# Research applications of primary biodiversity databases in the digital age

**Joan E. Ball-Damerow**[1]*, **Laura Brenskelle**[2], **Narayani Barve**[2], **Pamela S. Soltis**[2], **Petra Sierwald**[1], **Rüdiger Bieler**[1], **Raphael LaFrance**[2], **Arturo H. Ariño**[3], **Robert P. Guralnick**[2]

**1** Field Museum of Natural History, Chicago, IL, United States of America, **2** Florida Museum of Natural History, University of Florida, Gainesville, FL, United States of America, **3** Department of Environmental Biology, Universidad de Navarra, Pamplona, Spain

* joandamerow@gmail.com

## Abstract

Our world is in the midst of unprecedented change—climate shifts and sustained, widespread habitat degradation have led to dramatic declines in biodiversity rivaling historical extinction events. At the same time, new approaches to publishing and integrating previously disconnected data resources promise to help provide the evidence needed for more efficient and effective conservation and management. Stakeholders have invested considerable resources to contribute to online databases of species occurrences. However, estimates suggest that only 10% of biocollections are available in digital form. The biocollections community must therefore continue to promote digitization efforts, which in part requires demonstrating compelling applications of the data. Our overarching goal is therefore to determine trends in use of mobilized species occurrence data since 2010, as online systems have grown and now provide over one billion records. To do this, we characterized 501 papers that use openly accessible biodiversity databases. Our standardized tagging protocol was based on key topics of interest, including: database(s) used, taxa addressed, general uses of data, other data types linked to species occurrence data, and data quality issues addressed. We found that the most common uses of online biodiversity databases have been to estimate species distribution and richness, to outline data compilation and publication, and to assist in developing species checklists or describing new species. Only 69% of papers in our dataset addressed one or more aspects of data quality, which is low considering common errors and biases known to exist in opportunistic datasets. Globally, we find that biodiversity databases are still in the initial stages of data compilation. Novel and integrative applications are restricted to certain taxonomic groups and regions with higher numbers of quality records. Continued data digitization, publication, enhancement, and quality control efforts are necessary to make biodiversity science more efficient and relevant in our fast-changing environment.

## Introduction

Online databases with detailed information on organism occurrences collectively contain well over one billion records, and the numbers continue to grow. The digitization of natural history specimens [1,2] and development of online platforms for citizen science [3] have driven a steady accumulation of species occurrence records over the past decade. Each data point provides details on the taxonomic identification, date collected or observed, location, and name of the collector or observer for an organism. Applications of these primary biodiversity data are varied—such data have historically helped determine harmful effects of pesticides, document spread of infectious disease and invasive species, monitor environmental change, and much more [4–9]. The overall goal of this paper is to determine how researchers use open-access data in published work, focusing on the past decade, when growth of online biodiversity databases has been most rapid. As one illustration of that growth, the Global Biodiversity Information Facility (GBIF) has grown from provisioning just over 200 million records in 2010 to over 1.08 billion records today, a greater than fivefold increase [10].

Museums and funding agencies have invested considerable resources to digitize information from natural history specimens, make their data openly accessible [11,12], and sustain platforms to provide access to those data. Such efforts unlock previously inaccessible data and expand their availability to researchers around the world. However, the task of digitizing highly diverse groups, such as insects, has been particularly difficult. Estimates suggest that only 10% of biocollections worldwide are available in digital form [13,14], and it would take many decades to completely digitize estimated holdings at current rates [15]. While efforts towards workflow optimization will undoubtedly improve efficiency in certain areas [12,16–19], it is critical that the biocollections community prioritize efforts; we must advocate for continued digitization through production of innovative data products, tools, interdisciplinary collaborations, and by highlighting research that requires primary biodiversity data [3,20–22]. The greatest returns on digitization investments will result from expanded use of collections data and by linking a wide array of biotic and abiotic data [1,11]. Linked data environments are in high demand [23,24], are growing rapidly, and provide the greatest potential for data discovery and use [1].

The biggest obstacle for biodiversity data users is obtaining records of sufficient quantity and quality for the region and taxonomic group of interest [24,25]. Many taxa and regions are still highly under-sampled or completely unrepresented (e.g. rare taxa, regions that are difficult to access) in online databases [26–28], particularly for less known and highly diverse invertebrates [29,30]. Many records are also prone to missing important information or information loss over time, particularly the absence of geographic coordinates and associated uncertainty estimates [31]. When data are available, researchers must check for common errors and biases known to occur in opportunistic datasets that are often assembled over long time periods (e.g. [32])—a task that is labor-intensive [33]. Species identity and locality are the most error-prone aspects of collection information [7]. Estimates for rates of collection misidentification range from 5–60%, depending on the taxonomic group [11,34,35]. But if specimens exist, this information can be verified or corrected by taxonomic experts. Specimen images, while not always useful for diagnosis, can often help—particularly when they meet the criteria for taxonomic-grade imaging. Even with correct identification, names in species occurrence repositories may still be incorrect and need validation [36]. For many broad-scale studies, erroneous records primarily lead to overestimation of species richness in areas outside centers of diversity [33]. Geographic errors (or missing information) may be more readily corrected and associated with appropriate uncertainty estimates using standardized methods [31,37] and online tools (i.e. GEOLocate, www.geo-locate.org). Digitization of species occurrence records allows

researchers to explore the data relatively quickly and identify outliers. Further, data services are becoming more sophisticated in automatically addressing some data quality issues [38,39]. However, it is possible that many studies simply use available data and may not appropriately evaluate data quality.

Sources of potential biases in opportunistic occurrence data have also been well-documented in previous work and generally include variation in collection effort and taxonomic, spatial, and temporal biases [4,40–45]. Some examples of variables contributing to bias include socioeconomic factors [44,45], the exclusion of common species over rare and flashy ones [46–48], the selection of large and attractive specimens [49], seasonal bias [50], problematic distinction between living and dead-collected specimens and associated post-mortem transportation [51,52], and discarding worn specimens, which results in phenological bias or elimination of specimens with signs of disease [8]. Traditional methods for dealing with these issues may include subsampling, data aggregation, and additional surveys [7]. Effects of bias can be reduced for certain studies with higher numbers of records, by combining information from different institutions, and including observation records to supplement specimen data [8]. Newer statistical and modeling approaches to deal with biases in biodiversity data have also been developed [43,48,53,54]. However, it is unclear how often studies actually address issues of error and bias when using opportunistic records.

While several previous studies have reviewed uses of natural history collections data [4,6,8,55], and one study has analyzed field-specific usage for the GBIF index [56], to our knowledge no other study has quantitatively reviewed trends in how species occurrence databases are utilized in published research. Our overarching goal in this study is to determine how such usage has developed since 2010, during a time of unprecedented growth of online data resources. We also determine uses with the highest number of citations, how online occurrence data are linked to other data types, and if/how data quality is addressed. Specifically, we address the following questions: What primary biodiversity databases have been cited in published research, and which databases have been cited most often? Is the biodiversity research community citing databases appropriately, and are the cited databases currently accessible online? What are the most common uses, general taxa addressed, and data linkages, and how have they changed over time? What uses have the highest impact, as measured through the mean number of citations per year? Are certain uses applied more often for plants/invertebrates/ vertebrates? Are links to specific data types associated more often with particular uses? How often are major data quality issues addressed? What data quality issues tend to be addressed for the top uses?

## Literature search and characterization

We searched for papers that use online and openly accessible primary occurrence records or add data to an online database. Google Scholar (GS) provides full-text indexing, which was important for identifying data sources that often appear buried in the methods section of a paper. Our search was therefore restricted to GS and to the time period of 2010 through the date of the search (April 2017; note when looking at trends over time we remove 2017, as the year was not complete in our dataset). All authors discussed and agreed upon representative search terms, which were relatively broad to capture a variety of databases hosting primary occurrence records. The terms included: *"species occurrence" database* (8,800 results), *"natural history collection" database* (634 results), *herbarium database* (16,500 results), *"biodiversity database"* (3,350 results), *"primary biodiversity data" database* (483 results), *"museum collection" database* (4,480 results), *"digital accessible information" database* (10 results), and *"digital accessible knowledge" database* (52 results)–note that quotations are used as part of the search

terms where specific phrases are needed in whole. We downloaded the first 500 records (or all if there were fewer than 500 results), which are presumably the most relevant search returns, for each search term into a Zotero reference management database [57]. We obtained citation numbers for each paper from the GS search results at the time of downloading records (April 2017) [58]. After removing duplicates across search terms, the final database included 2,460 papers. We then randomly sorted papers into four separate sets of around 500 to allow sub-sampling of the dataset.

For a study to be relevant in this assessment, there must be an indication that the database used is publicly accessible online in a searchable database of biodiversity records. The databases used may include specimen and/or observation-based records from biodiversity data aggregators, online natural history collection databases, websites devoted to capturing citizen science observation records, or newly compiled data that are made available in online databases. Studies were not relevant if they *exclusively* used data that are not available online or from systematic surveys, government monitoring programs, or field data collected explicitly for the study in question. However, papers are relevant if they use these other types of occurrence data *in addition to* online databases of primary occurrence records (see section on data linkages, below), or if they compile these types of occurrence records and deposit them into an existing online biodiversity data aggregator (e.g. GBIF). Twenty-six percent (*n* = 501; see S1 File for citation information) of the papers in the final evaluated dataset (*n* = 1,934) were relevant according to these criteria. The full dataset is published and openly accessible [58].

Three of the authors with specialized knowledge of the field (J. Damerow, L. Brenskelle, and R. Guralnick) characterized relevant papers for the first 1000 papers using a standardized tagging protocol based on 14 key topics of interest with over 100 total tags. We developed a list of potential tags and descriptions for each topic; a full list with descriptions of tags is provided in S1 Table. J. Damerow subsequently checked each tagged paper from the first 1,000 papers to maintain consistency and became the sole tagger for an additional 934 papers. This process allowed the development of a more standardized tagging protocol. The database of tagged papers was then downloaded from Zotero for further data checking and analysis. We used OpenRefine (www.openrefine.org), an open source tool for data cleaning that aggregates similar records for efficient clean-up, to standardize tags from the final dataset.

## Trends in uses of primary biodiversity data

We characterize a variety of ways in which researchers are using species occurrence records by assessing the prevalence of individual tags corresponding to topics of interest. We identify the most commonly cited databases and most-studied taxa, number of taxa addressed, most common research uses, the types of data most often linked to species occurrence records, and aspects of data quality addressed in these papers. In addition, we determine prevalence of these tags over time to assess positive or negative trends. Some expected trends include the following:

- Data uses requiring large numbers of dispersed records, such as species distribution models and biodiversity studies, will be the most common applications of online databases.

- Data papers and those describing a new database will increase over time as new venues have grown supporting such publications.

- Uses involving other online data types (i.e. barcoding, citizen science, species interactions) that can be linked to species occurrence records will increase.

- The number of species addressed will increase over time as more data become available online and projects leverage broader-scale data.

- The most common data quality issues addressed will be checks for correct taxonomic nomenclature and georeferences, which can often be assessed with readily-available online resources.

## Primary biodiversity databases and accessibility of data

We identify 347 primary biodiversity databases used in papers from our dataset (S2 Table), the URL for each database, and the scale (institution, regional, global, taxa) and regional or taxonomic focus (e.g. Australia, fish) of each database. We then evaluate citation information provided in each paper, and assess whether the data are currently available online or not by visiting associated URLs. The most cited databases include: the Global Biodiversity Information Facility (GBIF [10]), Barcode of Life Data Systems that includes species occurrence and genetic data (BOLDSystems [59, 60]), SpeciesLink [61], Ocean Biogeographic Information System (OBIS [62]), Australasian Virtual Herbarium (AVH [63]), Tropicos [64], FishBase [65,66], Fishes of Texas [67], and CONABIO REMIB (Table 1, [68]); note that we did not find significant changes over the study time period (2010–2017) in usage of individual databases, likely due to insufficient data points per year.

