEN: Botanical Information and Ecology Network;
 719 BISON: Biodiversity Information Serving Our Nation; EOL: Encyclopedia of Life; GBIF:
 720 Global Biodiversity Information Facility; GIFT: Global Inventory of Floras and Traits; iDigBio:
 721 Integrated Digitized Biocollections; ITIS: Integrated Taxonomic Information System; IUCN:
 722 International Union for Conservation of Nature; MOL: Map of Life; TNRS: Taxonomic Name
 723 Resolution Service; TRY: TRY, a global database of plant traits; uBio: Universal Biological
 724 Indexer and Organizer. As the databases continue to grow and develop, this figure represents the
 725 best of our knowledge as of March 2021.

726



727
 728 **Figure 3.** Data integration among biodiversity databases. The status of data integration is
 729 classified as four categories: synced, lagged, impeded, and isolated. *Synced* refers to the status
 730 of full integration, in either one or multiple directions, between different databases in or near
 731 real-time. For example, data published to iDigBio is automatically published to GBIF. *Lagged*
 732 refers to the difference between otherwise fully integrated databases between two sync events.
 733 For example, BIEN imports and integrates data from GBIF and other sources (e.g., The Forest
 734 Inventory and Analysis or FIA) annually or at longer intervals and publishes the results as
 735 versioned database releases. The most recent data in those sources will not be available via BIEN
 736 until the next import and versioned release. *Impeded* refers to differences between databases
 737 caused by barriers that prevent subsets of the data from being shared. For example, iNaturalist
 738 only publishes data to GBIF that are properly licensed for open sharing (iNaturalist, 2018).
 739 Contrary to distribution databases, trait databases are generally isolated from one another in
 740 different databases, though there are flows/exchanges of plant trait data between TRY and GIFT,
 741 and TRY and EOL (Table S1). We caution that the data flow between or among databases is not
 742 well documented, and this figure represents the best of our knowledge as of March 2021.

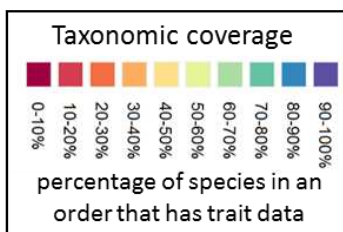
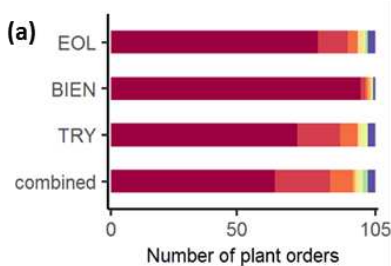
Plant Distribution Data



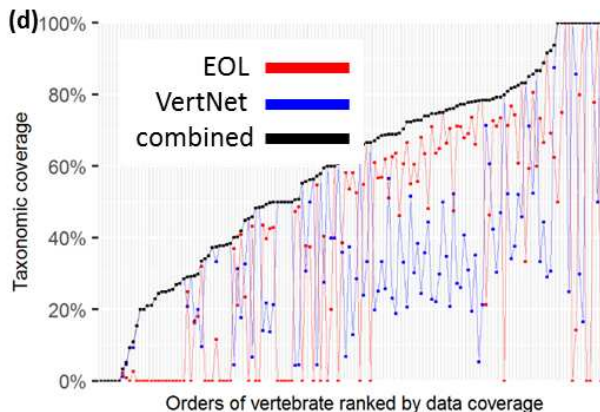
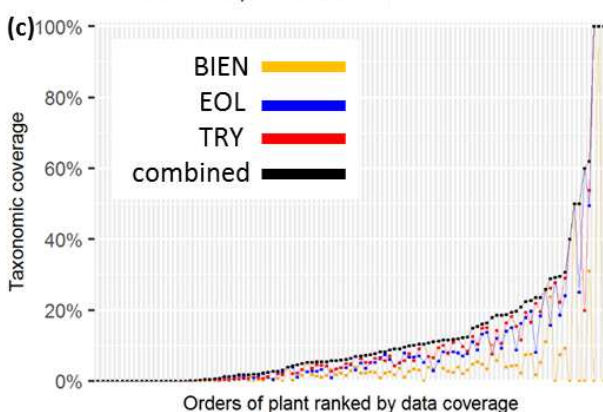
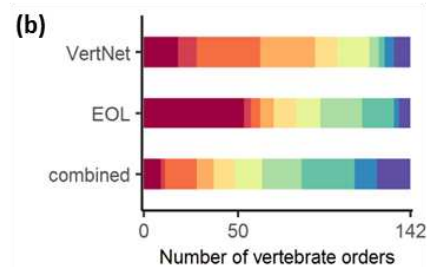
743
 744 **Figure 4.** Spatial and taxonomic coverage of terrestrial plant occurrence data. Georeferenced
 745 plant observations, as illustrated by observation dates in GBIF, the largest biodiversity
 746 informatics infrastructure, have increased exponentially over the past 200 years (panel a,b),
 747 though the number of species recorded in these databases is reaching saturation (panel c). By
 748 integrating additional data sources compiled by BIEN (i.e. non-GBIF sources; ~15 million
 749 records; panel d), the georeferenced plant observations in GBIF can be expanded by an
 750 additional ~4 million spatially unique records (panel e) and ~20 thousand species (panel f). Still,
 751 the gaps in plant distributions warrant our attention: large areas in Russia, central Asia, and
 752 northern Africa (red area in panel g) are missing publicly available occurrences. The black color
 753 in panel g represents ice covered areas.

