# State of biodiversity documentation in the Philippines: Metadata gaps, taxonomic biases, and spatial biases in the DNA barcode data of animal and plant taxa in the context of species occurrence data

**Carmela Maria Berba** [1], **Ambrocio Melvin Matias** [Corresp. 1]

[1] Institute of Biology, University of the Philippines Diliman, Quezon City, National Capital Region, Philippines

Corresponding Author: Ambrocio Melvin Matias
Email address: aamatias@up.edu.ph

Anthropogenic changes in the natural environment have led to alarming rates of biodiversity loss, resulting in a more urgent need for conservation. Although there is an increasing cognizance of the importance of incorporating biodiversity data into conservation, the accuracy of the inferences generated from these records can be highly impacted by gaps and biases in the data. Because of the Philippines' status as a biodiversity hotspot, the assessment of potential gaps and biases in biodiversity documentation in the country can be a critical step in the identification of priority research areas for conservation applications. In this study, we systematically assessed biodiversity data on animal and plant taxa found in the Philippines by examining the extent of metadata gaps, taxonomic biases, and spatial biases in DNA barcode data while using species occurrence data as a backdrop of the Philippines' biodiversity. These barcode and species occurrence data sets were obtained from public databases, namely: GenBank, Barcode of Life Data System and Global Biodiversity Information Facility. We found that much of the barcode data had missing information on either records and publishing, geolocation, or taxonomic metadata, which consequently, can limit the usability of barcode data for further analyses. We also observed that the amount of barcode data can be directly associated with the amount of species occurrence data available for a particular taxonomic group and location – highlighting the potential sampling biases in the barcode data. While the majority of barcode data came from foreign institutions, there has been an increase in local efforts in recent decades. However, much of the contribution to biodiversity documentation only come from institutions based in Luzon.

1 **State of biodiversity documentation in the Philippines:**

2 **metadata gaps, taxonomic biases, and spatial biases**

3 **in the DNA barcode data of animal and plant taxa in**

4 **the context of species occurrence data**

5

6 Carmela Maria P. Berba[1], Ambrocio Melvin A. Matias[1]

7 [1] Institute of Biology, University of the Philippines Diliman, Quezon City, NCR, Philippines

8

9 Corresponding Author:

10 Ambrocio Melvin A. Matias[1]

11 Institute of Biology, University of the Philippines Diliman, Quezon City, NCR, 1101 Philippines

12 Email address: aamatias@up.edu.ph

## Abstract

13

14         Anthropogenic changes in the natural environment have led to alarming rates of

15  biodiversity loss, resulting in a more urgent need for conservation. Although there is an

16  increasing cognizance of the importance of incorporating biodiversity data into conservation, the

17  accuracy of the inferences generated from these records can be highly impacted by gaps and

18  biases in the data. Because of the Philippines' status as a biodiversity hotspot, the assessment of

19  potential gaps and biases in biodiversity documentation in the country can be a critical step in the

20  identification of priority research areas for conservation applications. In this study, we

21  systematically assessed biodiversity data on animal and plant taxa found in the Philippines by

22  examining the extent of metadata gaps, taxonomic biases, and spatial biases in DNA barcode

23  data while using species occurrence data as a backdrop of the Philippines' biodiversity. These

24  barcode and species occurrence data sets were obtained from public databases, namely:

25  GenBank, Barcode of Life Data System and Global Biodiversity Information Facility. We found

26  that much of the barcode data had missing information on either records and publishing,

27  geolocation, or taxonomic metadata, which consequently, can limit the usability of barcode data

28  for further analyses. We also observed that the amount of barcode data can be directly associated

29  with the amount of species occurrence data available for a particular taxonomic group and

30  location – highlighting the potential sampling biases in the barcode data. While the majority of

31  barcode data came from foreign institutions, there has been an increase in local efforts in recent

32  decades. However, much of the contribution to biodiversity documentation only come from

33  institutions based in Luzon.

34

35  Keywords: online biodiversity database, genetic diversity, species diversity, sampling biases,

36  spatial analysis, comparative analysis, conservation

## Introduction

Biodiversity is the product of the interactions between many physical and biological processes across time (Boero & Bonsdorff, 2007; van der Plas, 2019). Unfortunately, recent anthropogenic activities have significantly impacted biodiversity resulting in its rapid decline (Halpern et al., 2008, 2015). If left unabated, this alarming biodiversity loss can potentially impair the capacity of ecosystems to support and sustain life over time (Ayyad, 2003; Butchart et al., 2010; Cardinale et al., 2012; Reich et al., 2012; Worm et al., 2006). Due to these anthropogenic impacts on biodiversity, conservation efforts have been implemented to mitigate biodiversity loss and to promote the recovery of affected ecosystems and species. These initiatives include prioritization and management of key areas that best represent biodiversity or the processes (i.e., ecological and evolutionary) sustaining it (Beger et al., 2014; Herrick et al., 2006; Hoffmann & Sgró, 2011; Moritz, 2002; Richardson & Whittaker, 2010; Selig et al., 2014; Sgrò et al., 2011). However, efforts to conserve biodiversity could potentially be ineffective, or even counterproductive, if there is a lack of understanding of the fundamental processes underlying biodiversity (e.g., Hoveka et al., 2020; Santangeli et al., 2013). Thus, an understanding of biodiversity and the processes underpinning it is necessary in order to improve the efficacy of conservation efforts.

Because biodiversity is organized at different levels (i.e., ecosystems, species, and genes), making inferences about biodiversity-generating processes that are relevant to conservation will require documentation and analysis of biodiversity at various levels (Laikre et al., 2010; Purvis & Hector, 2000; Sarkar & Margules, 2002). Although significant progress has been made regarding biodiversity documentation, there has always been a tendency for biodiversity data to be spatially and taxonomically biased. This bias is often in contrast with the natural patterns and distribution of biodiversity (Titley et al., 2017; Troudet et al., 2017). For example, globally, biodiversity documentation is biased towards developed countries within temperate regions despite the tropical regions being relatively more diverse (Meyer et al., 2015; Newbold, 2010; Titley et al., 2017). At regional scales, spatial bias is also prominent primarily because many biodiversity documentations are results of scientific research focused on answering specific questions. Consequently, sampling is associated with certain geographical features related to the research question (e.g., near or within protected areas). This bias potentially leads to the under-representation of many key habitats in biodiversity documentation (Fisher-Phelps et al., 2017; Newbold, 2010). Current knowledge on biodiversity is further biased towards more charismatic organisms (i.e., mostly plants and vertebrates) leaving significantly more diverse taxonomic groups, such as invertebrates, understudied (Titley et al., 2017; Troudet et al., 2017). Overall, the extent of biases in biodiversity documentation reflects the insufficient data in many regions and taxa, which are likely due to limited research topics brought by various historical, social, economic, and practical factors (dos Santos et al., 2020; Troudet et al., 2017).

76       The various spatial and taxonomic biases in biodiversity data can potentially affect key
77  inferences about biodiversity-related processes (e.g., Keyse et al., 2014; Matias & Riginos,
78  2018). Because these inferences are explicitly being incorporated in conservation, these biases
79  can potentially lead to poorly-advised decisions that may contribute to biodiversity decline.
80  Moreover, conservation entails costs at various stages of its implementation (i.e., opportunity,
81  acquisition, management, and maintenance), and providing for this cost involves allocation of
82  highly-constrained resources such as time and money (Margules & Pressey, 2000; Possingham &
83  Wilson, 2005). Thus, mitigating the impact of these biases can benefit conservation efforts by
84  making them cost-effective in the use of those valuable resources, particularly for countries
85  where such resources are limited but where conservation is in demand.
86
87       One example of countries that will benefit greatly from cost-effective conservation
88  efforts is the Philippines. This country is a tropical developing country that has been considered
89  as one of 17 megadiverse nations worldwide (Mittermeler & Mittermeler, 1997), largely due to
90  its rich diversity and endemism. It has been estimated that there are more than 38,000 species of
91  vertebrates and invertebrates in the country (Catibog-Shinha & Heaney, 2006) – a likely
92  conservative number given the variability in estimates across groups. For example, as new
93  species are being discovered, some reports have predicted that Philippine arthropod species
94  would eventually reach 50,000 to 100,000 in number (Gapud, 2002). For plant taxa, around
95  14,000 species are found in the Philippines (Madulid, 1985 as cited in Lagunzad et al., 2002)
96  along with 35 of 54 mangrove species (Tomlinson, 1986 as cited in Primavera, 2002), more than
97  1,000 seaweed species (E. Fortes, 2002), and 16 seagrass species (M. Fortes, 1986 as cited in M.
98  Fortes, 2002). Among the animal and plant species that have been described so far, more than
99  half of them are said to be endemic to the Philippines (Ong, 2002).
100
101       Despite the number of species that have already been described in the Philippines, there
102  are still a lot of uncertainties regarding the estimated biodiversity in the country. Moreover, there
103  is also growing threats on the local environment as the Philippines became one of the "hottest"
104  biodiversity hotspots in the world due to the amount and rate of loss and degradation in various
105  habitats (Halpern et al., 2015; Harvey et al., 2020; Myers et al., 2000). These threats to
106  biodiversity have increased the need for conservation. Yet, the gaps in biodiversity
107  documentation in the country can potentially constrain these efforts. Addressing this problem
108  will require the identification of biases present in biodiversity records. Thus, a comprehensive
109  and systematic assessment of the current biodiversity data is needed to ensure the efficacy of
110  future conservation efforts based on such information.
111
112       Previous works that have examined biases and gaps in biodiversity data have utilized
113  publications collected from search engines such as the Web of Science (dos Santos et al., 2020;
114  Titley et al., 2017) or certain biodiversity records obtained from public databases. For example,
115  DNA barcode data from GenBank identified through published work has been used to examine

116 the extent of DNA barcoding in the Philippines (Fontanilla et al., 2014). Similarly, for many
117 works, species occurrence data from the Global Biodiversity Information Facility (GBIF) is used
118 (Fisher-Phelps et al., 2017; Meyer et al., 2015; Oliveira et al., 2016; Troudet et al., 2017).
119 Importantly, in these previous examinations, species occurrence and DNA barcode data are
120 typically examined separately for biases and gaps. However, given that the components of
121 biodiversity and its underlying processes are fundamentally intertwined (e.g., genetic data
122 shedding light on cryptic species diversity), it becomes critical that species and genetic data are
123 examined side by side. This approach can potentially help identify common patterns of biases
124 and gaps in the documentation of biodiversity at both levels.
125
126 In this study, public databases are leveraged to systematically examine potential gaps and
127 biases present in current records and gain a better understanding of the state of biodiversity
128 documentation in the country. The study specifically focuses on public biodiversity data of
129 animal and plant taxa found in the Philippines that are accessible in three online databases,
130 namely: the Global Biodiversity Information Facility (GBIF), GenBank, and Barcode of Life
131 Data System (BOLD). These databases represent large repositories of biodiversity records that
132 are widely used among the scientific community – as well as citizen scientists mainly in the case
133 of GBIF (Petersen et al., 2021) – to publish data. Because these datasets are readily accessible,
134 they represent records more frequently processed and analyzed to generate inferences for
135 policymaking and conservation planning (Ball-Damerow et al., 2019). Thus, examining
136 biodiversity data from these databases will not only identify biases in the current data but can
137 also mitigate the risks posed by these biases to conservation efforts.  Although both species and
138 genetic data will be utilized, the analyses in this study will mainly focus on the genetic data with
139 species data serving as a background. Because species data from public database have prominent
140 biases (some inherent with citizen science and its opportunistic nature of collection), its
141 comparison with genetic data can potentially highlight biases in genetic data as well (Amano et
142 al., 2016; Petersen et al., 2021; Troudet et al., 2017). To systematically assess both datasets,
143 species and genetic data are examined for the following: (1) metadata gaps in relation to the
144 completeness of biodiversity records; (2) taxonomic biases at the species and genetic levels; and
145 (3) spatial biases in terms of sampled locations and origin of leading contributors. These
146 assessments are done to identify potential knowledge gaps present in Philippine biodiversity.
147 This approach is a key step in addressing biases to generate more accurate inferences and
148 develop better strategies on how to move forward in future efforts in biodiversity documentation
149 and conservation.
150
151 **Materials & Methods**
152 *Collecting and parsing of biodiversity data*
153 In examining the Philippines biodiversity data, we limited our collection of data to three
154 databases that are widely used and are easily accessible. Thus, our study represents information
155 that is likely to be used by many researchers or even policymakers. We obtained species