Our dataset includes 165 papers that involve compiling and publishing data online (117 data papers and 60 papers that describe a new database, some of these papers overlap). Previous work has outlined best practices for publication of biodiversity data [69–74] and scientific data more generally (e.g. [75]). However data are published, primary biodiversity data should also be integrated into an aggregate system with similar data, such as GBIF, OBIS, VertNet, iDigBio, or BOLDSystems [74].

Many researchers do not sufficiently cite databases used [76,77], and links to many databases become invalid over time [78–80]. We found that 34 percent of papers ($n$ = 170) had insufficient citation information for one or more databases; this meant that there was either no URL provided to access the database, or the URL was invalid. Twenty-six percent of databases ($n$ = 90) cited in one or more papers from our dataset were totally inaccessible at the time of this assessment. In some cases, researchers appropriately cited a database that is no longer in operation or has subsequently been integrated into an aggregate system. As a result of insufficient data citation practices and lack of data preservation, data are either completely lost or it is impossible to reproduce the dataset used and results. Study reproducibility, strongly linked to data persistence [78], is a key principle in the scientific process and a growing concern across scientific disciplines (e.g. [81]). Researchers who have compiled data from multiple

**Table 1. Top ten most used biodiversity databases (see S2 Table for a comprehensive list).**

| Database Name | Number of Papers Citing |
|---|---|
| GBIF | 155 |
| BOLDSystems | 27 |
| SpeciesLink | 21 |
| OBIS | 20 |
| Australia's Virtual Herbarium | 19 |
| Tropicos | 16 |
| FishBase | 14 |
| Fishes of Texas | 13 |
| CONABIO | 11 |

https://doi.org/10.1371/journal.pone.0215794.t001

sources for a particular analysis can better ensure that these data are accessible and get credit for the work involved in integrating datasets by formally publishing data with descriptive metadata and obtain a persistent DOI [75]. The prevalence of inaccessible databases and incomplete database citations indicates that many biodiversity researchers lack the resources to manage and preserve data for the long term and/or are unaware of best practices.
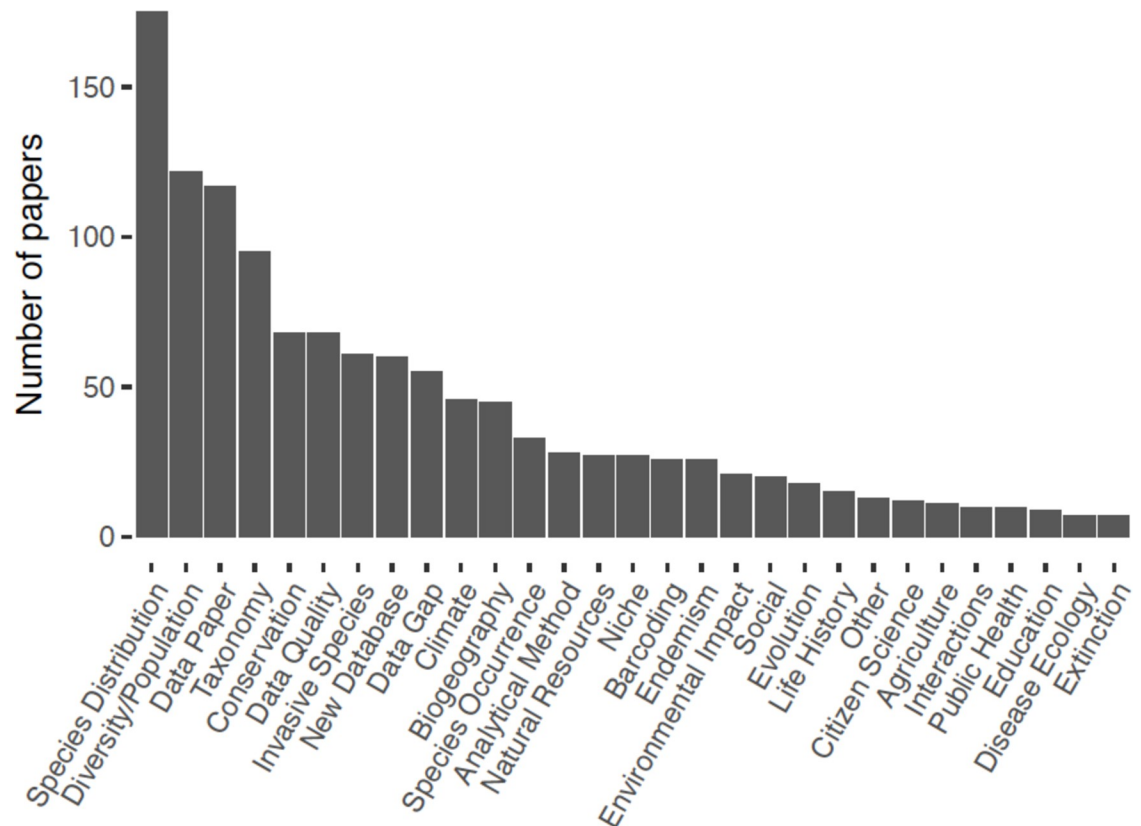
Guidance and infrastructure for citing online data sources have fairly recently emerged and are still evolving [76,82]. One major problem is that many papers using biodiversity data have obtained data from an aggregator, such as GBIF, which has potentially drawn from thousands of original data sources. Up to this point, researchers have most often cited GBIF in this case (usually in-text, not in the reference section) and neglect to credit original data sources [77]. Even for those who attempt to cite sources, many journals do not allow large numbers of citations in the reference section, and the only solution is to cite sources in a supplement or appendix which does not provide citation credit [77]. Data contributors who have submitted data to aggregators are not getting credit for the significant work spent on data management, standardization, and quality control. Ideally, data citations should include DOIs for datasets if they exist and citations of online databases both in text and in the reference section [76,77,83].

## Research uses

A primary goal for this work was to characterize research uses of the study databases. An initial list of use tags was developed based on usage outlined in [24], which surveyed needs of primary biodiversity data users. We subsequently split up certain aggregated topics and revised and added use categories based on important subject areas that arose during the tagging process. We ended with 31 potential research use tags, as listed and described in S1 Table. Most papers had multiple use tags assigned (mean = 2.5, max = 7). We then determined the average number of citations for papers involving each data use. Number of citations was extracted from the original web snapshots of the Google Scholar searches for each term in April 2017, and represent citations at that time [58].

The top research uses for online species occurrence databases—from our dataset of 501 relevant papers—were studies on species distribution ($n = 175$), diversity/population studies that usually assess species richness ($n = 122$), dataset description (i.e. data papers, $n = 117$), taxonomy ($n = 95$), conservation ($n = 68$), data quality ($n = 68$), invasive species ($n = 61$), and that described a new database ($n = 60$, Fig 1); see S1 Table for full descriptions of each category of research use. The prevalence of most uses did not change from 2010–2016, with the exception of data papers and taxonomy-related studies, which both increased (Fig 2); taxonomy studies usually involved developing regional species checklists. In the aforementioned survey assessment of user needs for primary biodiversity data [23,24], these same categories of use were among the top ways in which people listed that they use primary biodiversity data. Some exceptions were that a relatively large number of survey respondents claimed that they use biodiversity data for ecology/evolution studies, natural resources management, life history/phenology studies, and education/outreach, but relatively few published studies used occurrence data for these purposes in our dataset. It is possible that people use data for these purposes, but do not necessarily publish papers on the topic or may not cite databases for this work [84].

Some of the top research uses involved compiling and processing data, as reflected in the high numbers of data papers, papers describing new databases, and papers addressing data quality and data gaps (all of which were among the top ten uses, Fig 1). The biodiversity community is still in an active stage of compiling existing biodiversity data and dealing with issues of data quality. Data papers and papers describing a new database have increased over time (Fig 2), which is likely to be the result of the introduction and expansion of many data journals

**Fig 1. Frequency of major research uses in published papers (*n* = 501) that obtain data from species occurrence records available in online databases.** See S1 Table for detailed descriptions of each research type.

https://doi.org/10.1371/journal.pone.0215794.g001

[69,85], online platforms for reporting species occurrence observations such as iNaturalist [86] and eBird [3,87], and efforts over the past decade to digitize specimen records [1,13]. More journals accept papers or even focus on publishing high-quality data and recognize this as an important part of the scientific process [74,84,88,89].

Papers with the highest mean number of citations per year involved more applied studies in disease ecology (mean = 18, SD = 33), public health (mean = 8, SD = 7), documenting extinctions (mean = 7, SD = 7), developing a new analytical method to deal with species occurrence data (mean = 7, SD = 8), and citizen science (mean = 7, SD = 6; Table 2). Papers with the highest maximum number of citations per year focused on disease ecology, species diversity, and publishing data (each with a maximum of 97 citations/year; Table 2); we did not account for self-citation here.

## Taxa addressed

The third major topic for this work was to determine how often different taxonomic groups are represented in papers utilizing biodiversity databases. Taxa in relevant papers were coarsely characterized as plants, vertebrates, invertebrates, fungi, paleo, and/or all taxa; note that we addressed only macro-organisms because they are the focus of non-sequence-based species occurrence databases. These general taxonomic categories also correspond to common divisions for the organization of natural history collections and associated databases. Many papers include more than one taxon, and we use an "all taxa" categorization for studies that use all
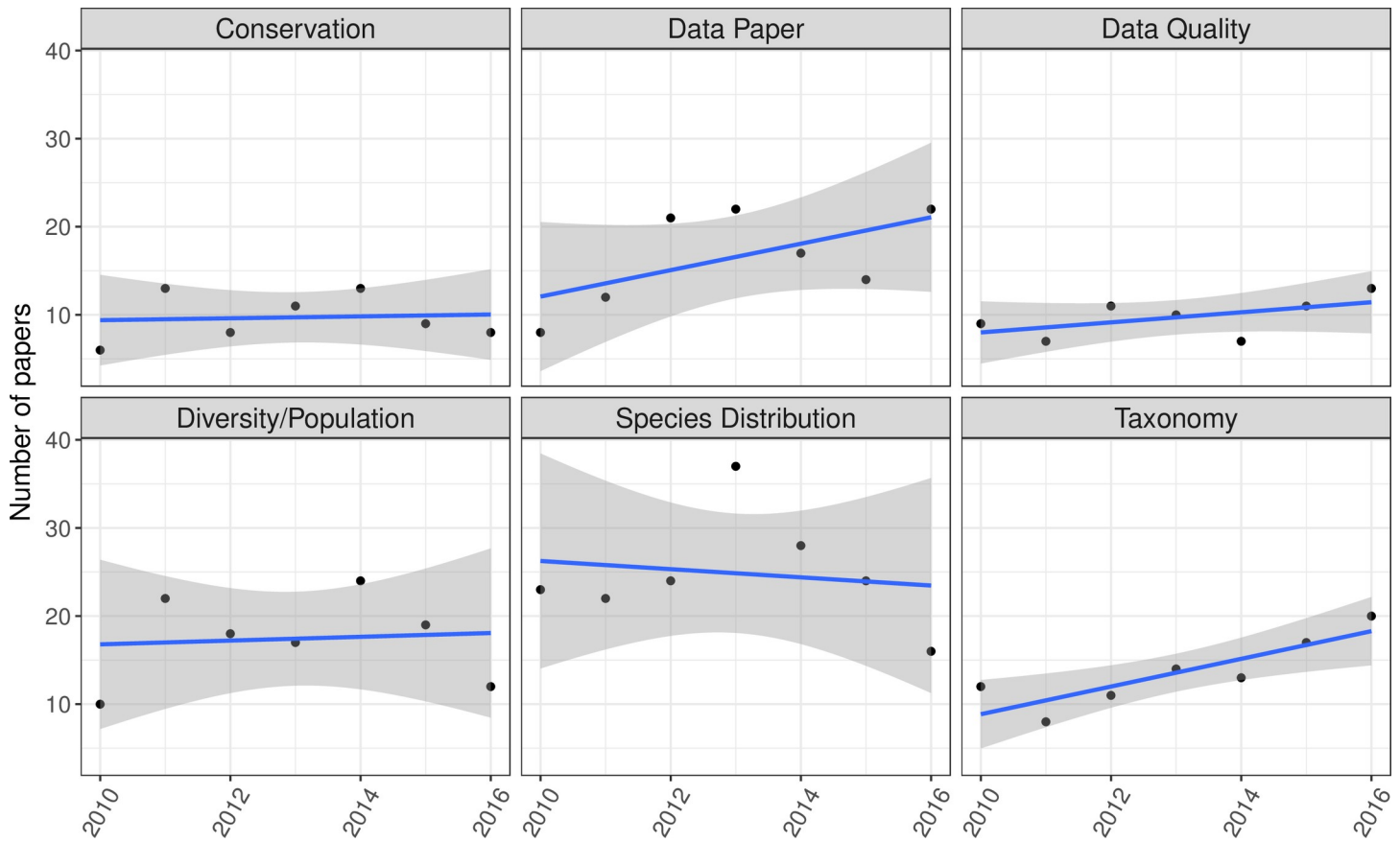
**Fig 2. Change in the number of papers from 2010–2016 involving the top six research applications for online species occurrence databases.**

available data within the species occurrence database(s), such as GBIF. We further categorized taxa addressed in each paper by adding one or more tag(s) for more specific taxonomic classifications (e.g. butterflies, *Danaus plexippus*). While an in-depth assessment of specific taxa is beyond the scope of the current paper, we did tag the number of taxa addressed in each paper, if that number was apparent. Our goals here were to characterize the most commonly studied taxonomic groups, the number of taxa addressed, and to determine uses associated with the three most common organismal groupings (plants, vertebrates, and invertebrates).