754

Plant Trait Data



Vertebrate Trait Data



755

756

757 **Figure 5.** Increased taxonomic coverage of plant and vertebrate trait data through data

758 integration. By combining trait databases, coverage could be expanded in 69-82 plant orders

759 (panel a) and 86-124 vertebrate orders (panel b) compared to individual data sources (panel c &

760 d). The taxonomic coverage of a database is measured as the percentage of the species in that

761 plant or vertebrate order that are represented. Panels c & d show the taxonomic coverages of

762 individual databases and the combined dataset; the positions of the points on the x-axis are re-

763 ordered from low to high based on the combined taxonomic coverage (orders with low coverage

764 on the left and orders with high coverage on the right).

765 **Table S1. Summary of data flow among biodiversity databases.**

From	To	Details	References/Links
ALA	GBIF	ALA is a GBIF publisher, though data hosted by ALA may not be fully available on GBIF because of, for example, data licenses.	https://web.archive.org/web/20210506151646/https://www.gbif.org/publisher/3c5e4331-7f2f-4a8d-aa56-81ece7014fc8
GBIF	ALA	ALA includes exported data from GBIF that occur in Australia.	https://web.archive.org/web/20210407034945/https://collections.ala.org.au/public/showDataResource/dr695
GBIF	MOL	MOL includes exported data from GBIF.	https://web.archive.org/web/20210506152723/https://mol.org/datasets/9905692e-6a28-4310-b01e-476a471e5bf8
BISON	GBIF	BISON is a product of the United States Geological Survey (USGS) (Administrator of the U.S. Node of GBIF), and thus works closely and shares data with GBIF.	https://bison.usgs.gov/#help
GBIF	BISON	The Canadian and U.S. data added directly to GBIF would become available through BISON.	https://bison.usgs.gov/#help
iNaturalist	GBIF	iNaturalist is a GBIF publisher.	https://web.archive.org/web/20210506161424/https://www.gbif.org/publisher/28eb1a3f-1c15-4a95-931a-4af90ecb574d
GBIF	iNaturalist	iNaturalist displays data from GBIF on the interactive map.	https://www.inaturalist.org/taxa/71130-Polyphaga
GBIF	EOL	EOL incorporates data from GBIF.	https://web.archive.org/web/20210506162446/https://opendata.eol.org/dataset/gbif-data-summaries
eBird	GBIF	eBird Observational Dataset is published on GBIF.	https://web.archive.org/web/20210329225357/https://ebird.org/news/gbif/
TRY	EOL	TRY summarized records are available from EOL.	https://web.archive.org/web/20210326174302/https://eol.org/resources/504
TRY	GIFT	Co-develop and exchange trait data on plant growth form.	(Kattge et al., 2020)
GIFT	TRY	Co-develop and exchange trait data on plant growth form.	(Kattge et al., 2020)
GBIF	BIEN	BIEN includes data exported from GBIF.	https://web.archive.org/web/20210506163327/https://bien.nceas.ucsb.edu/bien/biendata/bien-2/sources/

iDigBio	GBIF	iDigBio is a GBIF publisher.	https://web.archive.org/web/20210506164312/https://www.gbif.org/publisher/2053a639-84c3-4be5-b8bc-96b6d88a976c
VertNet	GBIF	VertNet is a GBIF publisher.	https://web.archive.org/web/20210329192932/http://vertnet.org/join/ipt.html
VertNet	iDigBio	The majority of the data in the datasets published by VertNet are available in other portals such as GBIF, Canadensys, and iDigBio.	https://web.archive.org/web/20201012204516/vertnet.org/resources/datalicensingguide.html