156   occurrence data directly from the Global Biodiversity Information Facility (GBIF,
157   https://www.gbif.org/) on October 18, 2020 (GBIF.org, 2020a, 2020b). The search was filtered
158   by country ("Philippines"), occurrence states ("Present"), and taxonomic key ("Animalia" and
159   "Plantae"). The barcode data was obtained directly from two separate databases, namely:
160   GenBank (https://www.ncbi.nlm.nih.gov/genbank/) on November 1 and 3, 2020 and Barcode of
161   Life Data System (BOLD, http://v4.boldsystems.org/) on November 3, 2020. In GenBank, four
162   searches were conducted using different sets of keywords to obtain barcode data based on the
163   gene marker of interest. The gene markers actively searched for in GenBank were the following:
164   cytochrome oxidase c subunit I (using the keywords, "COI OR co1 OR cox1 OR coxI OR
165   cytochrome oxidase OR cytochrome c oxidase AND Philippines"); cytochrome b (using the
166   keywords, "cytb OR cyt-b OR cyt b OR cytochrome b OR cytochrome-b AND Philippines");
167   ribulose-1,5-biphosphate carboxylase (using the keywords, "ribulose-1,5-bisphosphate
168   carboxylase OR rbcl OR rubisco OR ribulose-bisphosphate carboxylase AND Philippines");
169   maturase K (using the keywords, "matk OR MaturaseK OR maturase K AND Philippines"); and
170   lastly, internal transcribed spacer 2 (using the keywords, ""internal transcribed spacer 2" OR
171   ITS2 OR ITS AND Philippines"). Prior to downloading data from GenBank, the results of each
172   search were filtered based on species to only include "Animals" and "Plants". It is important to
173   note that the data obtained may have included entries labelled as "unverified" since our searches
174   were unfiltered for verification. In BOLD, several searches were conducted in the Public Data
175   Portal system based on geography (keyword, "Philippines") and taxonomy (using all taxonomic
176   groups listed under animals and plants in BOLD's Taxonomy Browser –
177   http://v4.boldsystems.org/index.php/TaxBrowser_Home).
178
179         We mainly utilized the data.table R package (Dowle & Srinivasan, 2020) to manage and
180   parse through the data we obtained. However, in the case of GenBank data, the downloaded data
181   had to be processed into more readable files for each data entry. We used our own set of R
182   functions – specifically made to parse through individual GenBank files – to pull out as much
183   information as possible and organize it into a more workable data frame. We created seven
184   functions that obtained the following information: (1) taxonomy of the specimen; (2) publishing
185   author; (3) publishing institution; (4) year submitted; (5) metadata associated with the "source";
186   (6) gene marker; and (7) barcoding sequence (made available in github.com/dinmatias). We also
187   conducted additional cleaning and fixing on the information pulled out from the GenBank files
188   on BOLD cross-reference, taxonomy, publishing institution, gene marker, and sampling location.
189   For the taxonomy information, we created a database derived from the unique species found in
190   GBIF to obtain only the information on phylum/division, class, order, family, and genus while
191   other taxonomic ranks were disregarded. To obtain the publishing institution, we manually
192   parsed through the unique publishing entries and narrowed the information down to two columns
193   that contained the name of the main institution involved (labelled as PublishingInstitution) as
194   well as the country where it is based (labelled as PublishingCountry). In the BOLD data, we
195   added an additional column for the country where the storing institution, copyright institution,

196 and sequencing center are based. Some of the gene markers entries initially pulled out were
197 unclear or vague due to the varying ways the information was laid out in the individual GenBank
198 files and how the markers were named (e.g., full name or different abbreviations). For these
199 reasons, these entries were manually parsed to standardize the names of the gene markers used.
200 While the coordinate entries for the sampling information required minimal cleaning, the
201 descriptive information on the locality where the specimen was sampled required intensive
202 manual parsing. This editing was done not only for GenBank data but also for BOLD data to
203 obtain the specific information on province, municipality, and/or barangay based on a location
204 database we derived from the Philippine Standard Geographic Code (PSGC) (Philippine
205 Statistics Authority, 2020). During the parsing and cleaning process, sampling information was
206 categorized based on the kind of issues we encountered during the parsing (if any) that made
207 them vague or inconclusive (see Table S1). Moreover, the descriptive information provided for
208 the sampling locality in the GBIF data was parsed through and cleaned such that it was organized
209 into province, municipality, and/or barangay.
210
211         After parsing and cleaning of the data, we obtained the subsets of the main datasets
212 containing the metadata associated with the following categories: records (i.e., entry ID and
213 collection date), taxonomy (i.e., phylum/division, class, order, family, genus, species),
214 geolocation (i.e., coordinates and administrative units where the specimen was sampled), and
215 publication (i.e., submission date, publishing institution and country) (see Table S2). For
216 taxonomy, we recognize that there are differences between animal and plant taxonomy,
217 particularly with regards to the taxonomic ranks lower than kingdom – e.g., phylum for animal
218 taxa and division for plant taxa. However, phylum and rank were placed in the same taxonomic
219 metadata in the species and genetic databases we collected from – generally being categorized as
220 "phylum". Hence, in this study, phylum and division were treated as one classification in the
221 analyses. For the downstream analyses, the GenBank and BOLD data sets were combined into
222 one barcode dataset after selecting the metadata of interest. In combining these two data sets, we
223 ensured that the columns (variables) were analogous between the two databases. We further
224 filtered our two main working datasets (i.e., species occurrence and barcode data) by excluding
225 the following entries: duplicates in barcode data based on accession number; gene markers that
226 were not part of the five markers actively searched for; barcode specimen sampled from foreign
227 countries; and species occurrence and barcode data on *Homo sapiens* and *H. luzonensis*.
228 Additionally, a substantial number of barcode records with missing information on the country of
229 collection was observed despite having filtered the searches based on geography. Because this
230 number was substantial, two sets of analyses were conducted: (1) one where *NA* was excluded
231 and (2) another where *NA* was included in the dataset. While it is likely that the latter approach
232 may have included a few sequences that are not actually from the Philippines, the results were
233 generally the same between the two sets of analyses. Thus, the results from the latter approach
234 were mainly presented.
235

236   *Examining for metadata gaps*
237        To assess the completeness of the metadata associated with the barcode data, we
238   quantified the number of records with missing information on the following categories:
239   publication and records, sampling location, and taxonomy. In the publication and records
240   category, the number of records that lacked information on the copyright institution, collection
241   year, and submission year were counted. In the sampling location category, we counted the
242   number of records that lacked coordinates (i.e., latitude and longitude) and within this data
243   subset, the proportion of records with (or without) additional information on the sampling
244   locality was examined. Additionally, we determined the frequency of each kind of issue
245   encountered while manually parsing through the descriptive information on the sampling locality
246   – with those having more than one issue being categorized as "mixed".
247
248        In the taxonomy category, we first assessed the entries that had information on the
249   species level but lacked information on one or more higher taxonomic ranks. Here, the original
250   entries for the species information that included the keywords, "sp." and "gen." were marked as
251   *NA* since the true species identity was not provided. For records with identified species but
252   incomplete taxonomic data, we attempted to fill in the missing entries using the same database
253   we derived from the taxonomy of unique species in GBIF. Because barcode data is mainly used
254   as a reference in "species identification", the use of sequences that are not identified to species
255   level is not maximized. Hence, to identify and examine the taxonomic groups with barcode data
256   with low species identification, we plotted the percent of identified species in barcode data
257   against the percent of species with available barcode records that are represented in species
258   occurrence data. This was done separately for animal and plant records at the phylum/division,
259   class, order, and family levels.
260
261   *Examining for taxonomic biases*
262        To compare the extent of species and genetic documentation among taxonomic groups,
263   we plotted the number of available records per taxon in barcode data against that of species
264   occurrence data. The data was first transformed using logarithmic function prior to plotting.
265   Similar to the previous section on taxonomic metadata gaps, this was done separately for animal
266   and plant records at the phylum/division, class, order, and family levels. Additionally, quantiles –
267   specifically, the $5^{th}$ and $95^{th}$ percentile – of both datasets were incorporated in the plots to
268   highlight taxonomic groups on the extreme 10% of the distribution of these two variables. Here,
269   the GBIF occurrence record was used as a measure of the commonness of a taxonomic group in
270   examining how well commonly recorded taxonomic groups are being barcoded.
271
272   *Examining for spatial biases*
273        To assess the sampling distribution of barcode and species data, we first obtained
274   shapefiles of the Philippine administrative boundaries, specifically, the Philippines - Subnational
275   Administrative shapefile (https://data.humdata.org/dataset/philippines-administrative-levels-0-to-