The most commonly studied taxa were plants ($n = 232$ papers, 46%), followed by invertebrates ($n = 125$, 25%), vertebrates ($n = 124$, 25%), "all taxa" ($n = 40$, 8%), fungi ($n = 16$, 3%), and paleontological specimens ($n = 14$, 3%; Table 3). However, the gap between number of papers addressing plants, vertebrates, and invertebrates closed in recent years (2014–2016, Fig 3). The overall prevalence of plants in this work corroborated a recent bibliometric study, which found that 56% of biodiversity-related papers addressed plants, compared to 29% for vertebrates and 23% for invertebrates [90]. The prevalence of plants in the field of biodiversity research may be the result of several factors. Plants are far more diverse than vertebrates (known to be relatively well-studied) and therefore generally require more taxonomic work. Herbarium sheets have also been the easiest historically to digitize, as sheets can be scanned and imaged using more automated processes [11,16]. The current prevalence of plants may also partially be the result of a strong history of plant research in Europe; this tendency is known as the "Matthew principle" whereby research concentrates on already well-studied

**Table 2. Summary statistics for the number of citations per year for each use of primary biodiversity data.** Note that not all papers had citation data available.

| Data Use | N | mean | sd | min | max |
|---|---|---|---|---|---|
| Disease Ecology | 8 | 18 | 33 | 2 | 97 |
| Public Health | 9 | 8 | 7 | 0 | 22 |
| Extinction | 6 | 8 | 7 | 1 | 17 |
| Analytical Method | 26 | 7 | 8 | 1 | 34 |
| Citizen Science | 7 | 7 | 6 | 1 | 17 |
| Species Distribution | 152 | 6 | 10 | 0 | 97 |
| Climate | 46 | 6 | 6 | 0 | 32 |
| Niche | 24 | 6 | 5 | 0 | 20 |
| Data Quality | 59 | 6 | 8 | 0 | 37 |
| Diversity/Population | 108 | 5 | 10 | 0 | 97 |
| Data Paper | 94 | 5 | 11 | 0 | 97 |
| Other(Paleontological) | 3 | 5 | 5 | 0 | 10 |
| Other(Behavior) | 1 | 5 | NA | 5 | 5 |
| Data Gap | 56 | 5 | 6 | 0 | 28 |
| Agriculture | 10 | 5 | 4 | 1 | 13 |
| Invasive Species | 55 | 5 | 5 | 0 | 32 |
| Conservation | 61 | 5 | 6 | 0 | 22 |
| Endemism | 23 | 5 | 5 | 0 | 20 |
| Evolution | 17 | 5 | 3 | 0 | 12 |
| Barcoding | 22 | 5 | 4 | 0 | 16 |
| Biogeography | 41 | 5 | 4 | 0 | 16 |
| New Database | 50 | 4 | 6 | 0 | 29 |
| Species Occurrence | 26 | 4 | 4 | 0 | 22 |
| Interactions | 7 | 3 | 3 | 1 | 9 |
| Natural Resources | 24 | 3 | 3 | 0 | 12 |
| Environmental Impact | 18 | 3 | 2 | 0 | 7 |
| Other(Movement) | 3 | 3 | 2 | 2 | 5 |
| Life History | 10 | 3 | 2 | 1 | 8 |
| Taxonomy | 72 | 2 | 3 | 0 | 16 |
| Other(Ethnobotany) | 1 | 2 | NA | 2 | 2 |
| Education | 5 | 2 | 2 | 0 | 5 |
| Social | 14 | 2 | 1 | 0 | 5 |
| Other(Reference) | 1 | 1 | NA | 1 | 1 |

subjects [90]. The total number of invertebrate studies was equivalent to the total number of vertebrate studies (Fig 3). However, invertebrates are much more diverse in terms of species

**Table 3. Total number of papers from dataset (501) addressing the major taxonomic groups and paleontological specimens.**

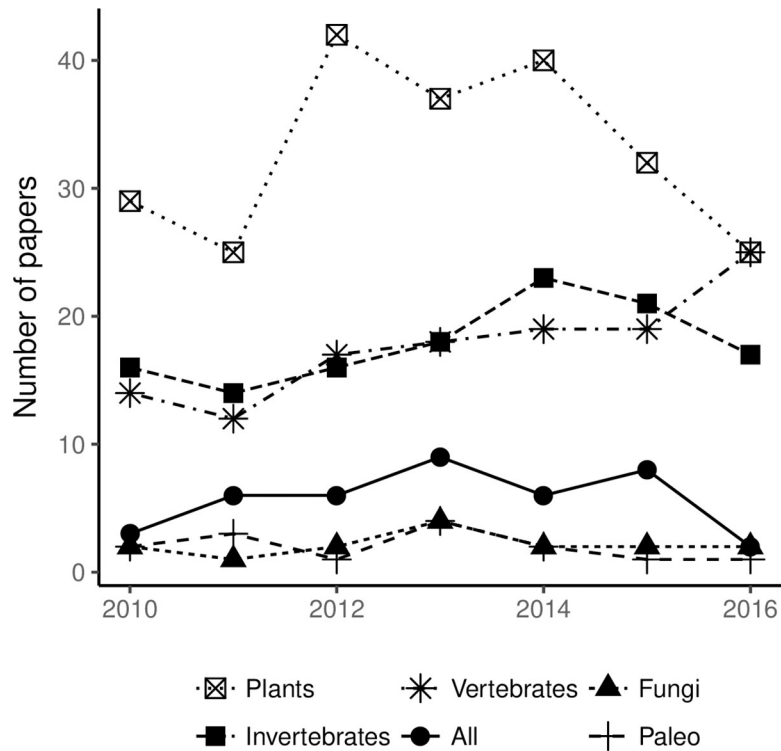| Taxa | Number of papers |
|---|---|
| Plants | 232 |
| Invertebrates | 125 |
| Vertebrates | 124 |
| All | 40 |
| Fungi | 16 |
| Paleo | 14 |

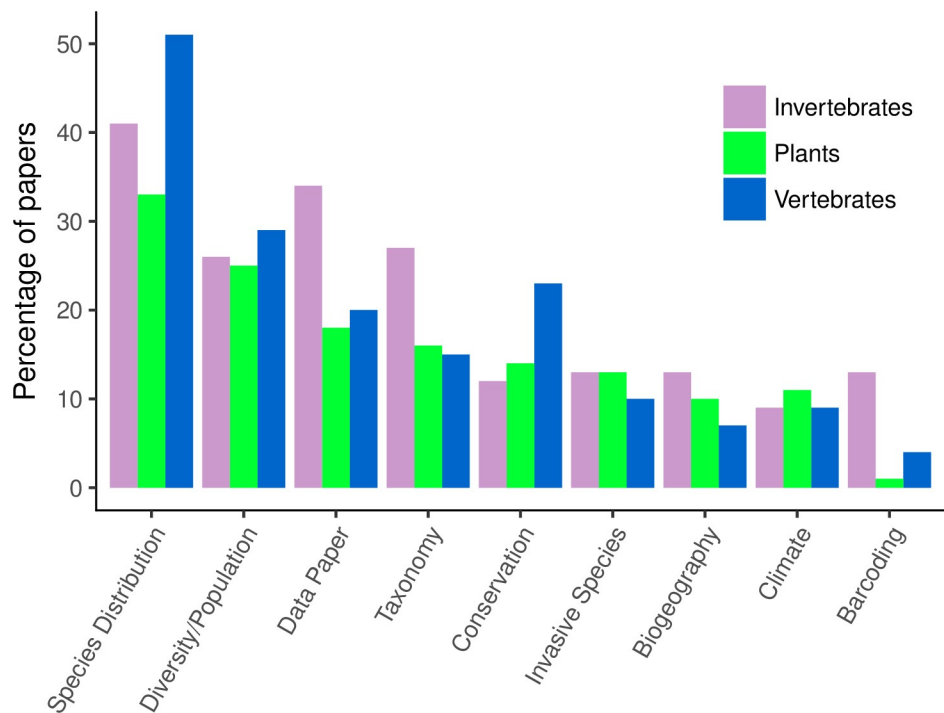**Fig 3. Number of papers addressing the major taxonomic groups and paleontological records over time.**

**Fig 4. Percentage of papers involving each of the major taxonomic groups (invertebrates, plants, and vertebrates) that use species occurrence databases for certain research applications: Species distribution, diversity/population, data paper, taxonomy, invasive species, biogeography, climate change, and barcoding.**

(estimated at 6,755,830 species, see [91]), and vertebrates are unquestionably more studied on a per-species basis. The numbers of papers addressing vertebrates and invertebrates has increased slightly and were roughly equivalent over time (Fig 3). The frequency of papers addressing "all taxa" from online databases has not changed significantly over time (Fig 3).

The most common data uses associated with the major taxonomic groups reflect the general maturity of data products associated with the respective group. Over 50% of vertebrate studies involved investigating species distribution (Fig 4); vertebrate data are generally more suitable for distribution studies because vertebrates are less diverse, many collections are completely digitized, and data for individual species are likely to contain sufficient numbers of records. Birds in particular have relatively good data available, in part because of online citizen science efforts and associated open data platforms, such as eBird [3]. While distribution studies were still the most common application across groups, significantly smaller percentages of plant (33%) and invertebrate (41%) studies dealt with species distribution. Plants and invertebrates are much more diverse, and the average species in these groups are less likely to have data of sufficient quantity and quality to estimate species distribution; however, growth in resources, especially for plants, is closing the gap. Data on insect distributions are less complete (or non-existent) for most species and hence may not be suitable for distribution and conservation studies [92,93].

A higher percentage of data papers, taxonomy, and barcoding papers involved invertebrates (Fig 4), reflecting in part the high taxonomic diversity for this group and need for more data. There are around 60,000 species of vertebrates, an estimated 400,000 plants, and an estimated 5–6 million species of insects; about one million insect species are currently described, which highlights the need for more taxonomic work in this group [20,94]. Other invertebrate phyla, such as Mollusca, are highly diverse as well (estimated 70,000–76,000 living species) [95]. Digitizing efforts for invertebrates have been particularly challenging, because many clades are so diverse, collections have much larger numbers of specimens, and the typically small specimens are difficult to digitize [96]. Automating digitization of such specimens, especially pinned insects and fluid-preserved invertebrates, faces significant obstacles [12,18,97–100].