766

767

768 **Table S2. Summary of taxonomic system of biodiversity databases.**

Name	Taxonomic system	References
GBIF	GBIF backbone	https://doi.org/10.15468/39omei
ALA	Biodiversity Information Explorer (BIE) backbone	https://web.archive.org/web/20210407032823/https://www.ala.org.au/blogs-news/updates-to-alas-name-and-taxonomy-index/
MOL	MOL developed a backbone that includes Catalogue of Life and manually curated taxonomic datasets for synonym issues.	<i>Anonymous reviewer</i>
BISON	Integrated Taxonomic Information System (ITIS)	https://web.archive.org/web/20210505185337/https://bison.usgs.gov/
iNaturalist	iNaturalist backbone is composed of global taxonomic authorities, regional taxonomic authorities, primary literature, and other name providers including Catalogue of Life and uBio.	https://web.archive.org/web/20210505185713/https://www.inaturalist.org/pages/curator+guide
EOL	The EOL Dynamic Hierarchy is curated by EOL staff based on a suite of classification providers (including Catalog of Life, the International Union for Conservation of Nature (IUCN), the National Center for Biotechnology Information (NCBI) and the World Register of Marine Species (WoRMS)) for different branches and layers of the tree of life, and can be manually patched and curated.	https://web.archive.org/web/20210505190456/https://eol.org/docs/what-is-eol/whats-new
TRY	Plant taxonomy of the TRY database is consolidated using the Taxonomic Names Resolution Service (TNRS) with a taxonomic backbone based on the Plant List, Tropicos, the Global Compositae Checklist, the International Legume Database and Information Service, and USDA's Plants Database.	<i>(Kattge et al., 2020)</i>
GIFT	The GIFT database standardized non-hybrid species names in The Plant List 1.1 and additional resources available via iPlant's Taxonomic Name Resolution Service (TNRS).	<i>(Weigelt et al., 2017)</i>
BIEN	Taxon names were corrected and standardized using the Taxonomic Name Resolution Service v5.0 (TNRS) with Tropicos, The Plant List and USDA Plants as taxonomic references, and all other options at their default settings.	<i>(Enquist et al.)</i>
eBird	eBird/Clements Checklist The eBird species and subspecies taxonomy follows the Clements Checklist. In addition to the formal taxonomic concepts that are included in the Clements Checklist, the eBird taxonomy includes an expanded list of other bird taxa that birders may report.	https://web.archive.org/web/20210505232653/https://ebird.org/science/use-ebird-data/the-ebird-taxonomy
iDigBio	The scientific names are matched to the GBIF backbone to correct typos and older names.	https://web.archive.org/web/20210505233105/https://www.idigbio.org/wiki/index.php/Data_Ingestion_Guidance
Vertnet	Flux system VertNet does not have a simple taxon resolution mechanism, and vertebrate species names are particularly in flux.	<i>(Zermoglio et al., 2016)</i>

770 **Table S3.** Summaries of model fitting for the temporal trend in plant distribution data.

Data source	Data	Model	AIC	Inflection point
combined	number of records	exponential	-1686	n/a
		linear	-239	n/a
		logistic	NA	NA
	number of spatially unique records	exponential	-1916	n/a
		linear	-258	n/a
		logistic	NA	NA
	number of species	exponential	-739	n/a
		linear	-510	n/a
		logistic	-1682	1947
GBIF	number of records	exponential	-1816	n/a
		linear	-315	n/a
		logistic	NA	2059
	number of spatially unique records	exponential	-1957	n/a
		linear	-301	n/a
		logistic	NA	NA
	number of species	exponential	-804	n/a
		linear	-552	n/a
		logistic	-1762	1949