276    3). Using this database, the province information of a given coordinate entry was determined
277    based on which defined boundaries of the administrative level 2 (i.e., province boundary) it falls
278    under. In the case of marine specimens with coordinates that do not fall within any province
279    boundary (because the boundary is based on land), the nearest province to them was assigned as
280    their province information. The nearest province was determined by first identifying the
281    "centroid" of each province and then measuring the distance of a data point to the centroid. The
282    province with the shortest distance from the data point was assigned as its province. For records
283    without any coordinates, only records with information on the province where the specimen was
284    sampled were included. These filtered data sets were then used to generate separate heatmaps for
285    the sampling distribution of barcode and species occurrence. Moreover, we also plotted the
286    number of records per province in barcode data against that of species occurrence data, with the
287    data transformed logarithmically prior to plotting and the 5th and 95th percentiles incorporated.
288
289         To examine the distribution of global contribution to Philippine barcode data, we focused
290    on the countries where the institutions that submitted or, in the case of BOLD, held the copyright
291    to the image data are from (i.e., copyright_institutions). Another metadata column in BOLD that
292    was considered to be examined for contribution was the institute that served as the storage place
293    of the voucher (i.e., institution_storing); however, the entries of the two columns were generally
294    the same. We quantified the number of barcode records published per country and visualized
295    their spatial distribution through the wrld_simpl shapefile from the maptools R package (Bivand
296    & Lewin-Koh, 2021). Additionally, the contribution of "local" and "foreign" efforts in
297    generation barcodes across time was compared. For this comparison, barcode records were
298    categorized as contributed by either "Foreign" or "Philippines" based on the copyright country.
299    This information was summarized into two plots showcasing the barcoding activity through time
300    in terms of year of collection (starting from the 1990s) and year of submission/publication
301    (starting from the 2000s). Note that we presented the barcoding activity across the year by
302    "smoothened" curved obtained through local regression (i.e., loess regression).
303
304         We then examined the contribution to barcode data at the national level – meaning
305    different institutes based in the Philippines. For each barcode record, we assigned the
306    "processing center" (i.e., region where the institute holding the copyright is located) and "region
307    sampled" (i.e., region where the specimen was collected). The total number of barcode records
308    generated by each "processing center" from a specific "region sampled" was used as its
309    contribution per "region sampled". The local contribution data was then summarized via a
310    correlation matrix heatmap, which plotted the region of sampling against the region of local
311    institutions. In this matrix, the regions were sorted according to their proximity to provide spatial
312    context. We utilized the following R packages to conduct our spatial analyses: sp (Bivand et al.,
313    2013; Pebesma & Bivand, 2005), raster (Hijmans, 2020), rgdal (Bivand et al., 2021), and
314    RColorBrewer (Neuwirth, 2014).
315

## Results

### Initial processing of biodiversity data

From the initial database searches conducted in late October to early November 2020, a total of 31,163 barcode records – 18,094 from GenBank and 13,069 from BOLD – and 1,557,709 species occurrence records were retrieved. Upon parsing through the raw datasets, duplicates, unwanted gene markers, and foreign samples in the barcode data as well as records involving *H. sapiens* and *H. luzonensis* in both barcode and species data were excluded. This initial filtering resulted in 20,482 barcode (16,719 excluding *NA* entries for country sampled) and 1,557,374 species records available for downstream analyses. For the barcode data, the majority of the records obtained are based on the COI gene marker (see Fig. 1A). This may be linked to the significantly higher number of animal records analyzed in comparison to the number of plant records (a trend also observed in the available species occurrence data, see Table S3) since gene markers are often utilized for certain organisms (e.g., COI for animals then rbcL and matK for plants).

### Metadata gaps in Philippine barcode data

Most of the barcode data used in the analyses were observed to have incomplete information in one or more categories of metadata. For the gaps in the records and publishing metadata, among the barcode data, 72.52% lacked information on the year of collection (66.73% excluding *NA* entries for country sampled), 22.01% on the year of submission (26.93% excluding *NA* entries for country sampled), and 18.51% on the publishing or copyright institution (22.64% excluding *NA* entries for country sampled). For the gaps in the geolocation metadata, approximately 65.78% had no coordinates (58.10% excluding *NA* entries for country sampled) and within that subset of data, more than half lacked any additional descriptive information on the sampling locality such as province, municipality, and barangay. Overall, 46.68% of barcode records lacked any kind of metadata on the sampling location (34.69% excluding *NA* entries for country sampled).  Records that did have metadata on the sampling locality in terms of administrative units were relatively difficult to parse through. Majority of them were vague in varying degrees depending on the kind of major issue encountered – with most being unspecified (see Fig. 1B). Additionally, there were several records wherein "Philippines" was indicated as the country sampled but upon further inspection of the description of the specific locality sampled, a mismatch was found.  Such entries were labelled as foreign and excluded from the analyses.

For the gaps in the taxonomic information, 3,793 records had no information on the specific group in one or more taxonomic ranks despite the specimen being identified at the species level. Using a taxonomic database derived from the species occurrence data, these gaps were filled in at the phylum/division, class, order, and family levels, narrowing down the number to 706 records with incomplete taxonomic information. The proportion of identified animal and plant species was also assessed in relation to the proportion of barcoded species per taxon at a

356    specific taxonomic rank – namely, phylum/division, class, order, and family (see Fig. 2). At the
357    phylum/division level, most of the taxa exhibited more than 50% percent species identification
358    except for Annelida and Rotifera (see Fig. 2A). However, at lower taxonomic ranks, there were
359    more taxa that had the majority (more than 50%) of their records unidentified at the species level
360    (see Fig. 2B to 2D). Moreover, while more taxa were being sampled, the rate at which these
361    groups were barcoded remains relatively low. Evidently, only a few groups exhibited a high
362    percentage of identified and barcoded species. It is important to note, however, that the identity
363    of the species was based on the information provided by the contributors who published the
364    barcode records. It was not verified if the species identities matched with the barcode sequences
365    associated with them. Additionally, in evaluating the proportion of barcoded species at the order
366    and family level, several taxa returned an undefined value (*NaN*). These were likely the result of
367    the absence of species occurrence records associated with those taxa despite having barcode
368    records available. There were eight (8) orders resulting in *NaN*, labelled as the following:
369    "Labriformes", "Ovalentaria", "Gobiiformes", "Trachiniformes", "Pristiformes", "Pulmonata",
370    "Vetigastropoda", and "Sebdeniales". On the other hand, there were five (5) resulting *NaN*
371    families, labelled as: "Pentanchidae", "Chilodontidae_gas", "Choristellidae", "Sebdeniaceae",
372    and "Areschougiaceae".
373
374    *Taxonomic biases in Philippine barcode data*
375        Examination of the taxonomic distribution of records collected revealed a general
376    increasing trend between the amount of barcode and species occurrence data for a particular
377    taxon (see Fig. 3). At the phylum/division level, the group with the highest record in both
378    barcode and species data was Chordata and accompanying it in the areas of either high genetic
379    data or high species data were Arthropoda, Mollusca, and Tracheophyta (see Fig. 3A). On the
380    other hand, the groups that had particularly low biodiversity records, particularly in terms of
381    barcode data, were Rotifera, Ctenophora, and Marchantiophtya. There were several taxa that had
382    species occurrence data but lacked barcode data, namely: Anthocerotophyta, Brachiopoda,
383    Bryozoa, Cephalorhyncha, Chaetognatha, Charophyta, Entoprocta, Hemichordata,
384    Nematomorpha, Phoronida, Sipuncula, and Xenacoelomorpha. Assessing the trends further down
385    the taxonomic hierarchy, it could be observed that while more groups had been sampled in terms
386    of species occurrence, many of them had little to no barcode records available (see Fig. 3B to
387    3D). Furthermore, groups that remained at or above the 95[th] percentile of genetic and species
388    data at the class, order, and family levels mostly belonged to Phylum Chordata.
389
390    *Spatial biases in Philippine barcode data*
391        Examination of the spatial distribution of records obtained showed a high similarity
392    between the sampling distributions of barcode and species occurrence data, particularly in terms
393    of the provinces wherein sampling was most and least concentrated (see Fig. 4A and 4B). In both
394    genetic and species data, the province that had been relatively more sampled (above the 95[th]
395    percentile) was Palawan. These similarities in sampling distribution meant that the amount of

396    barcode data could be directly related to the amount of species occurrence records sampled per
397    province (see Fig. 4C) – similar to the previous section on taxonomic bias. Furthermore, several
398    provinces were observed to fall under the 95th percentile of either dataset. For barcode data, in
399    particular, the provinces with the highest records (above 95th percentile) were Siquijor, Cavite,
400    Bohol, Aurora, and Palawan while the ones with the records (below 5th percentile) were Tarlac,
401    Basilan, Maguindanao, Zamboanga Sibugay, and Northern Samar.
402

403            Examination of the institutions contributing to the barcode data revealed that in provinces
404    where barcode sampling was most concentrated, the majority of the records were generated by
405    foreign institutions. A notable exemption was Pangasinan, the seventh most sampled location in
406    terms of barcoding data, majority of which were contributed by local institutes (~70.42% of the
407    records). A similar trend of high contribution by foreign institutions to barcoding was observed
408    when all barcode data were examined. While the Philippines had the most contribution to its
409    barcode records compared with other countries (see Fig. 5A), a comparison of the foreign and
410    local contributions showed that the Philippines had contributed only about 30.00% of the overall
411    barcode data on Philippine animal and plant biodiversity.
412

413            When foreign and local contribution of barcode data were examined across time –
414    specifically, the time of collection and submission – it was revealed that the Philippines had
415    increasingly collected and submitted more records by 2005. Moreover, at some point, the
416    Philippines had even surpassed the activity of foreign institutions (see Fig. 5B and 5C).
417    Additionally, though not represented in Fig. 5B, many of the specimens used by foreign
418    institutes in generating barcode data had been collected before the 1990s, even as far back as
419    1915, highlighting the importance of sample preservation in documenting not only species but
420    potentially genetic diversity as well.
421

422            Within the Philippines, there was also a substantial discrepancy in contributions of local
423    institutions to barcode data (see Fig. 6). When the regions of barcode-generating institutions
424    (termed as the "Processing Center") were compared with regions where sampling was
425    conducted, it was apparent that only six of seventeen regions were able to generate barcode data
426    for their local biodiversity (diagonals in Fig. 6). Furthermore, most local contributions were from
427    institutions found in the regions of Metro Manila and Central Luzon, and these institutes sampled
428    the most either within their local area or in nearby regions, which were situated mainly in Luzon.
429    It is important to note, however, that this analysis was based on the local institutions that hold the
430    copyright to the records, and collaborations with other local institutions were not considered.
431

## Discussion

433            In this study, biodiversity records on animal and plant taxa found in the Philippines were
434    systematically assessed by examining the extent of metadata gaps, taxonomic biases, and spatial
435    biases in barcode data while using species occurrence data mainly as a baseline. Results show

436     that much of the barcode data had missing information on records and publishing, geolocation, or
437     taxonomic information. Moreover, it was observed that the amount of barcode data can be
438     directly associated with the amount of species occurrence data available for a particular
439     taxonomic group and sampling locality. Lastly, the results also reveal that majority of the
440     barcode data came from foreign institutions and while local barcoding efforts have increased in
441     the recent decades, much of it is due to Philippine institutions being based within Luzon.
442
443     *Incompleteness of metadata in barcode data*
444         Biodiversity records have been used in various fields of study to further understand the
445     underlying processes that influence biodiversity. Barcode data, in particular, have broad
446     applications in various fields – e.g., in understanding the processes affecting regions with high
447     diversity (Crandall et al., 2019; Matias & Riginos, 2018), in assessing the quality and
448     authenticity of food products sold in markets (Barbuto et al., 2010; Maralit et al., 2013; Pazartzi
449     et al., 2019), in conservation (Deichmann et al., 2017), and in battling illegal wildlife trade
450     (Hartvig et al., 2015). Despite the various uses of barcode data, its overall utility can be reflected
451     by the completeness of its metadata. Publishing and records information, for instance, would be
452     useful in finding relevant references for future research and examining the global, national, or
453     local state of biodiversity documentation. For example, in a similar study that focused on animal
454     barcoding in the Philippines, they found that only about 20% of records on native species could
455     be traced back to local institutions (Fontanilla et al., 2014). With this kind of information, it
456     would be easier to objectively assess the progress of a particular institution or country in
457     contributing to DNA barcoding or, more generally, to biodiversity documentation. Additionally,
458     while metadata may not directly contribute new knowledge on biodiversity and its processes, it
459     can provide context on the records being generated – particularly in terms of who, when, and
460     possibly why they were published for a particular taxon and/or locality. As previously discussed,
461     many of the available barcode records have missing metadata. It might be possible to manually
462     retrieve this information from journal publications linked to these records but when dealing with
463     a large amount of data, this approach could become challenging.
464
465         Another example of highly useful metadata is geolocation. By providing this metadata,
466     barcode records could then be used for studies that examine the role of geography in biodiversity
467     – such as the case of biogeographic research. For example, existing barcode records made it
468     possible to examine the processes behind the rich marine diversity in the Indo-Pacific region,
469     particularly at the molecular level (Crandall et al., 2019; Matias & Riginos, 2018). These
470     inferences would not have been possible without the information on the location where the
471     specimens were collected. It is important to note that there is, however, a concern for accuracy
472     when dealing with this kind of information. In this study, two kinds of geolocation information
473     were encountered: the numerical coordinates and the descriptive information on the locality.
474     Evidently, coordinates are relatively more accurate compared to descriptive information since
475     they could be easily standardized and used in spatial analysis. However, most barcode records