The use of species occurrence data for conservation followed predicted trends. Vertebrate studies were more likely to address conservation; 23% of papers using vertebrate biodiversity records involved conservation, as compared to 14% of papers using plant records and 12% of papers using invertebrate records (Fig 4). Twenty percent of vertebrate species are currently classified as threatened, and that number is increasing [101]. While vertebrates have more data, they are by no means complete [102]; less-studied vertebrates (i.e. fish) also have much lower amounts of digitized data, as compared to birds [103]. Large species often receive more research and conservation funding, and very few conservation assessments exist for invertebrate taxa; most insect species are classified as "data deficient" (e.g. [104]). There is much need and potential for using primary biodiversity data to help determine conservation status of insects—perhaps starting with taxa known to be biological indicators of ecosystem health (e.g. [105,106]) and insects that provide important ecosystem services (e.g. [107]). However, identifying decline requires large numbers of records along with systematically collected surveys over time, which often do not exist for rare and potentially threatened species [108]. Opportunistic species occurrence records may therefore be best used to identify data gaps and promising areas for resurveys or standardized long-term monitoring studies when dealing with species decline [48].

Contrary to expectations, we found that studies addressing "all taxa" remained fairly consistent over time (Fig 3), and the maximum number of taxa addressed did not increase (Fig 5). However, this may be an effect of small sample sizes. Only four papers involved numbers of species in the hundreds of thousands over the period of 2010–2017 (Table 4). Most papers
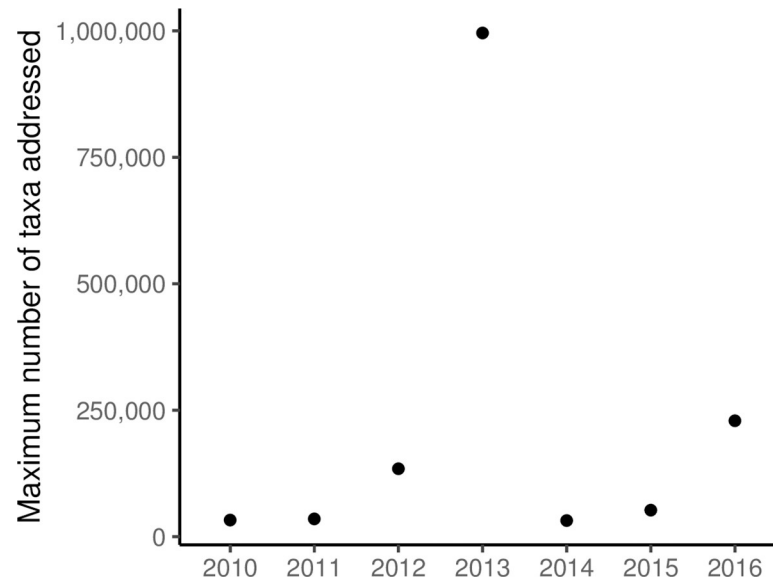
**Fig 5. Maximum number of taxa addressed in papers (*n* = 501) from 2010–2016.**

https://doi.org/10.1371/journal.pone.0215794.g005

focused on numbers of species in the single or double digits (Table 4). We found that the top data uses for papers that addressed "all taxa" involved data compilation and data quality (data quality assessments, data gap studies, data papers, and reporting on new databases, respectively). We argue that the scale of data that needs processing, along with issues of often sparse data, data obsolescence [109], and data of uncertain quality, make large-scale analyses challenging for anyone but a small group of data sciences-savvy end users. Additionally, effective large-scale assessments are often impossible without significant investments and active collaboration across study domains (e.g. taxonomy, ecology, biodiversity informatics) and geographical regions [110].

## Links to other data types

We determine how studies link primary biodiversity data to other data types by characterizing the variety of data compiled and used in each study (see S1 Table for full descriptions of 28 data linkage tags). We searched for information regarding other data types used within the methods section of each paper. Data link tags fall under four general categories of data types, including 1.) other types of occurrence data (i.e. data from literature, field surveys, species catalogues, private data); 2.) attributes of species occurrence data (e.g. information about the holding collections of specimens, species traits, conservation status, genetic data, associated image(s), species interactions, population data); 3.) environmental data (e.g. climate,

**Table 4. Number of taxa addressed by papers using online species occurrence records.**

| Number of taxa addressed | Number of papers |
| --- | --- |
| 1–9 | 113 |
| 10–99 | 106 |
| 100–999 | 82 |
| 1,000–9,999 | 68 |
| 10,000–99,999 | 22 |
| 100,000–999,999 | 4 |

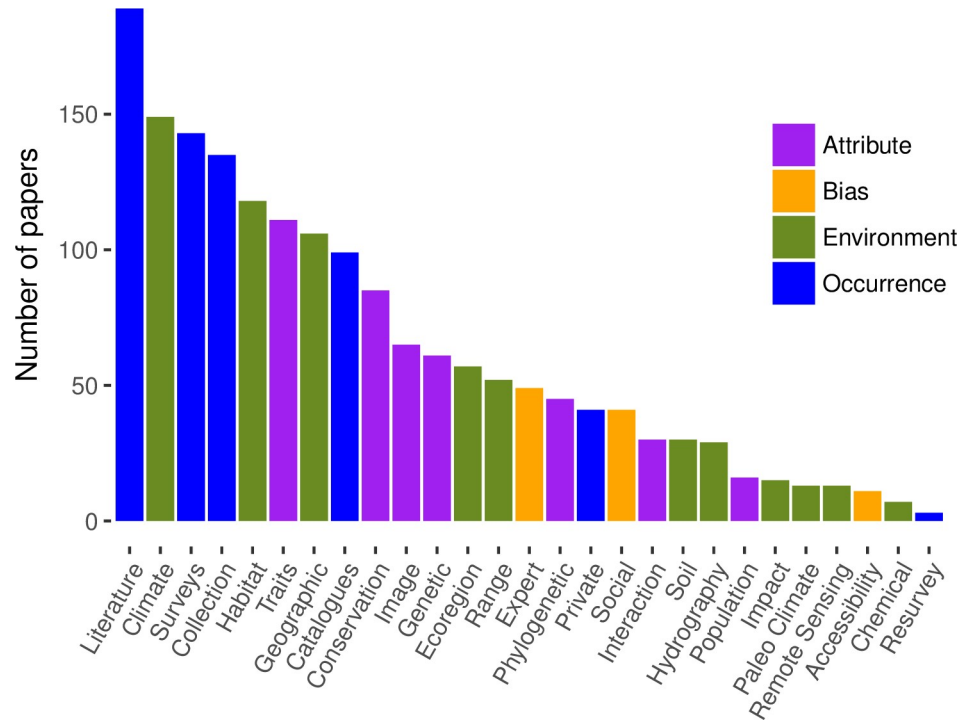https://doi.org/10.1371/journal.pone.0215794.t004

geographic information, habitat, ecoregion, etc.); and 4.) data that can be used to determine biases or gaps (socioeconomic data, expert knowledge, and accessibility of sites—with the last usually evaluated through proximity to roads or research institutions). We then determine the average number of data link tags associated with the six top uses, and the most common data type associated with each of these top uses.

Data types that were most often used in association with online species occurrence databases (out of 501 relevant papers) included occurrence records from previously published literature ($n = 189$), climate ($n = 149$), occurrence records from surveys ($n = 143$), collection information ($n = 135$), habitat ($n = 118$), traits ($n = 111$), and geographic data (e.g. elevation; $n = 106$, Fig 6). The only data types that changed over the time period of our dataset, 2010–2016, were collection, genetic, and phylogenetic data, which all increased (Fig 7). The average number of data linkages per paper was four (ranging from one to 11).

Table 5 summarizes top data linkages for different key uses. As predicted, climate is often a critical data type linked to occurrence records, especially for species distribution where it is the most commonly linked data type, and for diversity/population studies where it is a close second. For data papers and taxonomy studies, both collection data and literature data were often the most common data linkages. Conservation-focused studies most often linked occurrence records to conservation status, habitat, literature, and climatic data. Data quality studies often included a variety of data linkages, with little sorting of top linkages, likely representing the high dimensionality of data quality issues.
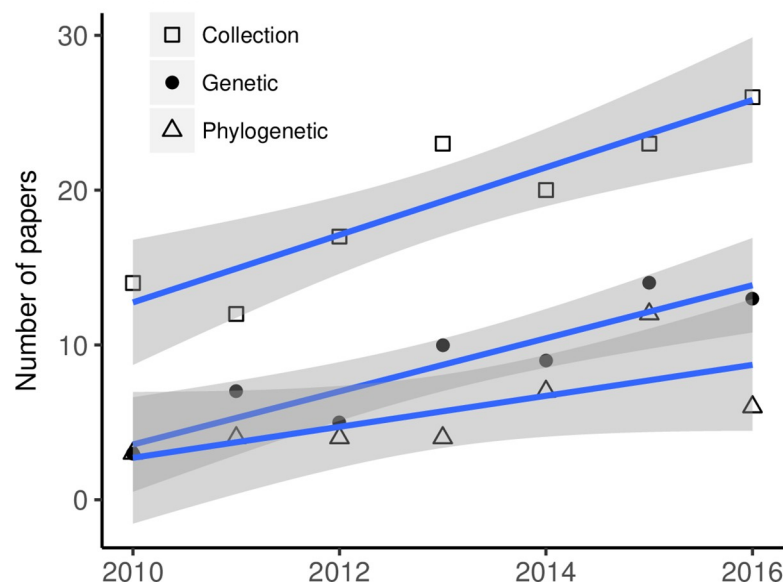
The high prevalence of studies compiling occurrence records from other sources indicates a continued demand for more and continued specimen sampling, and the need for more progress in getting these data into online databases (i.e. data papers and new database development). Three of the top five data types linked to online occurrence records included other types of occurrence data–literature-based occurrence data, surveys, and specimen data from natural history collections ($n = 189$, $n = 145$, and $n = 135$ papers used these data types, respectively). Sometimes the compiled data eventually make it into online data aggregators, such as GBIF, and sometimes they do not. Continued advocacy for data publication will be important to maximize the potential usability of all biodiversity data.

Environmental data used in conjunction with online biodiversity records are often applied in studies of species distribution. Specific environmental parameters used to predict distribution should be informed by expert knowledge of the requirements of a given species. Among environmental variables, climate data are perhaps the most readily available, relevant for the distribution of organisms on a global scale, and provide essential information for determining impacts of climate change on distribution [111,112]. Our data show that climate is indeed the most common environmental variable used in association with occurrence records (Fig 6; also documented in [56]). The second and third most common environmental data types used were geographic and habitat, which usually included GIS layers for elevation and land use and/or vegetation (see S1 Table). Elevation, land use, and vegetation data are also among the most readily available environmental data types, and are often relevant for evaluating species distribution at smaller spatial scales [113]. Despite increasing calls for incorporating relevant biotic interactions into models, only nine distribution studies incorporated data on interactions (i.e. competitive, consumptive, symbiotic, or pathogenic relationships), and 30 studies overall involved species interactions. The relatively low prevalence of species interaction information in these studies is thought to be primarily due to the large spatial scales usually considered in distribution models. Biotic interactions are often studied on a smaller scale by community ecologists, while distribution modeling is often done by macroecologists [114]. Primary species occurrences may provide needed data for studying biotic interactions on a larger scale, but these data are often not digitized, even if they exist in collections, and

**Fig 6. Number of papers that incorporate other data types to supplement or associate with online species occurrence records.** Data types fall within one of four categories, including 1.) attributes of occurrence information, 2.) data types that may help address bias in the data, 3.) environmental variables, and 4.) other kinds of occurrence data.

https://doi.org/10.1371/journal.pone.0215794.g006



**Fig 7. Data types linked to primary biodiversity data that increased over the period from 2010 through 2016.** These include data needed for taxonomic/phylogenetic studies, namely those from natural history specimens, genetic data, and phylogenetic data.

https://doi.org/10.1371/journal.pone.0215794.g007

**Table 5. Percentage of papers that associate online occurrence data with other data types—Separated by the six top uses of these databases.** Nine data types with the lowest percentages were removed from table. The top data type for each research use is bolded, and percentage values above 10% are highlighted yellow (10–29%), orange (30–49%), and red (>50%).

| Data Type | Species Distribution | Diversity/ Population | Data Paper | Taxonomy | Conservation | Data Quality |
|---|---|---|---|---|---|---|
| Climate | **58** | 37 | 7 | 2 | 32 | **26** |
| Literature | 41 | **40** | 29 | 52 | 40 | **26** |
| Geographic | 37 | 31 | 11 | 2 | 34 | 21 |
| Surveys | 36 | 36 | 29 | 32 | 32 | 13 |
| Habitat | 30 | 34 | 18 | 11 | 43 | 21 |
| Collection | 28 | 23 | **44** | 53 | 18 | 22 |
| Traits | 25 | 25 | 15 | 26 | 25 | 13 |
| Conservation | 20 | 29 | 9 | 15 | **75** | 15 |
| Expert | 15 | 7 | 9 | 3 | 22 | 7 |
| Private | 15 | 13 | 8 | 5 | 10 | 7 |
| Range | 14 | 12 | 6 | 5 | 22 | 13 |
| Catalogues | 11 | 18 | 20 | 25 | 19 | 22 |
| Hydrography | 11 | 12 | 3 | 2 | 16 | 1 |
| Soil | 11 | 11 | 2 | 0 | 10 | 3 |
| Ecoregion | 10 | 24 | 8 | 6 | 19 | 7 |
| Genetic | 10 | 13 | 24 | 26 | 6 | 6 |
| Social | 10 | 7 | 4 | 1 | 13 | 7 |
| Interaction | 9 | 5 | 4 | 8 | 6 | 0 |
| Paleo Climate | 7 | 5 | 1 | 0 | 1 | 0 |
| Image | 5 | 4 | 21 | 23 | 1 | 7 |
| Phylogenetic | 5 | 11 | 12 | 16 | 1 | 4 |

compiling data of sufficient quantity and quality for a given taxon remains an obstacle due to lack of automated data capture options for invertebrate collections.