771

772

773 **Appendix 1. Materials and Methods**

774 **Metadata review**

775 Many biodiversity databases have been built over the past decade, with varying emphases on
776 taxonomy, spatial location, and record type. Associated metadata for biodiversity databases is
777 typically found in publications or project websites. To synthesize the major attributes of existing
778 biodiversity databases, we selected 12 well-established biodiversity databases: Atlas of Living
779 Australia (ALA (Belbin & Williams, 2016)), Botanical Information and Ecology Network (BIEN
780 version 4.1 (Enquist *et al.*, 2016)), Biodiversity Information Serving Our Nation (BISON (U.S.
781 Geological Survey, 2018)), eBird (Sullivan *et al.*, 2014), Encyclopedia of Life (EOL (Parr *et al.*,
782 2014)), Global Biodiversity Information Facility (GBIF), Global Inventory of Floras and Traits
783 (GIFT (Weigelt *et al.*, 2017)), Integrated Digitized Biocollections (iDigBio (iDigBio, 2018a)),
784 iNaturalist (iNaturalist), Map of Life (MOL (Jetz *et al.*, 2012)), a global database of plant traits
785 (TRY version 1.0 (Kattge *et al.*, 2011)), and VertNet (Constable *et al.*, 2010). The twelve
786 databases we examined were chosen among the most commonly used, well-established, large-
787 scale biodiversity databases (MacFadden & Guralnick, 2016; Chandler *et al.*, 2017a; James *et*
788 *al.*, 2018; Singer *et al.*, 2018; Cornwell *et al.*, 2019; König *et al.*, 2019) to maximize the
789 generalizability of our results and conclusions. Selections were also limited to databases from
790 which we could either access the entirety of the data or the ones with clear documentations. We
791 compiled information from online documentation and relevant publications, though the design
792 and architecture of a database can be in continuous development. Specifically, we recorded
793 database name, taxonomic scope, taxonomic system, record type, number of records, and spatial
794 coverage. We classified the record types into three categories: geographic distribution, media
795 (image, audio, or video), and biological information (standardized trait databases or generalized
796 text descriptions). Within geographic distribution, we further classified the information as
797 specimen records, observations, checklists of geographic regions, and distribution maps.
798 Specimen records and observations both have information on species' geolocations, but only
799 specimen records are associated with physical specimens. Checklists usually contain lists of
800 species known to be present in certain geographic regions (e.g., political divisions or protected
801 areas). Distribution maps are either drawn by experts or generated through models. There are
802 frequent data exchanges among biodiversity databases, but many are not transparent to database
803 users. Consequently, we compiled data exchange information and assessed the status of data
804 integration between databases. We used geographic distribution and trait data as examples,
805 which are the most prominent record type among the reviewed databases. We assessed the
806 integration status by taxonomy groups, which are all organisms, plants, or vertebrates

807

808 **Improvement of data coverage by database integration**

809 To quantify the improvement gained by combining multiple databases, we compared leading
810 databases that focus on similar taxonomic groups and record type. We used terrestrial plants
811 (Embryophyta) and vertebrates as test cases, because these are the taxonomic groups that are
812 comparatively better collected and documented in biodiversity databases compared to other
813 taxonomic groups (Clark & May, 2002; Fazey *et al.*, 2005; Hecnar, 2009; Titley *et al.*, 2017;
814 Cornwell *et al.*, 2019; König *et al.*, 2019; Kattge *et al.*, 2020). We did not use taxa, such as
815 microbes, that account for large portions of biodiversity on Earth but face huge data gaps (Locey
816 & Lennon, 2016). More specifically, we compared (1) plant distribution data from GBIF and
817 non-GBIF sources compiled by BIEN (Enquist *et al.*, 2016), (2) plant trait data (i.e. plant height)

818 from BIEN, TRY, GIFT, and EOL, and (3) animal trait data (i.e. vertebrate body length) from
819 VertNet and EOL.

820
821 We obtained plant distribution data from BIEN (version 4.2; accessed March 2021) that
822 compiled plant distribution data from GBIF (<https://doi.org/10.15468/dl.87zyez>) and non-GBIF
823 sources, such as the *Forest Inventory and Analysis* (U.S. Department of Agriculture Forest
824 Service) (FIA) and *NeoTropTree* (Oliveira-Filho, 2017). The GBIF and non-GBIF sources have
825 been fused through a series of data scrubbing and standardization workflows (e.g. TNRS (Boyle
826 *et al.*, 2013)) and here we only included data with valid collection year and spatial coordinates.
827 We classified the data into three groups: data from GBIF, data from non-GBIF sources, and the
828 combined full dataset. We quantified the numbers of distribution records, numbers of spatially
829 unique records, and numbers of species with distribution records in all three data sources. A
830 spatially unique record is defined as a record of the distribution of a species (a pixel at 30 arc-
831 seconds resolution in WGS84 coordinate reference system that its coordinate corresponds to) that
832 is unique to a dataset. We standardized all species names against multiple reference taxonomies,
833 including *Tropicos* and *The Plant List*, through the *TNRS* (Boyle *et al.*, 2013). The
834 standardization process parses and corrects misspelled names and authorities, standardizes
835 variant spellings, and converts nomenclatural synonyms to currently accepted names. To reveal
836 the temporal trend of data accumulation, we quantified the cumulative numbers of observations
837 made over time, from 1750 to present (2020).