476    that were examined lacked coordinates. Contributors could have intentionally refrained from
477    including such information in their records or restricted access to it in the database since
478    coordinates – and geolocation in general – are considered to be "sensitive" data. Sensitive data is
479    any kind of information that, if made public, would cause an 'adverse effect' (e.g., illegal or
480    excessive collection, risk of disturbance) on the associated taxon or living individual (Chapman,
481    2020; Environmental Resources Information Network, 2016). Several governments – such as in
482    Australia (Andrews, 2009; Environmental Resources Information Network, 2016) and Canada
483    (AMEC Earth & Environmental, 2010) – have implemented legal policies that deal with
484    sensitive information of vulnerable species (e.g., plants and sessile animals, threatened or rare
485    species). These policies would then largely influence the guidelines of public databases – such as
486    GBIF (Chapman, 2020) – on managing the accessibility of sensitive metadata. With many
487    records lacking coordinates, the provinces pulled out from the descriptive information were
488    utilized for the analysis. Descriptions of the locality could also be informative. However, this
489    highly depends on how detailed and standardized they are which in turn, may depend on how
490    familiar the contributors were with the names and administrative units associated with the areas
491    being sampled. This may explain why the majority of entries with descriptive information (with
492    or without coordinates) were relatively more difficult to parse through (see Fig. 1B), with some
493    being unclear or inconclusive, while others were more informative.
494
495        While barcoding is a growing technique that has much potential in biodiversity studies,
496    one of its more popular applications is in species identification (Hebert & Gregory, 2005). Thus,
497    metadata on taxonomic information would prove essential for the DNA barcodes to be used as an
498    effective database, particularly for applications where organisms are not sampled (i.e.,
499    environmental DNA). While the results show that many taxonomic groups (see Fig. 2) had
500    incomplete taxonomic information or low species identification, they also identified potential
501    taxa for further taxonomic studies. Additionally, as the knowledge on taxonomy and
502    evolutionary relationships between different taxa grows, there is always a possibility for the
503    classification of a particular taxon to change. For examples, minor and major revisions have
504    recently been made in angiosperm (i.e., at the order and family levels) and annelid classification
505    (i.e., whole evolutionary tree) (Chase et al., 2016; Zrzavý et al., 2009). This changes in the
506    taxonomic classification may explain the anomalies observed in evaluating the percent of
507    barcoded species, as represented by the *NaN* orders and families. Upon further inspection, these
508    taxa mainly contained marine species, most of which were given the status of "Accepted" in the
509    World Register of Marine Species (https://www.marinespecies.org/). Moreover, the barcode
510    records associated with these *NaN* taxa were obtained specifically from BOLD. The current
511    taxonomic metadata of these records may also need to be updated. However, it is unclear
512    whether this responsibility falls with the contributors or the curators of the biodiversity data.
513
514        Overall, there were significant metadata gaps present in the current barcode records on
515    Philippine biodiversity that were retrieved from GenBank and BOLD – particularly, the

516  information on the sampling location and identity of the species. Regardless of whether these
517  kinds of information are being collected by researchers, if they are not included in the
518  submission to these public databases, they can be perceived as missing. In this study, due to the
519  extent of missing information, not all barcode records were deemed useful in some of the
520  analyses. This does not necessarily imply that barcode records with incomplete metadata are
521  unusable but highlights how the completeness of metadata allows these records to be used in
522  various kinds of analyses. Because of the importance of metadata, its collection and publication
523  have been strongly advocated and have inspired the creation of a database for metadata (Deck et
524  al., 2017). Thus, researchers and contributors need to acknowledge the importance of metadata
525  and be aware that in order to increase the utility of current biodiversity records, there is a need to
526  also increase the availability of metadata by collecting and properly sharing this information with
527  public databases. With regards to sensitive data (e.g., coordinates of vulnerable species), it may
528  be possible to acquire authorization from the contributors to access the metadata (Chapman,
529  2020). Otherwise, the sampling locality description may be a sufficient substitute for
530  coordinates, provided that the entries are more standardized and informative up to the province
531  level, at least.
532
533  *Barcode data favoring commonly documented taxa*
534      In examining for taxonomic biases, it was observed that the rate of barcoding of taxa was
535  associated with how commonly they were observed (see Fig. 3). Given that species occurrence
536  records are largely opportunistic in nature (Petersen et al., 2021), the strong association between
537  species and genetic datasets may indicate certain biases that are inherent to barcode data as well.
538  Other than commonness, other factors might contribute to the variability in barcoding effort
539  across taxonomic groups in the Philippines. For example, popular research likely influenced
540  interest in barcoding of specific taxonomic groups. These topics include high endemism of
541  vertebrates and vascular plants, and the high marine biodiversity in the Philippines, which led to
542  efforts of barcoding vertebrates, endemic plants, and reef fishes, respectively (Carpenter &
543  Springer, 2005; Ong, 2002; Posa et al., 2008). The limited number of experts available in the
544  Philippines could potentially contribute to the observed taxonomic bias (Arayata, 2019; Senate
545  of the Philippines, 2017). This lack of expertise is evident, for example, in the online roster of
546  experts provided by the Department of Environment and Natural Resources – Biodiversity
547  Management Bureau (https://bmb.gov.ph/index.php/resources/roster-of-experts), where it is
548  evident that not all plant and animal taxa are well-represented. Furthermore, in relation to the
549  findings on spatial bias, most DNA barcoding is processed in institutions based in Metro Manila.
550  Each of these universities has a limited number of researchers with research interest focused only
551  on certain taxa. Although there may be local experts specialized in less-represented groups, these
552  experts may be based in regions where there is limited access to molecular approaches. In this
553  case, collaborations become essential in providing these experts access to molecular facilities.
554  Due to these factors, more attention in Philippine barcoding may have been given to certain

555    groups belonging to the following phyla/divisions: Chordata, Arthropoda, Mollusca, and
556    Tracheophyta.
557
558        Some exceptions were observed from the general trend that high genetic data can be
559    expected with high species data. For example, there is currently no barcode data for the Family
560    Ceratobatrachidae (Phylum Chordata, Class Amphibia) despite the more than 20 species of
561    limestone-forest frogs (*Platymantis*) recorded in the Philippines (Siler et al., 2009). Given the
562    high endemicity and potential cryptic species diversity among this group (Siler et al., 2009),
563    DNA barcode data can prove valuable in documenting the diversity within this taxon. Among
564    plants, the Family Dipterocarpaceae (Division Tracheophyta, Class Magnoliopsida) is an
565    example of taxon that lacks barcode data. This family contains ecologically important yet
566    exploited and endangered tree species. Examples are species of the genus *Parashorea*, *Shorea*,
567    and *Hopea*, which largely contribute to the tree diversity and richness in many Philippine forests
568    such as Mt. Apo Natural Park and Rajah Sikatuna Protected Landscape (Aureo et al., 2020;
569    Zapanta et al., 2019). Unfortunately, some species in this family have become vulnerable to
570    exploitation brought by logging, leading to some being critically endangered (Aureo et al., 2020;
571    Zapanta et al., 2019). The lack of barcode data for these animal and plant taxa translates to
572    missed opportunities in obtaining valuable information for this group – information that could be
573    used in understanding the diversity of these groups and in the conservation of vulnerable species.
574

575    *Barcode data favoring areas with high species documentation & foreign contributors*
576        In examining for spatial biases, a similar trend was observed with the taxonomic biases.
577    Specifically, examination of location information showed that barcode sampling is more likely
578    conducted in areas where documentation of species is commonly done (see Fig. 4). The results
579    revealed that the five provinces with the highest barcode sampling were Siquijor, Cavite, Bohol,
580    Aurora, and Palawan. Three of these provinces are found in Luzon making them accessible to
581    institutes that had the capacity to barcode. This accessibility however does not only pertain to
582    proximity to barcoding institutions, but also to protected areas as well as the availability of routes
583    to sampling locations (Fisher-Phelps et al., 2017; Oliveira et al., 2016). Indeed, in the
584    Philippines, local biodiversity more frequently sampled are situated in provinces with more
585    developed travel routes (or relatively near to urban areas). Security and safety are also linked to
586    accessibility of an area. Governments often provide travel advisories that restrict access to
587    certain areas due to the high risk of threats such as disease outbreaks, natural disasters, civil
588    unrest, or terrorist attacks (Foreign, Commonwealth & Development Office, 2013). For instance,
589    foreign researchers, who have been observed to generate a large portion of Philippine barcode
590    data, are often strongly advised against travelling to many provinces in Mindanao due to "crime,
591    terrorism, civil unrest and kidnapping" (Government of Canada, 2021; U.S. Department of State,
592    2021). As a result, provinces that are deemed to have lower risks to local and foreign researchers
593    are more likely to be sampled compared to other provinces.
594