The only data types that have increased over time were specimen collection, genetic, and phylogenetic data (Fig 7). We expected to see an increase in use of genetic data in particular, as these data are known to have expanded with the growth of databases, such as the Barcode of Life Data Systems (BOLDSystems) that links molecular, morphological, and distribution data [59]; the number of records in BOLDSystems increased from about 0.5 million in 2007 to 1.5 million today [60]. Further, large-scale phylogenetic resources, such as Open Tree of Life [115] that launched in 2015, have made it easier than ever before to phylogenies with other species data. The increasingly available collections, genetic, and phylogenetic data are highly relevant in taxonomy-related studies and data papers, which increased over time (Fig 2).

Both taxonomy and data papers used collection data most frequently in addition to data already available in online databases. Taxonomy-related uses of online species occurrence databases sometimes involve describing new species, but more commonly involve compilation of regional species checklists. The most traditional use of collections data is for taxonomy, so it is not surprising that over 50% of taxonomy papers also involve collections and literature data. The relatively high percentage of data papers that involve collections data (44%) reflects recent digitization efforts for natural history collections [1,9,13,116].

## Data quality

We characterize papers that address major data quality issues known to be associated with species occurrence data, including both common errors and biases. Data quality tags involve

improving data quality for a particular purpose addressed in the paper. Taxonomic nomenclature, species identification, spatial, and temporal data quality tags represent adjustments to the dataset used in a study that at least partially corrects the associated errors (see S1 Table). We also characterize studies that exclude certain inappropriate records, remove records with high georeferencing uncertainty, remove outliers, and those that address collection effort—see S1 Table. In addition to errors, some studies address specific biases known to be a problem in opportunistic datasets, including taxonomic, spatial, temporal, and environmental biases. Finally, we have a "detection" tag to represent use of statistical methods to estimate detection probability [53]. We assess the average number of quality tags associated with papers overall, and the most common data quality issues addressed within each of the top uses.

Overall, 69% of studies from our dataset that used online species occurrence records addressed one or more aspects of data quality. The biggest data quality concerns cited by users of primary biodiversity data in a recent survey [24] were georeference quality and taxonomic quality—we found that studies addressed these issues in 24% (spatial error in georeferences), 39% (taxonomic nomenclature), and 19% (species identifications) of published papers from our dataset (Table 6). Two data quality checks increased from 2010 to 2016: correcting taxonomic nomenclature and specimen identification (Fig 8), reflecting also the increase in taxonomy-related and data papers.
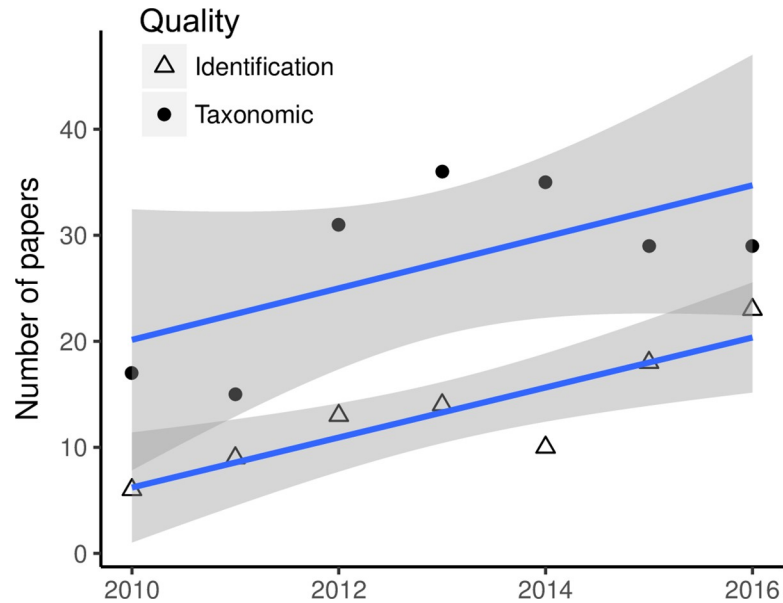
Spatial errors and taxonomic nomenclature are generally the easiest data quality errors to correct. Non-experts can check for spatial outliers or incorrect georeferences using standardized methods and online georeferencing tools [37,117]. Depending on data needs, one may also use existing uncertainty radii associated with georeferenced coordinates to select appropriate records for a study. However, most records in GBIF, for example, still do not have uncertainty radii; in a recent assessment of GBIF records for Odonata, Ephemeroptera, Plecoptera, and Trichoptera from the U.S.A., we found that the percentage of records with uncertainty radii associated with them was only 7–36% for these aquatic insect groups (as of April 2017). Of the 6.2 million catalogued molluscan lots in U.S. and Canadian collections, 4.5 million have undergone some form of data digitization. Of these, about 1.1 million (24%) of digitized records have been georeferenced, which represents 18% of all catalogued lots [49]. However, only a subset of these have uncertainty radii associated. Many digitization efforts for insects in particular have prioritized transcribing and publishing specimen label information and have not yet begun or completed georeferencing.

**Table 6. Papers from dataset (*n* = 501) that addressed data quality issues associated with species occurrence records.**

| Quality Tag | Number of Papers | Percentage |
|---|---|---|
| Taxonomic | 193 | 39% |
| Spatial | 121 | 24% |
| Identification | 94 | 19% |
| Spatial Bias | 59 | 12% |
| Exclusion | 57 | 11% |
| Effort | 50 | 10% |
| Precision | 30 | 6% |
| Temporal | 18 | 4% |
| Outliers | 17 | 3% |
| Temporal Bias | 11 | 2% |
| Taxonomic Bias | 9 | 2% |
| Environmental Bias | 6 | 1% |
| Detection | 4 | 1% |

**Fig 8. Number of papers that address identification errors and/or update taxonomic nomenclature from 2010–2016; note that these were the only two data quality issues that changed significantly over time.**

https://doi.org/10.1371/journal.pone.0215794.g008

Online taxonomic catalogues and tools to check records against updated catalogues are available for correcting taxonomic nomenclature [118,119]. However, we still have not reached the major goal of having online taxonomic data sources that are consistently updated by taxonomic experts for all species, although community-supported resources such as FishBase [65], WoRMS [120], and the latter's affiliated databases such as MilliBase [121], and MolluscaBase [122] are approaching that goal for many taxonomic groups. Other groups may lack online sources or have sources that are significantly out of date [123]. Unfortunately, the decline in resources devoted to the field of taxonomy does not bode well for achieving a unified taxonomic backbone usable for resolving all taxonomic issues [124,125]. Given the speed of taxonomic concept changes [126], lack of updated resources is a significant impediment to proper data integration. The best way for taxonomic experts to help ensure that nomenclature for their group is current is to engage with the community-supported and specialist-edited taxonomic database projects in their respective fields. The combined data of massive authority file efforts spanning multiple taxon groups, such as those covered by WoRMS, allow for novel approaches to data analysis [127].

Correcting species identifications requires taxonomic expertise for many organisms, particularly high-diversity groups, such as insects. Many users outside of the community of trained collection scientists may not understand or be interested in taxonomic concepts [1]. Therefore, despite misidentification being a well-known problem, this issue is less often directly addressed in papers. For those who are not taxonomic experts, some possible approaches to address misidentifications include: choosing taxonomic groups that are relatively easy to identify and less likely to have identification error, or including only records identified by reliable experts. For broad-scale biodiversity studies it may be appropriate to check occurrence locations against known ranges (where those exist); one may then identify outliers in the data where species are found in regions where they are not known to occur. Such efforts require both taxonomic and geospatial skills, although some automation may be possible [128].

Biases that result from variation in collection effort across space, time, taxonomic groups, and environments are also well-known problems in opportunistic biodiversity records [32,41,42,92]. The most commonly addressed bias in our dataset was spatial (addressed in 12% of papers, Table 7), as it is important for accurate species distribution modeling, and some methods to deal with spatial bias have been developed [41]. Other forms of bias were rarely addressed in only 1–2% of papers and include temporal bias (usually seasonal bias for certain times of year, or bias for certain years where specialists are active), taxonomic bias (e.g. preference for endangered species, charismatic taxa, avoiding common species or pests [47]), and environmental bias (e.g. preference for collecting in certain habitats or climates [41]).

Data quality issues are often dictated by the specific use. The most commonly checked data quality issues for papers involving species distribution were spatial errors (28% of distribution studies), taxonomic nomenclature (27%), spatial bias (24%), specimen identification (21%), and excluding inappropriate records (19%; Table 6). Taxonomic nomenclature was the most commonly checked data quality issue for all other top uses, ranging from 40% of papers (conservation and data quality uses) to 56% (taxonomy). In general, taxonomy papers only check issues related to nomenclature and identification. Data quality papers tend to focus evenly on the two most easily corrected issues (spatial and taxonomic, each 40% of data quality papers), followed by accounting for spatial bias (29% of data quality papers), effort (25%), and correcting specimen identification (18%). Diversity/population and conservation papers both also address taxonomic nomenclature and spatial errors most frequently (Table 7).

Automated data quality annotations are growing within the major online data aggregators (e.g. GBIF, iDigBio), but there is still much room to improve upon methods to easily tag data and highlight errors, biases, and uncertainty levels in the data. We need better methods to document confidence in data at a record and dataset level [23]. When data quality is addressed, it is usually done manually, and workflows are difficult to document, extend, and share. More recently, programs to automate and document data cleaning workflows have been developed, such as Kurator, a Kepler data curation package [38], but are not yet widely used due to the highly technical user interface, and have uncertain future support. Biodiversity databases allow efficient access to data that can expedite work, but care is still needed when using these resources. Data quality improvements on a large scale will require additional investment in data enhancements (e.g. collaborative georeferencing using standardized point-radius method) and quality control (e.g. efficiently identifying records that may need correction or attention from taxonomic experts).

## Conclusions and next steps

1. A high proportion of studies did not sufficiently cite databases, and many databases were no longer accessible at the time of this study; in most cases it was unclear whether the data were lost or moved to an aggregator. Continued efforts in data preservation and promoting best practices in data citation are essential for advancing scientific reproducibility, sustaining data resources, and encouraging publication of high-quality biodiversity data.

2. The increasing number of data papers over time reflects progress in digitization and online platforms for reporting observations through citizen science, as well as increases in journals that support data publication. Continued growth of data publications will enhance the efficiency and relevance of the field in addressing biodiversity conservation and environmental management.