838
839 To describe and quantify those temporal trends, we fitted the cumulative numbers (dependent
840 variable) and years (independent variable) with simple linear (eqn 1), exponential (eqn 2), and
841 logistic regression (eqn 3) using ordinary least squares (“nls” function in stats package version
842 3.4.2 in R version 3.4.2):

$$843 \quad y = a + b * x \text{ (eqn 1)}$$

$$844 \quad y = e^{a+b*x} \text{ (eqn 2)}$$

$$845 \quad y = \frac{a}{1 + e^{-b-c*x}} \text{ (eqn 3)}$$

846 where x represents time and y represents either number of records, number of spatially unique
847 records, or the number of species. We determined the best model fit from the lowest Akaike
848 Information Criterion value (AIC). To reveal the contribution of GBIF or non-GBIF sources to
849 the combined dataset, we quantified the commonalities and uniqueness of GBIF and non-GBIF
850 subsets in terms of number of records, number of spatially unique records, and number of species
851 with distribution data. For our quantification of the temporal trend in the number of species
852 observed, we also retained only currently accepted names to reduce uncertainty (Berendsohn,
853 1997; Franz & Peet, 2009; Boyle *et al.*, 2013), which yield comparable temporal pattern.
854 We identified knowledge gaps in two ways. We showed the pixels (at 30 arc-seconds resolution
855 in WGS84 coordinate reference system) for which there were no valid plant geolocation data,
856 and quantified the geographic area of those pixels (in Eckert IV equal area projection). We
857 caution that the gap here may be an overestimation because the plant distribution data compiled
858 by BIEN (including the data exported from GBIF) do not include all possible data sources, but
859 rather shareable data that are mainly publicly available. We then calculated the taxonomic
860 completeness of the distribution data at the level of plant orders. We obtained a list of accepted
861 names of extant terrestrial plant species from the *Catalogue of Life* (Catalogue of Life, 2021) and
862 considered that as the master list of known species. All taxonomic names were standardized

863 through TNRS (Boyle *et al.*, 2013). We obtained the order level completeness by calculating the
864 percentage of species in a plant order that have distribution information in the combined dataset.
865

866 In addition to distribution data, we also investigated the improvement in taxonomic coverage of
867 trait data through database integration, specifically terrestrial plant height and vertebrate body
868 length. We downloaded plant height data from BIEN, EOL, and TRY (accessed March 2021).
869 We also obtained a list of accepted names of extant terrestrial plant species from *Catalogue of*
870 *Life* (accessed March 2021) and considered that as the master list of known species. All
871 taxonomic names were standardized through TNRS (Boyle *et al.*, 2013). We calculated the
872 taxonomic completeness of species trait information at the species and order levels. We obtained
873 the species level completeness by checking species whose heights were recorded in BIEN, EOL,
874 TRY, or the combined dataset, against the names recorded in COL. We obtained the order level
875 completeness by calculating the percentage of species in a plant order that have height
876 information in either dataset. We calculated the improvement in percentages by comparing
877 individual datasets to the combined dataset. The improvement in taxonomic coverage represents
878 the benefit of using multiple databases.
879

880 Following the same workflow, we quantified the taxonomic coverage of animal trait and
881 percentage improvement between individual dataset and the combined dataset. Body length of
882 vertebrates were downloaded from VertNet and EOL (accessed March 2021). Accepted names of
883 extant vertebrates were obtained from *Catalogue of Life*. The taxonomic names were
884 standardized through Global Names Resolver using the *Taxize* package (Chamberlain & Szocs,
885 2013) (version 0.9.4.9100) in R (version 3.4.2). The Global Names Resolver resolves names
886 against specific name databases, which is *Catalogue of Life* in this study. The resolution process
887 includes a series of exact and fuzzy matches based on the full or part of the name input (see more
888 details in <https://resolver.globalnames.org/about>). The matching process also considers the
889 context of taxonomy and reduces the likelihood of matches to taxonomic homonyms. The
890 matching process yields a series of confidence scores for all possible matches; here we only kept
891 the best matching records. However, the creation of a single authoritative list of names will take
892 time; full reconciliation of synonyms and distinct taxon concepts may take decades (Berendsohn,
893 1997; Franz & Peet, 2009; Boyle *et al.*, 2013). The standardization of taxonomic names based on
894 either TNRS or Global Names Resolver will not solve all issues of taxonomic name integration,
895 but this step represents the state-of-the-art in standardizing taxonomy names in biodiversity
896 databases and provides a baseline for the comparisons of different biodiversity databases.
897