595    Another aspect of spatial bias examined in this study was the origins of contributors. It
596    must be noted, however, that in this study, institutions holding the image copyright (specifically
597    for BOLD entries with images associated with them) were assumed to be the submitter of the
598    barcode data. In contrast to BOLD, submitter information is more explicitly indicated in
599    GenBank entries. From a global perspective, most of the current barcode data of Philippine
600    biodiversity was generated by foreign institutions with researchers from the United States being
601    the most active contributors (see Fig. 5A). The high contribution of foreign institutions is likely
602    due to the high research capacity of foreign institutions, especially in terms of funding and
603    barcoding facilities. For example, there exists a grant known as the "PIRE: Centennial Genetic
604    and Species Transformations in the Epicenter of Marine Biodiversity" that enables researchers
605    from various institutes based in the United States to conduct marine expeditions in the
606    Philippines (Carpenter et al., 2017). Moreover, foreign institutions may also have access to more
607    extensive specimen collections. For example, the Smithsonian National Museum of Natural
608    History houses over 126 million specimens in their catalog. Additionally, the United States has
609    about 1,500 other institutions that may also house a significant number of cataloged specimens
610    but often with restricted access (Page et al., 2015). It is likely that many of their specimens, not
611    exclusive to the United States, had been sampled even during the early years of exploration,
612    which may explain why there are several barcode records generated from older samples.
613
614    Examination of contribution to barcode data across time showed that Philippines has
615    become more active in barcoding in recent decades, particularly in terms of collecting samples
616    and submitting barcode data (see Fig. 5B and 5C). The upward trend in both collection of
617    samples and submission of barcode data seemed to have started between 2005 and 2010, around
618    the time DNA barcoding was slowly being adopted in the Philippines. For example, the UP
619    Institute of Biology initiated the creation of a public DNA barcode database in 2008 and several
620    years later, partnered with the Department of Environment and Natural Resources, to use DNA
621    barcoding against illegal wildlife trade (Encarnacion, 2019).
622
623    While local contributions to barcode data have increased over the years, spatial bias was
624    still prominent when the origins of contributors were examined from a national perspective.
625    Specifically, there was a mismatch between the localities producing (or processing) the barcode
626    data and the areas that were being sampled. This mismatch is likely due to the limited number of
627    local institutions with the capacity to process and generate barcode data, whether in terms of
628    facilities, funding, equipment, or expertise. Most of the local contributions are processed by a
629    small group of institutions located in the regions of Metro Manila and Central Luzon (see Fig. 6)
630    – many of which, if not all, have their own well-equipped DNA barcode laboratories. In line with
631    this, it may be possible to increase the capacity of local institutions found in regions where there
632    is currently minimal to no processing of barcode data by establishing the appropriate facilities
633    and conducting professional training. While this will require funding and time, it could empower
634    more local institutions to take initiative in barcoding their own local biodiversity – particularly

635  those based in regions that remain relatively unexplored. This would be ideal as these local
636  institutions are in the best position to sample their local biodiversity. Alternatively,
637  collaborations with other local institutions (e.g., local government agencies, non-governmental
638  organizations, etc.) can facilitate barcoding of local biodiversity. Indeed, many of the current
639  barcode records are a product of collaborations between institutions based in Metro Manila and
640  various local groups across the Philippines. While these may be indicated in the publications
641  linked to these records, there is no clear metadata information on collaborative works provided
642  on the raw barcode data obtained. The present limitation in the contributor metadata of these
643  public databases potentially under-represents the role of local institutions in the documentation
644  of Philippine biodiversity. For barcoding in particular, it is essential to acknowledge that both
645  sampling and barcode generating efforts are equally important. Hence, institutions who
646  contributed to either or both efforts in collaborations must also be credited equally – whether in
647  publications or databases. Thus, a more explicit acknowledgment of the roles of local
648  collaborators in the metadata associated with barcode data would increase the visibility of these
649  local institutes, which could potentially foster further collaborations in biodiversity
650  documentation.
651  
## Conclusions

653  By conducting a systematic assessment of the barcode data on animal and plant taxa, the
654  state of barcoding in the Philippines was examined, giving insight on the extent of metadata
655  gaps, taxonomic biases, and spatial biases present in current records. In analyzing the data, many
656  barcode records were found to have missing information for publishing, records, geolocation, or
657  taxonomic metadata. These gaps resulted in the exclusion of those records in some of the
658  analyses, demonstrating that incompleteness of metadata can limit the usability of barcode data
659  for different kinds of analyses. Also, the presence of metadata gaps makes biodiversity data more
660  tedious to work with. Philippine barcoding is more often conducted on taxa and provinces that
661  are associated with high documentation of species occurrence, with most records generated by
662  foreign countries with generally high research capacity. Moving forward with the findings of this
663  study, future contributors of barcode data are encouraged to increase the availability of metadata
664  by collecting and sharing this information to online databases upon submission to maximize the
665  potential utility of these records in various kinds of analyses. Additionally, future barcoding
666  efforts should prioritize areas where biodiversity documentation is currently lacking such as
667  documenting taxa and sampling regions that are under-represented in Philippine biodiversity
668  data. This approach of sampling under-represented taxa and regions may be done by
669  collaborating with institutions active in DNA barcoding and biodiversity experts specializing in
670  less-represented animal or plant taxa and by conducting field sampling in locations that currently
671  have limited data. Furthermore, it is essential to highlight the importance of empowering more
672  local institutions to take part in Philippine barcoding whether by increasing their capacity to
673  generate barcode data or collaborating with groups from different regions in the Philippines. For

674  future studies on the biases and gaps in biodiversity data, collaborations with data scientists are
675  also recommended to mitigate the tedious work involved in processing large amounts of data.
676

687

## Data Accessibility

689       The R scripts used in this work are deposited in the following:
690  https://github.com/miaberba/2021_PH_BiodiversityAssessment and
691  https://github.com/dinmatias/GeneBankParse. All data needed to run the analyses (i.e., files
692  needed for the various R scripts) are accessible through https://doi.org/10.5281/zenodo.6153441.
693

## References

695  Amano, T., Lamming, J. D. L., & Sutherland, W. J. (2016). Spatial Gaps in Global Biodiversity
696       Information and the Role of Citizen Science. *BioScience*, *66*(5), 393–400.
697       https://doi.org/10.1093/biosci/biw022
698  AMEC Earth & Environmental. (2010). *Best Practices For Sharing Sensitive Environmental*
699       *Geospatial Data*. Natural Resources Canada.
700       https://publications.gc.ca/site/eng/9.694895/publication.html
701  Andrews, J. (2009). *Sensitive Species Data Policy*. Department of Environment, Climate Change
702       and Water NSW. https://www.environment.nsw.gov.au/-/media/OEH/Corporate-
703       Site/Documents/Animals-and-plants/Wildlife-management/sensitive-species-data-
704       policy.pdf
705  Arayata, M. C. (2019, July 12). *PH needs more scientists: NAST*. Philippine News Agency.
706       https://www.pna.gov.ph/articles/1074747
707  Aureo, W. A., Reyes, T. D., Francis, F. C., Jose, R. P., & Sarnowski, M. B. (2020). Diversity and
708       composition of plant species in the forest over limestone of Rajah Sikatuna Protected
709       Landscape, Bohol, Philippines. *Biodiversity Data Journal*, *8*, 1–23.
710       https://doi.org/10.3897/BDJ.8.E55790
711  Ayyad, M. A. (2003). Case studies in the conservation of biodiversity: degradation and threats.
712       *Journal of Arid Environments*, *54*(1), 165–182. https://doi.org/10.1006/jare.2001.0881
713  Ball-Damerow, J. E., Brenskelle, L., Barve, N., Soltis, P. S., Sierwald, P., Bieler, R., LaFrance,
714       R., Ariño, A. H., & Guralnick, R. P. (2019). Research applications of primary biodiversity
715       databases in the digital age. *PLoS ONE*, *14*(9), 1–26.
716       https://doi.org/10.1371/journal.pone.0215794

717 Barbuto, M., Galimberti, A., Ferri, E., Labra, M., Malandra, R., Galli, P., & Casiraghi, M.
718     (2010). DNA barcoding reveals fraudulent substitutions in shark seafood products: The
719     Italian case of "palombo" (Mustelus spp.). *Food Research International*, *43*(1), 376–381.
720     https://doi.org/10.1016/j.foodres.2009.10.009
721 Beger, M., Selkoe, K. A., Treml, E. A., Barber, P. H., von der Heyden, S., Crandall, E. D.,
722     Toonen, R. J., & Riginos, C. R. (2014). Evolving coral reef conservation with genetic
723     information. *Bulletin of Marine Science*, *90*(1), 159–185.
724     https://doi.org/https://doi.org/10.5343/bms.2012.1106
725 Bivand, R., Keitt, T., & Rowlingson, B. (2021). *rgdal: Bindings for the "Geospatial" Data
726     Abstraction Library* (R package version 1.5-23). https://cran.r-project.org/package=rgdal
727 Bivand, R., & Lewin-Koh, N. (2021). *maptools: Tools for Handling Spatial Objects* (R package
728     version 1.1-1). https://cran.r-project.org/package=maptools
729 Bivand, R. S., Pebesma, E., & Gomez-Rubio, V. (2013). *Applied spatial data analysis with R*
730     (2nd ed.). Springer, NY. https://asdar-book.org/
731 Boero, F., & Bonsdorff, E. (2007). A conceptual framework for marine biodiversity and
732     ecosystem functioning. *Marine Ecology*, *28*(SUPPL. 1), 134–145.
733     https://doi.org/10.1111/j.1439-0485.2007.00171.x
734 Butchart, S. H. M., Walpole, M., Collen, B., van Strien, A., Scharlemann, J. P. W., Almond, R.
735     E. A., Baillie, J. E. M., Bomhard, B., Brown, C., Bruno, J., Carpenter, K. E., Carr, G. M.,
736     Chanson, J., Chenery, A. M., Csirke, J., Davidson, N. C., Dentener, F., Foster, M., Galli, A.,
737     … Watson, R. (2010). Global Biodiversity: Indicators of Recent Declines. *Science*,
738     *328*(5982), 1164–1168. https://doi.org/10.1126/science.1187512
739 Cardinale, B. J., Duffy, J. E., Gonzalez, A., Hooper, D. U., Perrings, C., Venail, P., Narwani, A.,
740     MacE, G. M., Tilman, D., Wardle, D. A., Kinzig, A. P., Daily, G. C., Loreau, M., Grace, J.
741     B., Larigauderie, A., Srivastava, D. S., & Naeem, S. (2012). Biodiversity loss and its impact
742     on humanity. *Nature*, *486*(7401), 59–67. https://doi.org/10.1038/nature11148
743 Carpenter, K., Barshis, D., Bird, C., Pinsky, M., & Polidoro, B. (2017). *NSF Award Search:
744     Award # 1743711 - PIRE: Centennial Genetic and Species Transformations in the
745     Epicenter of Marine Biodiversity*.
746     https://www.nsf.gov/awardsearch/showAward?AWD_ID=1743711
747 Carpenter, K. E., & Springer, V. G. (2005). The center of the center of marine shore fish
748     biodiversity: The Philippine Islands. *Environmental Biology of Fishes*, *72*(4), 467–480.
749     https://doi.org/10.1007/s10641-004-3154-4
750 Catibog-Shinha, C., & Heaney, L. R. (2006). *Philippine Biodiversity: Principles and Practice*.
751     Haribon Foundation.
752 Chapman, A. D. (2020). Current Best Practices for Generalizing Sensitive Species Occurrence
753     Data. *Copenhagen: GBIF Secretariat*. https://doi.org/https://doi.org/10.15468/doc-5jp4-
754     5g10
755 Chase, M. W., Christenhusz, M. J. M., Fay, M. F., Byng, J. W., Judd, W. S., Soltis, D. E.,
756     Mabberley, D. J., Sennikov, A. N., Soltis, P. S., Stevens, P. F., Briggs, B., Brockington, S.,
757     Chautems, A., Clark, J. C., Conran, J., Haston, E., Möller, M., Moore, M., Olmstead, R., …
758     Weber, A. (2016). An update of the Angiosperm Phylogeny Group classification for the
759     orders and families of flowering plants: APG IV. *Botanical Journal of the Linnean Society*,
760     *181*(1), 1–20. https://doi.org/10.1111/boj.12385
761 Crandall, E. D., Riginos, C., Bird, C. E., Liggins, L., Treml, E., Beger, M., Barber, P. H.,
762     Connolly, S. R., Cowman, P. F., DiBattista, J. D., Eble, J. A., Magnuson, S. F., Horne, J. B.,