3. Our study corroborated a recent bibliometric analysis of the larger field of biodiversity research, finding that more studies address plants (46% of studies using biodiversity

**Table 7. Percentage of papers that check aspects of data quality for online occurrence data—Separated by the six top uses of these databases.** Nine data types with the lowest percentages were removed from table. The top data type for each research use is bolded, and percentage values above 10% are highlighted yellow (10–29%), orange (30–49%), and red (>50%).

| Data Quality Check | Species Distribution | Diversity/ Population | Data Paper | Taxonomy | Conservation | Data Quality |
|---|---|---|---|---|---|---|
| Spatial | **28** | 27 | 26 | 9 | 29 | **40** |
| Taxonomic | 27 | **48** | **48** | **56** | 40 | 40 |
| Spatial Bias | 24 | 15 | 4 | 2 | 16 | 29 |
| Identification | 21 | 14 | 38 | 40 | 9 | 18 |
| Exclusion | 19 | 20 | 5 | 1 | 15 | 9 |
| Effort | 14 | 19 | 9 | 2 | 12 | 25 |
| Precision | 9 | 7 | 3 | 0 | 12 | 15 |
| Outliers | 5 | 1 | 1 | 1 | 3 | 10 |
| Temporal Bias | 4 | 3 | 2 | 1 | 1 | 4 |
| Temporal | 3 | 2 | 5 | 1 | 1 | 13 |
| Environmental Bias | 2 | 1 | 1 | 1 | 0 | 6 |
| Taxonomic Bias | 2 | 4 | 2 | 0 | 1 | 4 |
| Detection | 1 | 0 | 0 | 0 | 1 | 1 |

https://doi.org/10.1371/journal.pone.0215794.t007

databases) than vertebrates (25%) and invertebrates (25%). The prevalence of plants in studies that use online biodiversity databases may be due to a strong history of plant diversity work in Europe in particular, and the relative ease with which herbarium records can be digitized by scanning herbarium sheets.

4. While studies overall were less common for vertebrates than for plants, vertebrates may generally be more suitable for distribution studies because the group is less diverse, many collections are completely digitized, there are prolific citizen science communities reporting bird observations in particular, and data for individual species are more likely to contain sufficient numbers of records. Conservation studies are also more common for vertebrates, likely because they are disproportionately represented in threat assessments. In contrast, highly diverse invertebrates are more likely to be the subject of foundational biodiversity studies, such as taxonomy, barcoding, and data papers.

5. It is concerning that a relatively large proportion of studies does not explicitly address data quality—only 69% of studies in our dataset reported addressing one or more aspects of data quality. Authors who do address data quality are most likely to standardize nomenclature using online resources or to correct spatial errors. For nearly all uses of these data, there are errors and biases that can compromise results when using opportunistic records. Improving upon automated solutions to flag errors, and efficient mechanisms to report and correct data quality issues is critical in advancing the relevance and broadest use of this type of biodiversity data [129].

6. Significant investments in data enhancement and quality control are needed. This may be one limiting factor holding back studies that utilize all data currently held within biodiversity databases and studies that address very large numbers of taxa within clades. We found only four studies since 2010 that address hundreds of thousands of taxa, and most papers address numbers of taxa in the single or double digits. Large-scale improvements in data availability and fitness will require interdisciplinary effort and collaboration.

7. To limit the scope of the present paper, we focused efforts here on data citation, research uses, general taxa addressed, data linkages, and data quality issues addressed. However, we are also utilizing the dataset of tagged papers to address additional questions regarding

author connectedness and collaboration across institutions, countries, and disciplines. Such next-step efforts will provide additional context about the nature and scope of collaborations and resources that coalesce around digitally accessible primary biodiversity data.

## Supporting information

**S1 Table. Description of tags used to characterize papers, and number of papers assigned to each tag.**
(XLSX)

**S2 Table. Online biodiversity databases cited in published research and information on database scale, accessibility, and subject focus of the database (region, institution, and/or taxa included).**
(XLSX)

**S1 File. Paper metadata.** File in csv format containing citation information for 501 relevant journal articles analyzed in this review.
(XLSX)

## Acknowledgments

## Author Contributions

**Conceptualization:** Joan E. Ball-Damerow, Laura Brenskelle, Narayani Barve, Pamela S. Soltis, Arturo H. Ariño, Robert P. Guralnick.

**Data curation:** Joan E. Ball-Damerow, Laura Brenskelle, Narayani Barve, Robert P. Guralnick.

**Formal analysis:** Joan E. Ball-Damerow.

**Funding acquisition:** Joan E. Ball-Damerow.

**Investigation:** Joan E. Ball-Damerow, Robert P. Guralnick.

**Methodology:** Joan E. Ball-Damerow, Laura Brenskelle, Narayani Barve, Raphael LaFrance, Arturo H. Ariño, Robert P. Guralnick.

**Project administration:** Joan E. Ball-Damerow, Robert P. Guralnick.

**Software:** Raphael LaFrance.

**Supervision:** Petra Sierwald, Rüdiger Bieler, Robert P. Guralnick.

**Validation:** Joan E. Ball-Damerow, Robert P. Guralnick.

**Visualization:** Joan E. Ball-Damerow.

**Writing – original draft:** Joan E. Ball-Damerow.

**Writing – review & editing:** Joan E. Ball-Damerow, Laura Brenskelle, Pamela S. Soltis, Petra Sierwald, Rüdiger Bieler, Arturo H. Ariño, Robert P. Guralnick.

# References

1. Beaman R, Cellinese N. Mass digitization of scientific collections: New opportunities to transform the use of biological specimens and underwrite biodiversity science. ZooKeys. 2012; 209: 7–17. https://doi.org/10.3897/zookeys.209.3313 PMID: 22859875

2. Matsunaga A, Thompson A, Figueiredo RJ, Germain-Aubrey CC, Collins M, Beaman RS, et al. A Computational- and Storage-Cloud for Integration of Biodiversity Collections. 2013 IEEE 9th International Conference on e-Science. 2013. pp. 78–87. https://doi.org/10.1109/eScience.2013.48

3. Sullivan BL, Aycrigg JL, Barry JH, Bonney RE, Bruns N, Cooper CB, et al. The eBird enterprise: an integrated approach to development and application of citizen science. Biol Conserv. 2014; 169: 31–40.

4. Shaffer HB, Fisher RN, Davidson C. The role of natural history collections in documenting species declines. Trends Ecol Evol. 1998; 13: 27–30. https://doi.org/10.1016/S0169-5347(97)01177-4 PMID: 21238186

5. Ristaino JB. Tracking historic migrations of the Irish potato famine pathogen, Phytophthora infestans. Microbes Infect. 2002; 4: 1369–1377. https://doi.org/10.1016/S1286-4579(02)00010-2 PMID: 12443902

6. Suarez AV, Tsutsui ND. The Value of Museum Collections for Research and Society. BioScience. 2004; 54: 66–74. https://doi.org/10.1641/0006-3568(2004)054[0066:TVOMCF]2.0.CO;2

7. Graham CH, Ferrier S, Huettman F, Moritz C, Peterson AT. New developments in museum-based informatics and applications in biodiversity analysis. Trends Ecol Evol. 2004; 19: 497–503. https://doi.org/10.1016/j.tree.2004.07.006 PMID: 16701313

8. Pyke GH, Ehrlich PR. Biological collections and ecological/environmental research: a review, some observations and a look to the future. Biol Rev. 2010; 85: 247–266. https://doi.org/10.1111/j.1469-185X.2009.00098.x PMID: 19961469

9. Baird RC. Leveraging the fullest potential of scientific collections through digitisation. Biodivers Inform. 2010; 7. https://doi.org/10.17161/bi.v7i2.3987

10. GBIF [Internet]. [cited 5 Apr 2019]. Available: https://www.gbif.org/

11. Baker B. New Push to Bring US Biological Collections to the World's Online Community Advances in technology put massive undertaking within reach. BioScience. 2011; 61: 657–662. https://doi.org/10.1525/bio.2011.61.9.4

12. Blagoderov V, Kitching I, Livermore L, Simonsen T, Smith V. No specimen left behind: industrial scale digitization of natural history collections. ZooKeys. 2012; 209: 133–146. https://doi.org/10.3897/zookeys.209.3178 PMID: 22859884

13. Page LM, MacFadden BJ, Fortes JA, Soltis PS, Riccardi G. Digitization of Biodiversity Collections Reveals Biggest Data on Biodiversity. BioScience. 2015; 65: 841–842. https://doi.org/10.1093/biosci/biv104

14. Ariño AH. Approaches to estimating the universe of natural history collections data. Biodivers Inform. 2010;7. https://doi.org/10.17161/bi.v7i2.3991

15. Ariño A. Putting your Finger upon the Simplest Data. Biodivers Inf Sci Stand. 2018; 2: e26300. https://doi.org/10.3897/biss.2.26300

16. Nelson G, Paul D, Riccardi G, Mast A. Five task clusters that enable efficient and effective digitization of biological collections. ZooKeys. 2012; 209: 19–45. https://doi.org/10.3897/zookeys.209.3135 PMID: 22859876

17. Tulig M, Tarnowsky N, Bevans M, Kirchgessner A, Thiers B. Increasing the efficiency of digitization workflows for herbarium specimens. ZooKeys. 2012; 209: 103–113. https://doi.org/10.3897/zookeys.209.3125 PMID: 22859882

18. Hudson LN, Blagoderov V, Heaton A, Holtzhausen P, Livermore L, Price BW, et al. Inselect: Automating the Digitization of Natural History Collections. PLOS ONE. 2015; 10: e0143402. https://doi.org/10.1371/journal.pone.0143402 PMID: 26599208

19. Allan EL, Livermore L, Price B, Shchedrina O, Smith V. A Novel Automated Mass Digitisation Workflow for Natural History Microscope Slides. Biodivers Data J. 2019; 7: e32342. https://doi.org/10.3897/BDJ.7.e32342 PMID: 30863197

20. Pimm SL, Jenkins CN, Abell R, Brooks TM, Gittleman JL, Joppa LN, et al. The biodiversity of species and their rates of extinction, distribution, and protection. Science. 2014; 344: 1246752. https://doi.org/10.1126/science.1246752 PMID: 24876501

21. Alroy J. Current extinction rates of reptiles and amphibians. Proc Natl Acad Sci. 2015; 112: 13003–13008. https://doi.org/10.1073/pnas.1508681112 PMID: 26438855

**22.** Régnier C, Achaz G, Lambert A, Cowie RH, Bouchet P, Fontaine B. Mass extinction in poorly known taxa. Proc Natl Acad Sci. 2015; 112: 7761–7766. https://doi.org/10.1073/pnas.1502350112 PMID: 26056308

**23.** Faith D, Collen B, Ariño A, Koleff PKP, Guinotte J, Kerr J, et al. Bridging the biodiversity data gaps: Recommendations to meet users' data needs. Biodivers Inform. 2013; 8. Available: https://journals.ku.edu/index.php/jbi/article/view/4126

**24.** Ariño AH, Chavan V, Faith DP. Assessment of user needs of primary biodiversity data: Analysis, concerns, and challenges. Biodivers Inform. 2013; 8. https://doi.org/10.17161/bi.v8i2.4094

**25.** Guralnick R, Hill A. Biodiversity informatics: automated approaches for documenting global biodiversity patterns and processes. Bioinformatics. 2009; 25: 421–428. https://doi.org/10.1093/bioinformatics/btn659 PMID: 19129210

**26.** Sousa-Baena MS, Garcia LC, Peterson AT. Knowledge behind conservation status decisions: data basis for "Data Deficient" Brazilian plant species. Biol Conserv. 2014; 173: 80–89.

**27.** Feeley K. Are We Filling the Data Void? An Assessment of the Amount and Extent of Plant Collection Records and Census Data Available for Tropical South America. PLOS ONE. 2015; 10: 1–17. https://doi.org/10.1371/journal.pone.0125629 PMID: 25927831

**28.** Meyer C, Kreft H, Guralnick R, Jetz W. Global priorities for an effective information basis of biodiversity distributions. Nat Commun. 2015; 6. https://doi.org/10.1038/ncomms9221 PMID: 26348291

**29.** Beck J, Ballesteros-Mejia L, Buchmann CM, Dengler J, Fritz SA, Gruber B, et al. What's on the horizon for macroecology? Ecography. 2012; 35: 673–683. https://doi.org/10.1111/j.1600-0587.2012.07364.x

**30.** Beck J, Ballesteros-Mejia L, Nagel P, Kitching IJ. Online solutions and the Wallacean shortfall what does GBIF contribute to our knowledge of species ranges? Divers Distrib. 2013; 19: 1043–1050.

**31.** Peterson AT, Asase A, Canhos D, Souza S de, Wieczorek J. Data Leakage and Loss in Biodiversity Informatics. Biodivers Data J. 2018; 6: e26826. https://doi.org/10.3897/BDJ.6.e26826 PMID: 30473617

**32.** Daru BH, Park DS, Primack RB, Willis CG, Barrington DS, Whitfeld TJS, et al. Widespread sampling biases in herbaria revealed from large-scale digitization. New Phytol. 2018; 217: 939–955. https://doi.org/10.1111/nph.14855 PMID: 29083043

**33.** Maldonado C, Molina CI, Zizka A, Persson C, Taylor CM, Alban J, et al. Estimating species diversity and distribution in the era of Big Data: to what extent can we trust public databases? Glob Ecol Biogeogr. 2015; 24: 973–984. https://doi.org/10.1111/geb.12326 PMID: 27656106

**34.** Meier R, Dikow T. Significance of Specimen Databases from Taxonomic Revisions for Estimating and Mapping the Global Species Diversity of Invertebrates and Repatriating Reliable Specimen Data. Conserv Biol. 2004; 18: 478–488. https://doi.org/10.1111/j.1523-1739.2004.00233.x

**35.** Goodwin ZA, Harris DJ, Filer D, Wood JRI, Scotland RW. Widespread mistaken identity in tropical plant collections. Curr Biol CB. 2015; 25: R1066–1067. https://doi.org/10.1016/j.cub.2015.10.002 PMID: 26583892

**36.** Zermoglio PF, Guralnick RP, Wieczorek JR. A Standardized Reference Data Set for Vertebrate Taxon Name Resolution. PLOS ONE. 2016; 11: e0146894. https://doi.org/10.1371/journal.pone.0146894 PMID: 26760296

**37.** Wieczorek J, Guo Q, Hijmans R. The point-radius method for georeferencing locality descriptions and calculating associated uncertainty. Int J Geogr Inf Sci. 2004; 18: 745–767. https://doi.org/10.1080/13658810412331280211

**38.** Dou L, Cao G, Morris PJ, Morris RA, Ludäscher B, Macklin JA, et al. Kurator: A Kepler package for data curation workflows. Procedia Comput Sci. 2012; 9: 1614–1619. https://doi.org/10.1016/j.procs.2012.04.177

**39.** Mathew C, Güntsch A, Obst M, Vicario S, Haines R, Williams A, et al. A semi-automated workflow for biodiversity data retrieval, cleaning, and quality control. Biodivers Data J. 2014; 2: 1–12.

**40.** Ponder W, Carter G, Flemons P, R. Chapman R. Evaluation of Museum Collection Data for Use in Biodiversity Assessment. Conserv Biol. 2001;15. https://doi.org/10.1046/j.1523-1739.2001.015003648.x

**41.** Boakes EH, McGowan PJ, Fuller RA, Chang-qing D, Clark NE, O'Connor K, et al. Distorted views of biodiversity: spatial and temporal bias in species occurrence data. PLOS Biol. 2010; 8: e1000385. https://doi.org/10.1371/journal.pbio.1000385 PMID: 20532234

**42.** Isaac NJ, Strien AJ, August TA, Zeeuw MP, Roy DB. Statistics for citizen science: extracting signals of change from noisy ecological data. Methods Ecol Evol. 2014; 5: 1052–1060.

**43.** Ruete A. Displaying bias in sampling effort of data accessed from biodiversity databases using ignorance maps. Biodivers Data J. 2015; 1–15.

44. Meyer C, Weigelt P, Kreft H. Multidimensional biases, gaps and uncertainties in global plant occurrence information. Ecol Lett. 2016; 19: 992–1006. https://doi.org/10.1111/ele.12624 PMID: 27250865

45. Meyer C, Jetz W, Guralnick RP, Fritz SA, Kreft H. Range geometry and socio-economics dominate species-level biases in occurrence information. Glob Ecol Biogeogr. 2016; 25: 1181–1193. https://doi.org/10.1111/geb.12483

46. Guralnick R, Van Cleve J. Strengths and weaknesses of museum and national survey data sets for predicting regional species richness: comparative and combined approaches. Divers Distrib. 2005; 11: 349–359. https://doi.org/10.1111/j.1366-9516.2005.00164.x

47. Ball-Damerow JE, Oboyski PT, Resh VH. California dragonfly and damselfly (Odonata) database: temporal and spatial distribution of species records collected over the past century. ZooKeys. 2015; 67.

48. Rapacciuolo G, Ball-Damerow JE, Zeilinger AR, Resh VH. Detecting long-term occupancy changes in Californian odonates from natural history and citizen science records. Biodivers Conserv. 2017; 26: 2933–2949. https://doi.org/10.1007/s10531-017-1399-4

49. Sierwald P, Bieler R, Shea EK, Rosenberg G. Mobilizing Mollusks: Status Update on Mollusk Collections in the U.S.A. and Canada. Am Malacol Bull. 2018; 36: 177–214. https://doi.org/10.4003/006.036.0202

50. ter Steege H, A. Persaud C. The phenology of Guyanese timber species—A compilation of a century of observations. Plant Ecol. 1991; 95: 177–198. https://doi.org/10.1007/BF00045216

51. Peterson CH. Relative abundances of living and dead molluscs in two Californian lagoons. Lethaia. 1976; 9: 137–148. https://doi.org/10.1111/j.1502-3931.1976.tb00958.x

52. Boag DA. Overcoming sampling bias in studies of terrestrial gastropods. Can J Zool. 1982; 60: 1289–1292. https://doi.org/10.1139/z82-173

53. Dorazio RM. Accounting for imperfect detection and survey bias in statistical analysis of presence-only data. Glob Ecol Biogeogr. 2014; 23: 1472–1484. https://doi.org/10.1111/geb.12216

54. Zeilinger AR, Rapacciuolo G, Turek D, Oboyski PT, Almeida RPP, Roderick GK. Museum specimen data reveal emergence of a plant disease may be linked to increases in the insect vector population. Ecol Appl Publ Ecol Soc Am. 2017; 27: 1827–1837. https://doi.org/10.1002/eap.1569 PMID: 28459124

55. Chapman AD. Uses of Primary Species-Occurrence Data, version 1.0. Report for the Global Biodiversity Information Facility. [Internet]. Copenhagen; 2005. Available: Http://www.gbif.org/orc/?doc_id=1300.

56. Ariño A, Noesgaard D, Hjarding A, Schigel D. Biodiversity Information Services: A (not-so-) little knowledge that acts. Biodivers Inf Sci Stand. 2018; 2: e25738. https://doi.org/10.3897/biss.2.25738

57. Roy Rosenzweig Center for History and New Media. Zotero [Internet]. 2017. Available: www.zotero.org/download

58. Ball-Damerow JE, Brenskelle L, Barve N, LaFrance R, Soltis PS, Sierwald P, et al. Bibliographic dataset characterizing studies that use online biodiversity databases [Internet]. Zenodo; 2019. https://doi.org/10.5281/zenodo.2589439

59. Ratnasingham S, Hebert PDN. bold: The Barcode of Life Data System (http://www.barcodinglife.org). Mol Ecol Notes. 2007; 7: 355–364. https://doi.org/10.1111/j.1471-8286.2007.01678.x PMID: 18784790

60. BOLDSystems v4 [Internet]. [cited 5 Apr 2019]. Available: http://www.boldsystems.org/

61. speciesLink: Sistema de Informação Distribuído para Coleções Biológicas [Internet]. 2019 [cited 8 Jun 2019]. Available: http://splink.cria.org.br/

62. Ocean Biogeographic Information System [Internet]. 2019 [cited 8 Jun 2019]. Available: https://obis.org/

63. AVH | The Australasian Virtual Herbarium [Internet]. [cited 8 Jun 2019]. Available: https://avh.chah.org.au/

64. Tropicos—Home [Internet]. 2019 [cited 8 Jun 2019]. Available: https://www.tropicos.org/

65. Froese R, Pauly D. FishBase. World Wide Web electronic publication. 2014; Available: https://www.scienceopen.com/document?vid=dc419213-0ca3-48cc-901c-2934ecf4441e

66. FishBase [Internet]. 2019 [cited 8 Jun 2019]. Available: https://www.fishbase.in/search.php

67. Hendrickson DA, Cohen AE. Fishes of Texas Project Database (Version 2.0) [Internet]. 1 Sep 2015 [cited 8 Jun 2019]. Available: http://www.fishesoftexas.org/documentation/

68. Collections of the REMIB [Internet]. [cited 8 Jun 2019]. Available: http://www.conabio.gob.mx/remib_ingles/doctos/remibnodosdb.html?

69. Chavan V, Penev L. The data paper: a mechanism to incentivize data publishing in biodiversity science. BMC Bioinformatics. 2011; 12: S2. https://doi.org/10.1186/1471-2105-12-S15-S2 PMID: 22373175

70. Moritz T, Krishnan S, Roberts D, Ingwersen P, Agosti D, Penev L, et al. Towards mainstreaming of biodiversity data publishing: recommendations of the GBIF Data Publishing Framework Task Group. BMC Bioinformatics. 2011; 12: S1. https://doi.org/10.1186/1471-2105-12-S15-S1 PMID: 22373150

71. Whitlock MC. Data archiving in ecology and evolution: best practices. Trends Ecol Evol. 2011; 26: 61–65. https://doi.org/10.1016/j.tree.2010.11.006 PMID: 21159406

72. Smith V, Penev L. E-Infrastructures for Data Publishing in Biodiversity Science. PenSoft Publishers LTD; 2011.

73. Costello MJ, Michener WK, Gahegan M, Zhang Z-Q, Bourne PE. Biodiversity data should be published, cited, and peer reviewed. Trends Ecol Evol. 2013; 28: 454–461. https://doi.org/10.1016/j.tree.2013.05.002 PMID: 23756105

74. Costello MJ, Wieczorek J. Best practice for biodiversity data management and publication. Biol Conserv. 2014; 173: 68–73. https://doi.org/10.1016/j.biocon.2013.10.018

75. Wilkinson MD, Dumontier M, Aalbersberg IjJ, Appleton G, Axton M, Baak A, et al. The FAIR Guiding Principles for scientific data management and stewardship. Sci Data. 2016; 3: 160018. https://doi.org/10.1038/sdata.2016.18 PMID: 26978244

76. Mooney H, Newton M. The Anatomy of a Data Citation: Discovery, Reuse, and Credit. J Librariansh Sch Commun. 2012; 1: eP1035. https://doi.org/10.7710/2162-3309.1035

77. Escribano N, Galicia D, Ariño AH. The tragedy of the biodiversity data commons: a data impediment creeping nigher? Database J Biol Databases Curation. 2018;2018. https://doi.org/10.1093/database/bay033 PMID: 29688384

78. Vines TH, Albert AYK, Andrew RL, Débarre F, Bock DG, Franklin MT, et al. The Availability of Research Data Declines Rapidly with Article Age. Curr Biol. 2014; 24: 94–97. https://doi.org/10.1016/j.cub.2013.11.014 PMID: 24361065

79. Klump J, Huber R. 20 Years of Persistent Identifiers–Which Systems are Here to Stay? Data Sci J. 2017; 16: 9. https://doi.org/10.5334/dsj-2017-009

80. McMurry JA, Juty N, Blomberg N, Burdett T, Conlin T, Conte N, et al. Identifiers for the 21st century: How to design, provision, and reuse persistent identifiers to maximize utility and impact of life science data. PLOS Biol. 2017; 15: e2001414. https://doi.org/10.1371/journal.pbio.2001414 PMID: 28662064

81. Stark PB. Before reproducibility must come preproducibility. Nature. 2018; 557: 613. https://doi.org/10.1038/d41586-018-05256-0 PMID: 29795524

82. Cousijn H, Kenall A, Ganley E, Harrison M, Kernohan D, Lemberger T, et al. A data citation roadmap for scientific publishers. Sci Data. 2018; 5: 180259. https://doi.org/10.1038/sdata.2018.259 PMID: 30457573

83. Force MM, Robinson NJ. Encouraging data citation and discovery with the Data Citation Index. J Comput Aided Mol Des. 2014; 28: 1043–1048. https://doi.org/10.1007/s10822-014-9768-5 PMID: 24980647

84. Costello MJ, Appeltans W, Bailly N, Berendsohn WG, de Jong Y, Edwards M, et al. Strategies for the sustainability of online open-access biodiversity databases. Biol Conserv. 2014; 173: 155–165.