763     Kochzius, M., Lessios, H. A., Liu, S. Y. V., Ludt, W. B., Madduppa, H., Pandolfi, J. M., …
764     Gaither, M. R. (2019). The molecular biogeography of the Indo-Pacific: Testing hypotheses
765     with multispecies genetic patterns. *Global Ecology and Biogeography*, *28*(7), 943–960.
766     https://doi.org/10.1111/geb.12905
767   Deck, J., Gaither, M. R., Ewing, R., Bird, C. E., Davies, N., Meyer, C., Riginos, C., Toonen, R.
768     J., & Crandall, E. D. (2017). The Genomic Observatories Metadatabase (GeOMe): A new
769     repository for field and sampling event metadata associated with genetic samples. *PLoS*
770     *Biology*, *15*(8), 1–7. https://doi.org/10.1371/journal.pbio.2002925
771   Deichmann, J. L., Mulcahy, D. G., Vanthomme, H., Tobi, E., Wynn, A. H., Zimkus, B. M., &
772     McDiarmid, R. W. (2017). How many species and under what names? Using DNA
773     barcoding and GenBank data for west Central African amphibian conservation. In *PLoS*
774     *ONE* (Vol. 12, Issue 11). https://doi.org/10.1371/journal.pone.0187283
775   dos Santos, J. W., Correia, R. A., Malhado, A. C. M., Campos-Silva, J. V., Teles, D., Jepson, P.,
776     & Ladle, R. J. (2020). Drivers of taxonomic bias in conservation research: a global analysis
777     of terrestrial mammals. *Animal Conservation*, *2*, 1–10. https://doi.org/10.1111/acv.12586
778   Dowle, M., & Srinivasan, A. (2020). *data.table: Extension of `data.frame`* (R package version
779     1.13.6). https://cran.r-project.org/package=data.table
780   Encarnacion, A. D. P. (2019). *Using DNA barcodes against the wildlife black market –*
781     *University of the Philippines System Website*. University of the Philippines: Research and
782     Breakthroughs. https://up.edu.ph/using-dna-barcodes-against-the-wildlife-black-
783     market/?fbclid=IwAR0BZo-
784     IK8wRXO4w2CGsVMkqNQ8saYajyV9XiLOtH7Vmja0sxbhjp-rUsrQ
785   Environmental Resources Information Network. (2016). *Sensitive Ecological Data Access and*
786     *Management Policy*. Australian Government Department of the Environment.
787   Fisher-Phelps, M., Cao, G., Wilson, R. M., & Kingston, T. (2017). Protecting bias: Across time
788     and ecology, open-source bat locality data are heavily biased by distance to protected area.
789     *Ecological Informatics*, *40*, 22–34. https://doi.org/10.1016/j.ecoinf.2017.05.003
790   Fontanilla, I. K. C., Torres, A. F., Cañasa, J. A. D., Yap, S. L., & Ong, P. S. (2014). State of
791     animal DNA barcoding in the Philippines: A review of COI sequencing of Philippine native
792     fauna. *Philippine Science Letters*, *7*(1), 104–137.
793   Foreign, Commonwealth & Development Office. (2013, August 22). *About Foreign,*
794     *Commonwealth & Development Office travel advice*. GOV.UK.
795     https://www.gov.uk/guidance/about-foreign-commonwealth-development-office-travel-
796     advice#when-we-advise-against-foreign-travel
797   Fortes, E. (2002). State-of-Knowledge Assessments of Each Thematic Group – SEAWEEDS. In
798     P. S. Ong, L. E. Afuang, & R. G. Rosell-Ambal (Eds.), *Philippine Biodiversity*
799     *Conservation Priorities: A Second Iteration of the National Biodiversity Strategy and*
800     *Action Plan* (p. 49). Department of Environment and Natural Resources, Conservation
801     International, University of the Philippines Center for Integrative and Development Studies,
802     and Foundation for the Philippine Environment.
803   Fortes, M. (2002). State-of-Knowledge Assessments of Each Thematic Group – SEAGRASSES.
804     In P. S. Ong, L. E. Afuang, & R. G. Rosell-Ambal (Eds.), *Philippine Biodiversity*
805     *Conservation Priorities: A Second Iteration of the National Biodiversity Strategy and*
806     *Action Plan* (pp. 49–50). Department of Environment and Natural Resources, Conservation
807     International, University of the Philippines Center for Integrative and Development Studies,
808     and Foundation for the Philippine Environment.

809 Gapud, V. P. (2002). State-of-Knowledge Assessments of Each Thematic Group –
810     ARTHROPODS. In P. S. Ong, L. S. Afuang, & R. G. Rosell-Ambal (Eds.), *Philippine*
811     *Biodiversity Conservation Priorities: A Second Iteration of the National Biodiversity*
812     *Strategy and Action Plan* (pp. 25–26). Department of Environment and Natural Resources,
813     Conservation International, University of the Philippines Center for Integrative and
814     Development Studies, and Foundation for the Philippine Environment.
815 GBIF.org. (2020a). *GBIF Occurrence Download*.
816     https://doi.org/https://doi.org/10.15468/dl.ks8qmp
817 GBIF.org. (2020b). *GBIF Occurrence Download*.
818     https://doi.org/https://doi.org/10.15468/dl.v7bjna
819 Government of Canada. (2021, August 21). *Travel advice and advisories for Philippines*.
820     https://travel.gc.ca/destinations/philippines
821 Halpern, B. S., Frazier, M., Potapenko, J., Casey, K. S., Koenig, K., Longo, C., Lowndes, J. S.,
822     Rockwood, R. C., Selig, E. R., Selkoe, K. A., & Walbridge, S. (2015). Spatial and temporal
823     changes in cumulative human impacts on the world's ocean. *Nature Communications*,
824     *6*(May), 1–7. https://doi.org/10.1038/ncomms8615
825 Halpern, B. S., Walbridge, S., Selkoe, K. A., Kappel, C. V., Micheli, F., D'Agrosa, C., Bruno, J.
826     F., Casey, K. S., Ebert, C., Fox, H. E., Fujita, R., Heinemann, D., Lenihan, H. S., Madin, E.
827     M. P., Perry, M. T., Selig, E. R., Spalding, M., Steneck, R., & Watson, R. (2008). A global
828     map of human impact on marine ecosystems. *Science*, *319*(5865), 948–952.
829     https://doi.org/10.1126/science.1149345
830 Hartvig, I., Czako, M., Kjær, E. D., Nielsen, L. R., & Theilade, I. (2015). The use of DNA
831     barcoding in identification and conservation of rosewood (Dalbergia spp.). *PLoS ONE*,
832     *10*(9). https://doi.org/10.1371/journal.pone.0138231
833 Harvey, M. G., Bravo, G. A., Claramunt, S., Cuervo, A. M., Derryberry, G. E., Battilana, J.,
834     Seeholzer, G. F., Shearer McKay, J., O'Meara, B. C., Faircloth, B. C., Edwards, S. V.,
835     Pérez-Emán, J., Moyle, R. G., Sheldon, F. H., Aleixo, A., Smith, B. T., Chesser, R. T.,
836     Silveira, L. F., Cracraft, J., … Derryberry, E. P. (2020). The evolution of a tropical
837     biodiversity hotspot. *Science*, *370*(6522), 1343–1348.
838     https://doi.org/10.1126/science.aaz6970
839 Hebert, P. D. N., & Gregory, T. R. (2005). The promise of DNA barcoding for taxonomy.
840     *Systematic Biology*, *54*(5), 852–859. https://doi.org/10.1080/10635150500354886
841 Herrick, J. E., Schuman, G. E., & Rango, A. (2006). Monitoring ecological processes for
842     restoration projects. *Journal for Nature Conservation*, *14*(3–4), 161–171.
843     https://doi.org/10.1016/j.jnc.2006.05.001
844 Hijmans, R. J. (2020). *raster: Geographic Data Analysis and Modeling* (R package version 3.4-
845     5). https://cran.r-project.org/package=raster
846 Hoffmann, A. A., & Sgró, C. M. (2011). Climate change and evolutionary adaptation. *Nature*,
847     *470*(7335), 479–485. https://doi.org/10.1038/nature09670
848 Hoveka, L. N., van der Bank, M., Bezeng, B. S., & Davies, T. J. (2020). Identifying biodiversity
849     knowledge gaps for conserving South Africa's endemic flora. *Biodiversity and*
850     *Conservation*, *29*(9–10), 2803–2819. https://doi.org/10.1007/s10531-020-01998-4
851 Keyse, J., Crandall, E. D., Toonen, R. J., Meyer, C. P., Treml, E. A., & Riginos, C. (2014). The
852     scope of published population genetic data for Indo-Pacific marine fauna and future
853     research opportunities in the region. *Bulletin of Marine Science*, *90*(1), 47–78.
854     https://doi.org/10.5343/bms.2012.1107