85. Huang X, Hawkins BA, Qiao G. Biodiversity data sharing: Will peer-reviewed data papers work? BioScience. 2013; 63: 5–6.

86. Pimm SL, Alibhai S, Bergl R, Dehgan A, Giri C, Jewell Z, et al. Emerging Technologies to Conserve Biodiversity. Trends Ecol Evol. 2015; 30: 685–696. https://doi.org/10.1016/j.tree.2015.08.008 PMID: 26437636

87. Wood KR. Rediscovery, conservation status and taxonomic assessment of Melicope degeneri (Rutaceae), Kaua 'i, Hawai 'i. Endanger Species Res. 2011; 14: 61–68.

88. Costello MJ. Motivating Online Publication of Data. BioScience. 2009; 59: 418–427. https://doi.org/10.1525/bio.2009.59.5.9

89. Costello MJ, Bouchet P, Boxshall G, Fauchald K, Gordon D, Hoeksema BW, et al. Global coordination and standardisation in marine biodiversity through the World Register of Marine Species (WoRMS) and related databases. PLOS ONE. 2013; 8: e51629. https://doi.org/10.1371/journal.pone.0051629 PMID: 23505408

90. Tydecks L, Jeschke JM, Wolf M, Singer G, Tockner K. Spatial and topical imbalances in biodiversity research. PLOS ONE. 2018; 13: e0199327. https://doi.org/10.1371/journal.pone.0199327 PMID: 29975719

91. Chapman AD. Numbers of Living Species in Australia and the World: A Report for the Australian Biological Resources Study [Internet]. Toowoomba, Australia: Australian Government Department of the Environment and Energy; 2009. Report No.: ISBN: 978 0 642 56861 8. Available: http://www.environment.gov.au/science/abrs/publications/other/numbers-living-species/contents#copyright

92. Sánchez-Fernández D, Lobo JM, Abellán P, Ribera I, Millán A. Bias in freshwater biodiversity sampling: the case of Iberian water beetles. Divers Distrib. 2008; 14: 754–762. https://doi.org/10.1111/j.1472-4642.2008.00474.x

93. Ballesteros-Mejia L, Kitching IJ, Jetz W, Nagel P, Beck J. Mapping the biodiversity of tropical insects: species richness and inventory completeness of African sphingid moths. Glob Ecol Biogeogr. 2013; 22: 586–595. https://doi.org/10.1111/geb.12039

94. Costello MJ, Wilson S, Houlding B. Predicting total global species richness using rates of species description and estimates of taxonomic effort. Syst Biol. 2012; 61: 871–883. https://doi.org/10.1093/sysbio/syr080 PMID: 21856630

95. Rosenberg G. A New Critical Estimate of Named Species-Level Diversity of the Recent Mollusca*. Am Malacol Bull. 2014; 32: 308–322. https://doi.org/10.4003/006.032.0204

96. Schuh RT, Hewson-Smith S, Ascher JS. Specimen databases: A case study in entomology using web-based software. Am Entomol. 2010; 56: 206–216.

97. Mantle B, LaSalle J, Fisher N. Whole-drawer imaging for digital management and curation of a large entomological collection. ZooKeys. 2012; 209: 147–163. https://doi.org/10.3897/zookeys.209.3169 PMID: 22859885

98. Holovachov O, Zatushevsky A, Shydlovsky I. Whole-Drawer Imaging of Entomological Collections: Benefits, Limitations and Alternative Applications. J Conserv Mus Stud. 2014; 12: Art. 9. https://doi.org/10.5334/jcms.1021218

99. Hereld M, Ferrier NJ, Agarwal N, Sierwald P. Designing a High-Throughput Pipeline for Digitizing Pinned Insects. 2017 IEEE 13th International Conference on e-Science (e-Science). 2017. pp. 542–550. https://doi.org/10.1109/eScience.2017.88

100. Price BW, Dupont S, Allan EL, Blagoderov V, Butcher AJ, Durrant J, et al. ALICE: Angled Label Image Capture and Extraction for high throughput insect specimen digitisation. 2018; None

101. Hoffmann M, Hilton-Taylor C, Angulo A, Böhm M, Brooks TM, Butchart SHM, et al. The Impact of Conservation on the Status of the World's Vertebrates. Science. 2010; 330: 1503–1509. https://doi.org/10.1126/science.1194442 PMID: 20978281

102. Pino-del-Carpio A, Ariño AH, Miranda R. Data exchange gaps in knowledge of biodiversity: implications for the management and conservation of Biosphere Reserves. Biodivers Conserv. 2014; 23: 2239–2258.

103. Pino-Del-Carpio A, Villarroya A, Ariño AH, Puig J, Miranda R. Communication gaps in knowledge of freshwater fish biodiversity: implications for the management and conservation of Mexican biosphere reserves. J Fish Biol. 2011; 79: 1563–1591. https://doi.org/10.1111/j.1095-8649.2011.03073.x PMID: 22136240

104. Ball J, Beche L, Mendez P, H. Resh V. Biodiversity in Mediterranean-climate streams of California. Hydrobiologia. 2013; 719. https://doi.org/10.1007/s10750-012-1368-6

105. Dewalt E, Favret C, W. Webb D. Just how imperiled are aquatic insects? A case study of stoneflies (Plecoptera) in Illinois. Ann Entomol Soc Am. 2005; 98: 941–950. https://doi.org/10.1603/0013-8746(2005)098[0941:JHIAAI]2.0.CO;2

106. Ball-Damerow JE, M'Gonigle LK, Resh VH. Changes in occurrence, richness, and biological traits of dragonflies and damselflies (Odonata) in California and Nevada over the past century. Biodivers Conserv. 2014; 23: 2107–2126. https://doi.org/10.1007/s10531-014-0707-5

107. Colla SR, Gadallah F, Richardson L, Wagner D, Gall L. Assessing declines of North American bumble bees (Bombus spp.) using museum specimens. Biodivers Conserv. 2012; 21: 3585–3595.

108. Hallmann CA, Sorg M, Jongejans E, Siepel H, Hofland N, Schwan H, et al. More than 75 percent decline over 27 years in total flying insect biomass in protected areas. PLOS ONE. 2017; 12: e0185809. https://doi.org/10.1371/journal.pone.0185809 PMID: 29045418

109. Escribano N, Ariño AH, Galicia D. Biodiversity data obsolescence and land uses changes. PeerJ. 2016; 4: 1–15.

110. Peterson AT, Soberón J, Krishtalka L. A global perspective on decadal challenges and priorities in biodiversity informatics. BMC Ecol. 2015; 15: 15. https://doi.org/10.1186/s12898-015-0046-8 PMID: 26022532

111. Austin M, Van Niel K. Improving species distribution models for climate change studies: Variable selection and scale. J Biogeogr. 2010; 38: 1–8. https://doi.org/10.1111/j.1365-2699.2010.02416.x

112.    Stanton JC, Pearson RG, Horning N, Ersts P, Reşit Akçakaya H. Combining static and dynamic variables in species distribution models under climate change. Methods Ecol Evol. 2012; 3: 349–357.

113.    Fournier A, Barbet-Massin M, Rome Q, Courchamp F. Predicting species distribution combining multiscale drivers. Glob Ecol Conserv. 2017; 12: 215–226. https://doi.org/10.1016/j.gecco.2017.11.002

114.    Staniczenko PPA, Sivasubramaniam P, Suttle KB, Pearson RG. Linking macroecology and community ecology: refining predictions of species distributions using biotic interaction networks. Ecol Lett. 2017; 20: 693–707. https://doi.org/10.1111/ele.12770 PMID: 28429842

115.    Hinchliff CE, Smith SA, Allman JF, Burleigh JG, Chaudhary R, Coghill LM, et al. Synthesis of phylogeny and taxonomy into a comprehensive tree of life. Proc Natl Acad Sci. 2015; 112: 12764–12769. https://doi.org/10.1073/pnas.1423041112 PMID: 26385966

116.    Chavan V, Berents P, Hamer M. Towards demand driven publishing: approaches to the prioritisation of digitisation of natural history collections data. Biodivers Inform. 2010;7. https://doi.org/10.17161/bi.v7i2.3990

117.    Rios NE, Bart HL. GEOLocate. Belle Chasse, LA: Tulane University Museum of Natural History; 2018.

118.    Boyle B, Hopkins N, Lu Z, Raygoza Garay JA, Mozzherin D, Rees T, et al. The taxonomic name resolution service: an online tool for automated standardization of plant names. BMC Bioinformatics. 2013; 14: 16. https://doi.org/10.1186/1471-2105-14-16 PMID: 23324024

119.    Chamberlain SA, Szöcs E. taxize: taxonomic search and retrieval in R. F1000Research. 2013; 2. https://doi.org/10.12688/f1000research.2-191.v2 PMID: 24555091

120.    WoRMS Editorial Board. World Register of Marine Species. Available from http://www.marinespecies.org at VLIZ. Accessed yyyy-mm-dd. [Internet]. VLIZ; 2017. doi:10.14284/170

121.    MilliBase [Internet]. [cited 5 Apr 2019]. Available: http://www.millibase.org/

122.    MolluscaBase—Introduction [Internet]. [cited 5 Apr 2019]. Available: http://www.molluscabase.org/

123.    Ball-Damerow JE, Mendez PK, Sierwald P, Bieler R, Yoder M, DeWalt RE. Taxonomic data quality in GBIF: a case study of aquatic macroinvertebrate groups. Ann Arbor, MI; 2017.

124.    Wägele H, Klussmann-Kolb A, Kuhlmann M, Haszprunar G, Lindberg D, Koch A, et al. The taxonomist —an endangered race. A practical proposal for its survival. Front Zool. 2011; 8: 25. https://doi.org/10.1186/1742-9994-8-25 PMID: 22029904

125.    Drew LW. Are We Losing the Science of Taxonomy?: As need grows, numbers and training are failing to keep up. BioScience. 2011; 61: 942–946. https://doi.org/10.1525/bio.2011.61.12.4

126.    Vaidya G, Lepage D, Guralnick R. The tempo and mode of the taxonomic correction process: How taxonomists have corrected and recorrected North American bird species over the last 127 years. PLoS ONE. 2018; 13. https://doi.org/10.1371/journal.pone.0195736 PMID: 29672539

127.    Arvanitidis CD, Warwick RM, Somerfield PJ, Pavloudi C, Pafilis E, Oulas A, et al. Research Infrastructures offer capacity to address scientific questions never attempted before: Are all taxa equal? [Internet]. PeerJ Inc.; 2018 Aug. Report No.: e26819v2. https://doi.org/10.7287/peerj.preprints.26819v2

128.    Otegui J, Guralnick RP. The geospatial data quality REST API for primary biodiversity data. Bioinformatics. 2016; 32: 1755–1757. https://doi.org/10.1093/bioinformatics/btw057 PMID: 26833340

129.    Paul D, Fisher N. Challenges For Implementing Collections Data Quality Feedback: synthesizing the community experience. Biodivers Inf Sci Stand. 2018; 2: e26003. https://doi.org/10.3897/biss.2.26003