855 Lagunzad, D. A., Co, L. L., & Navarro, J. M. (2002). State-of-Knowledge Assessments of Each
856     Thematic Group – TERRESTRIAL PLANTS. In P. S. Ong, L. E. Afuang, & R. G. Rosell-
857     Ambal (Eds.), *Philippine Biodiversity Conservation Priorities: A Second Iteration of the*
858     *National Biodiversity Strategy and Action Plan* (pp. 24–25). Department of Environment
859     and Natural Resources, Conservation International, University of the Philippines Center for
860     Integrative and Development Studies, and Foundation for the Philippine Environment.
861 Laikre, L., Allendorf, F. W., Aroner, L. C., Baker, C. S., Gregovich, D. P., Hansen, M. M.,
862     Jackson, J. A., Kendall, K. C., McKelvey, K., Neel, M. C., Olivieri, I., Ryman, N.,
863     Schwartz, M. K., Bull, R. S., Stetz, J. B., Tallmon, D. A., Taylor, B. L., Vojta, C. D.,
864     Waller, D. M., & Waples, R. S. (2010). Neglect of genetic diversity in implementation of
865     the convention on biological diversity: Conservation in practice and policy. *Conservation
866     Biology*, *24*(1), 86–88. https://doi.org/10.1111/j.1523-1739.2009.01425.x
867 Maralit, B. A., Aguila, R. D., Ventolero, M. F. H., Perez, S. K. L., Willette, D. A., & Santos, M.
868     D. (2013). Detection of mislabeled commercial fishery by-products in the Philippines using
869     DNA barcodes and its implications to food traceability and safety. *Food Control*, *33*(1),
870     119–125. https://doi.org/10.1016/j.foodcont.2013.02.018
871 Margules, C. R., & Pressey, R. L. (2000). Systematic conservation planning. *Nature*, *405*, 243–
872     253. https://doi.org/https://doi.org/10.1038/35012251
873 Matias, A. M. A., & Riginos, C. (2018). Revisiting the "Centre Hypotheses" of the Indo-West
874     Pacific: Idiosyncratic genetic diversity of nine reef species offers weak support for the Coral
875     Triangle as a centre of genetic biodiversity. *Journal of Biogeography*, *45*(8), 1806–1817.
876     https://doi.org/10.1111/jbi.13376
877 Meyer, C., Kreft, H., Guralnick, R., & Jetz, W. (2015). Global priorities for an effective
878     information basis of biodiversity distributions. *Nature Communications*, *6*.
879     https://doi.org/10.1038/ncomms9221
880 Mittermeler, R. A., & Mittermeler, C. G. (1997). *Megadiversity: Earth's Biologically Wealthiest
881     Nations* (R. A. Mittermeler & C. G. Mittermeler (eds.)). CEMEX.
882 Moritz, C. (2002). Strategies to protect biological diversity and the evolutionary processes that
883     sustain it. *Systematic Biology*, *51*(2), 238–254. https://doi.org/10.1080/10635150252899752
884 Myers, N., Mittermeler, R. A., Mittermeler, C. G., Da Fonseca, G. A. B., & Kent, J. (2000).
885     Biodiversity hotspots for conservation priorities. *Nature*, *403*(6772), 853–858.
886     https://doi.org/10.1038/35002501
887 Neuwirth, E. (2014). *RColorBrewer: ColorBrewer Palettes* (R package version 1.1-2).
888     https://cran.r-project.org/package=RColorBrewer
889 Newbold, T. (2010). Applications and limitations of museum data for conservation and ecology,
890     with particular attention to species distribution models. *Progress in Physical Geography*,
891     *34*(1), 3–22. https://doi.org/10.1177/0309133309355630
892 Oliveira, U., Paglia, A. P., Brescovit, A. D., de Carvalho, C. J. B., Silva, D. P., Rezende, D. T.,
893     Leite, F. S. F., Batista, J. A. N., Barbosa, J. P. P. P., Stehmann, J. R., Ascher, J. S., de
894     Vasconcelos, M. F., De Marco, P., Löwenberg-Neto, P., Dias, P. G., Ferro, V. G., & Santos,
895     A. J. (2016). The strong influence of collection bias on biodiversity knowledge shortfalls of
896     Brazilian terrestrial biodiversity. *Diversity and Distributions*, *22*(12), 1232–1244.
897     https://doi.org/10.1111/ddi.12489
898 Ong, P. S. (2002). Current Status and Prospects of Protected Areas in the Light of the Philippine
899     Biodiversity Conservation Priorities. *Proceedings of IUCN/WCPA-EA-4 Taipei Conference*,
900     95–126.

901    Page, L. M., Macfadden, B. J., Fortes, J. A., Soltis, P. S., & Riccardi, G. (2015). Digitization of
902        Biodiversity Collections Reveals Biggest Data on Biodiversity. *BioScience*, *65*(9), 841–842.
903        https://doi.org/10.1093/biosci/biv104
904    Pazartzi, T., Siaperopoulou, S., Gubili, C., Maradidou, S., Loukovitis, D., Chatzispyrou, A.,
905        Griffiths, A. M., Minos, G., & Imsiridou, A. (2019). High levels of mislabeling in shark
906        meat – Investigating patterns of species utilization with DNA barcoding in Greek retailers.
907        *Food Control*, *98*(September 2018), 179–186.
908        https://doi.org/10.1016/j.foodcont.2018.11.019
909    Pebesma, E. J., & Bivand, R. S. (2005). Classes and methods for spatial data in R. *R News*, *5*(2).
910        https://cran.r-project.org/doc/Rnews/
911    Petersen, T. K., Speed, J. D. M., Grøtan, V., & Austrheim, G. (2021). Species data for
912        understanding biodiversity dynamics: The what, where and when of species occurrence data
913        collection. *Ecological Solutions and Evidence*, *2*(1), 1–17. https://doi.org/10.1002/2688-
914        8319.12048
915    Philippine Statistics Authority. (2020). *Philippine Standard Geographic Code (PSGC)*.
916        https://psa.gov.ph/classification/psgc/
917    Posa, M. R. C., Diesmos, A. C., Sodhi, N. S., & Brooks, T. M. (2008). Hope for threatened
918        tropical biodiversity: Lessons from the Philippines. *BioScience*, *58*(3), 231–240.
919        https://doi.org/10.1641/B580309
920    Possingham, H. P., & Wilson, K. A. (2005). Biodiversity: Turning up the heat on hotspots.
921        *Nature*, *436*(7053), 919–920. https://doi.org/10.1038/436919a
922    Primavera, J. (2002). State-of-Knowledge Assessments of Each Thematic Group –
923        MANGROVES. In P. S. Ong, L. E. Afuang, & R. G. Rosell-Ambal (Eds.), *Philippine*
924        *Biodiversity Conservation Priorities: A Second Iteration of the National Biodiversity*
925        *Strategy and Action Plan* (p. 49). Department of Environment and Natural Resources,
926        Conservation International, University of the Philippines Center for Integrative and
927        Development Studies, and Foundation for the Philippine Environment.
928    Purvis, A., & Hector, A. (2000). Getting the measure of biodiversity. *Nature*, *405*(6783), 212–
929        219. https://doi.org/10.1038/35012221
930    Reich, P. B., Tilman, D., Isbell, F., Mueller, K., Hobbie, S. E., Flynn, D. F. B., & Eisenhauer, N.
931        (2012). Impacts of biodiversity loss escalate through time as redundancy fades. *Science*,
932        *336*(6081), 589–592. https://doi.org/10.1126/science.1217909
933    Richardson, D. M., & Whittaker, R. J. (2010). Conservation biogeography - foundations,
934        concepts and challenges. *Diversity and Distributions*, *16*(3), 313–320.
935        https://doi.org/10.1111/j.1472-4642.2010.00660.x
936    Santangeli, A., Wistbacka, R., Hanski, I. K., & Laaksonen, T. (2013). Ineffective enforced
937        legislation for nature conservation: A case study with Siberian flying squirrel and forestry in
938        a boreal landscape. *Biological Conservation*, *157*, 237–244.
939        https://doi.org/10.1016/j.biocon.2012.09.012
940    Sarkar, S., & Margules, C. (2002). Operationalizing biodiversity for conservation planning.
941        *Journal of Biosciences*, *27*(4 SUPPL. 2), 299–308. https://doi.org/10.1007/BF02704961
942    Selig, E. R., Turner, W. R., Troëng, S., Wallace, B. P., Halpern, B. S., Kaschner, K., Lascelles,
943        B. G., Carpenter, K. E., & Mittermeier, R. A. (2014). Global priorities for marine
944        biodiversity conservation. *PLoS ONE*, *9*(1), 1–11.
945        https://doi.org/10.1371/journal.pone.0082898
946    Senate of the Philippines. (2017, May 13). *Press Release - Sen. Bam: PH lacks 19,000 scientists,*

947     *lags behind in R&D*. Senate of the Philippines - 18th Congress.
948         http://legacy.senate.gov.ph/press_release/2017/0513_aquino1.asp
949     Sgrò, C. M., Lowe, A. J., & Hoffmann, A. A. (2011). Building evolutionary resilience for
950         conserving biodiversity under climate change. *Evolutionary Applications*, *4*(2), 326–337.
951         https://doi.org/10.1111/j.1752-4571.2010.00157.x
952     Siler, C. D., Alcala, A. C., Diesmos, A. C., & Brown, R. M. (2009). A new species of limestone-
953         forest frog, genus platymantis (Amphibia: Anura: Ceratobatrachidae) from Eastern Samar
954         Island, Philippines. *Herpetologica*, *65*(1), 92–104. https://doi.org/10.1655/08-040R.1
955     Titley, M. A., Snaddon, J. L., & Turner, E. C. (2017). Scientific research on animal biodiversity
956         is systematically biased towards vertebrates and temperate regions. *PLoS ONE*, *12*(12), 1–
957         14. https://doi.org/10.1371/journal.pone.0189577
958     Troudet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R., & Legendre, F. (2017). Taxonomic bias
959         in biodiversity data and societal preferences. *Scientific Reports*, *7*(9132), 1–14.
960         https://doi.org/10.1038/s41598-017-09084-6
961     U.S. Department of State. (2021, June 16). *Philippines Travel Advisory*. Travel.State.Gov.
962         https://travel.state.gov/content/travel/en/traveladvisories/traveladvisories/philippines-travel-
963         advisory.html
964     van der Plas, F. (2019). Biodiversity and ecosystem functioning in naturally assembled
965         communities. *Biological Reviews*, *94*(4), 1220–1245. https://doi.org/10.1111/brv.12499
966     Worm, B., Barbier, E. B., Beaumont, N., Duffy, J. E., Folke, C., Halpern, B. S., Jackson, J. B. C.,
967         Lotze, H. K., Micheli, F., Palumbi, S. R., Sala, E., Selkoe, K. A., Stachowicz, J. J., &
968         Watson, R. (2006). Impacts of biodiversity loss on ocean ecosystem services. *Science*,
969         *314*(5800), 787–790. https://doi.org/10.1126/science.1132294
970     Zapanta, B. R., Achondo, M. J. M. M., Raganas, A. F. M., Camino, F. A., Delima, A. G. D.,
971         Mantiquilla, J. A., Puentespina, R. P., & Salvaña, F. R. P. (2019). Species richness of trees
972         in disturbed habitats within a protected area and its implications for conservation: The case
973         of Mt. Apo Natural Park, Mindanao Island, Philippines. *Biodiversitas*, *20*(7), 2081–2091.
974         https://doi.org/10.13057/biodiv/d200740
975     Zrzavý, J., Říha, P., Piálek, L., & Janouškovec, J. (2009). Phylogeny of Annelida
976         (Lophotrochozoa): Total-evidence analysis of morphology and six genes. *BMC*
977         *Evolutionary Biology*, *9*(1), 1–14. https://doi.org/10.1186/1471-2148-9-189

**Figure Legends**

**Figure 1. Summary of barcode records associated with specific gene markers and issues encountered while manually parsing through the descriptive information on sampling locality.** For graph **A**, the genetic summary of the available barcode records focuses on the gene markers of interest used in the examination for metadata gaps, taxonomic biases, and spatial biases in DNA barcode data on animal and plant taxa sampled in the Philippines were the following: cytochrome b (CYTB), cytochrome oxidase c subunit I (COI), internal transcribed spacer 2 (ITS2), ribulose-1,5-biphosphate carboxylase (rbcL), and maturase K (matK). For graph **B**, the geolocation issues resulted in the descriptions of the sampling location (particularly in terms of administrative units) being unclear or in some cases, inconclusive. The categories include misspelled (incorrect spelling), none (no major issue), mixed (more than one issue), unspecified (somewhat informative but still vague), unknown (completely not informative), multiple (provided more than one location), and mismatch (discrepancies between the administrative units provided). This dataset includes the records with *NA* entries for country sampled (for **A** and **B**) and those that had additional information on the geolocation other than the coordinates (for **B** only).

**Figure 2. Relationship between the percentage of barcode records identified at the species level and the proportion of documented species (represented in species occurrence data) that currently have DNA barcode data available.** This relationship was evaluated for each known animal (orange) and plant (green) taxonomic group represented in the Philippine barcode data at the phylum/division (**A**), class (**B**), order (**C**), and family (**D**) levels. This dataset includes the records with *NA* entries for country sampled.

**Figure 3. Relationship between the amount of genetic and species data associated with each known animal and plant taxonomic group represented in the Philippine biodiversity data at different taxonomic levels.** This relationship was evaluated for each known animal (orange) and plant (green) taxonomic group represented in the Philippine barcode data at the phylum/division (**A**), class (**B**), order (**C**), and family (**D**) levels. Values were transformed logarithmically prior to plotting however, taxa with zero (0) records in either genetic or species data were assigned the value of negative one (-1). Dashed lines represent the 5th and 95th percentiles for genetic (horizontal) and species (vertical) data. This dataset includes the records with *NA* entries for country sampled.

**Figure 4. Maps of the sampling distribution of barcode and species occurrence data on animal and plant taxa across the Philippines and the relationship between the two datasets in terms of province.** For both maps (**A** – barcode data and **B** – species occurrence data), records on marine specimens were assigned to a specific province based on which corresponding centroid has the shortest distance from the given sampling coordinates (if available). Also, values presented in the maps represent the number of records in the thousands. In the scatter plot (**C**), values were transformed logarithmically and provinces with zero (0) records in either genetic or species data were assigned the value of negative one (-1). Dashed lines represent the 5th and 95th percentiles for genetic (horizontal) and species (vertical) data. The barcode dataset includes the records with *NA* entries for country sampled.

**Figure 5. Map of the distribution of barcode data on Philippine animal and plant biodiversity contributed by different countries across the world and their contribution to documenting**

1024    **efforts across the years.** For map **A**, contribution was based on the institution that holds the
1025    copyright to the image associated with the records while for the graphs, it was based on the
1026    collection of samples, starting from the 1990s (**B**) and submission of barcode data, starting from
1027    the 2000s (**C**) by foreign countries (violet) and the Philippines (red). Trendlines in the graphs
1028    represent the average, "best" fitted line. This dataset includes the records with *NA* entries for
1029    country sampled.
1030
1031    **Figure 6. Heatmap matrix showcasing the relationship between the number of barcode**
1032    **records associated with regions that have been sampled and the regions of local institutions**
1033    **that contributed the data.** There are officially seventeen regions in the Philippines, represented
1034    by the Philippine map (**A**), with non-numerical regions labelled as follows: *ca*, Cordillera
1035    Administrative Region (CAR); *mm*, National Capital Region (NCR or also referred to as Metro
1036    Manila); and *br*, Bangsamoro Autonomous Region in Muslim Mindanao (BARMM). Regions are
1037    also divided based on their island groups – namely Luzon (red), Visayas (yellow), and Mindanao
1038    (blue). For matrix **B**, contribution was based on the institution that holds the copyright to the image
1039    associated with the records. Regions along the x- and y-axis are sorted to provide spatial context,
1040    with the map as a reference. The diagonal line represents the "ideal" scenario wherein the region
1041    serving as the processing center of barcode data can sufficiently sample its own local area. This
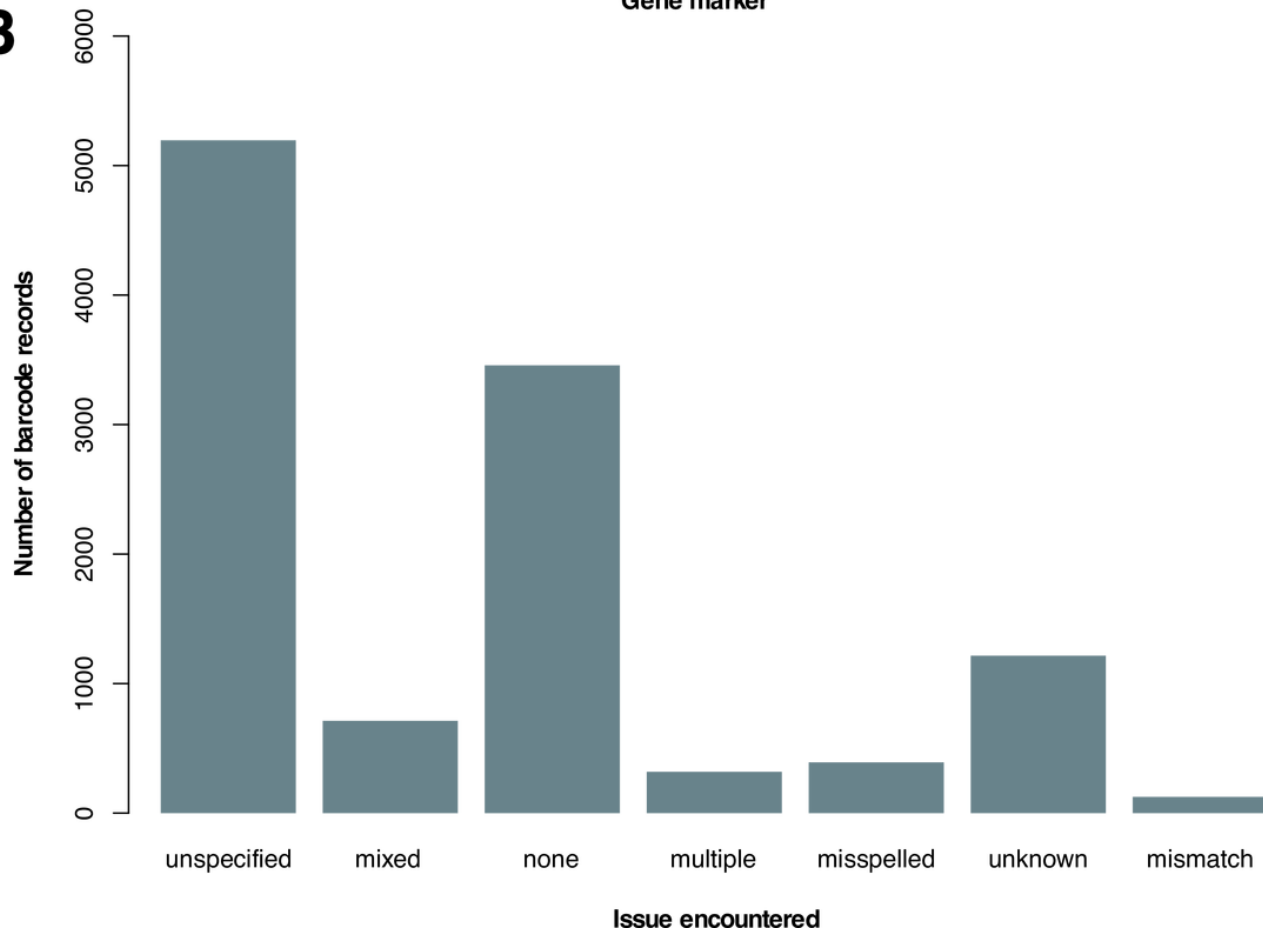1042    dataset includes the records with *NA* entries for country sampled.

# Figure 1

Figure 1

Summary of barcode records associated with specific gene markers and issues encountered while manually parsing through the descriptive information on sampling locality
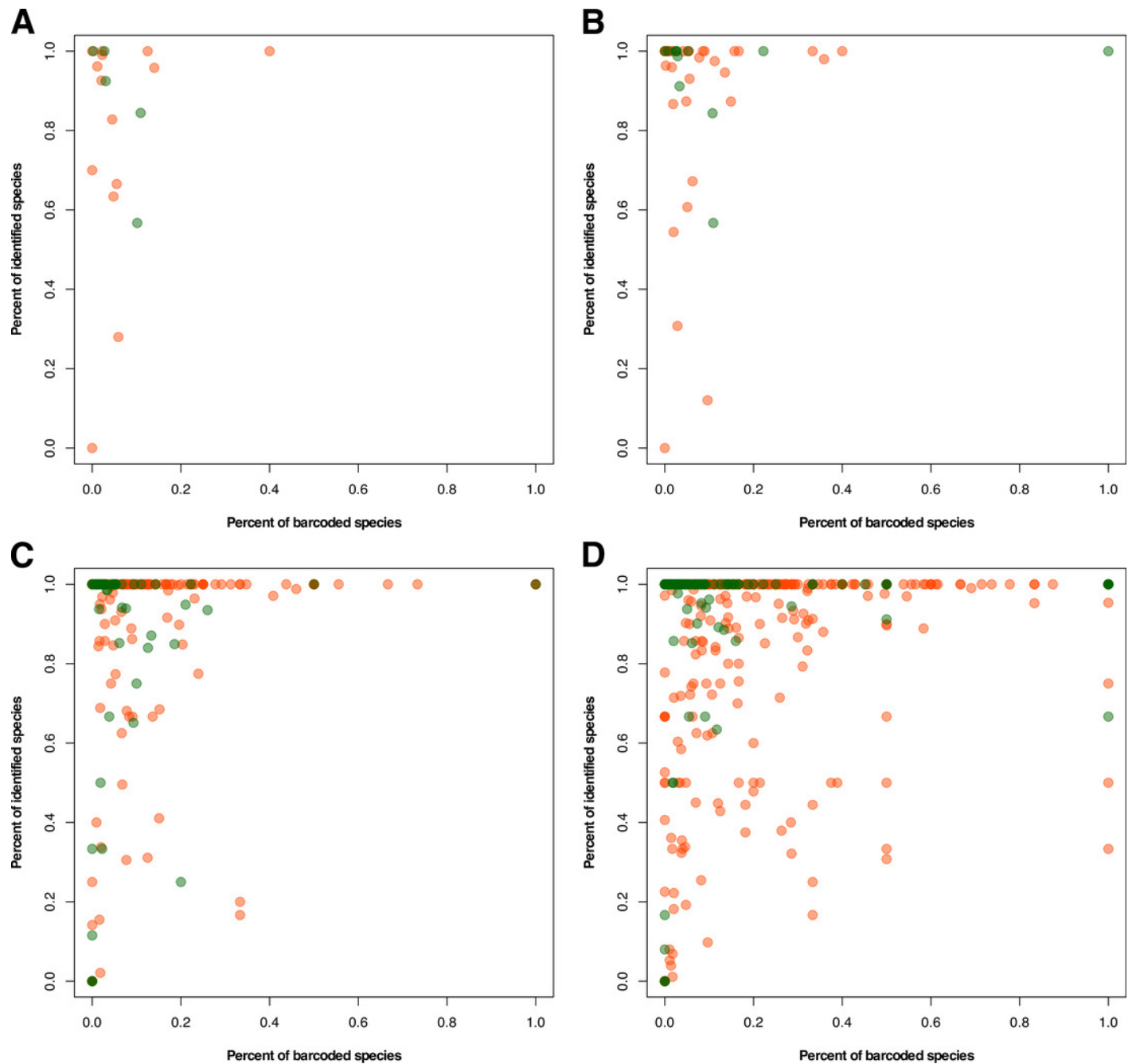
# Figure 2

Figure 2

Relationship between the percentage of barcode records identified at the species level and the proportion of documented species (represented in species occurrence data) that currently have DNA barcode data available
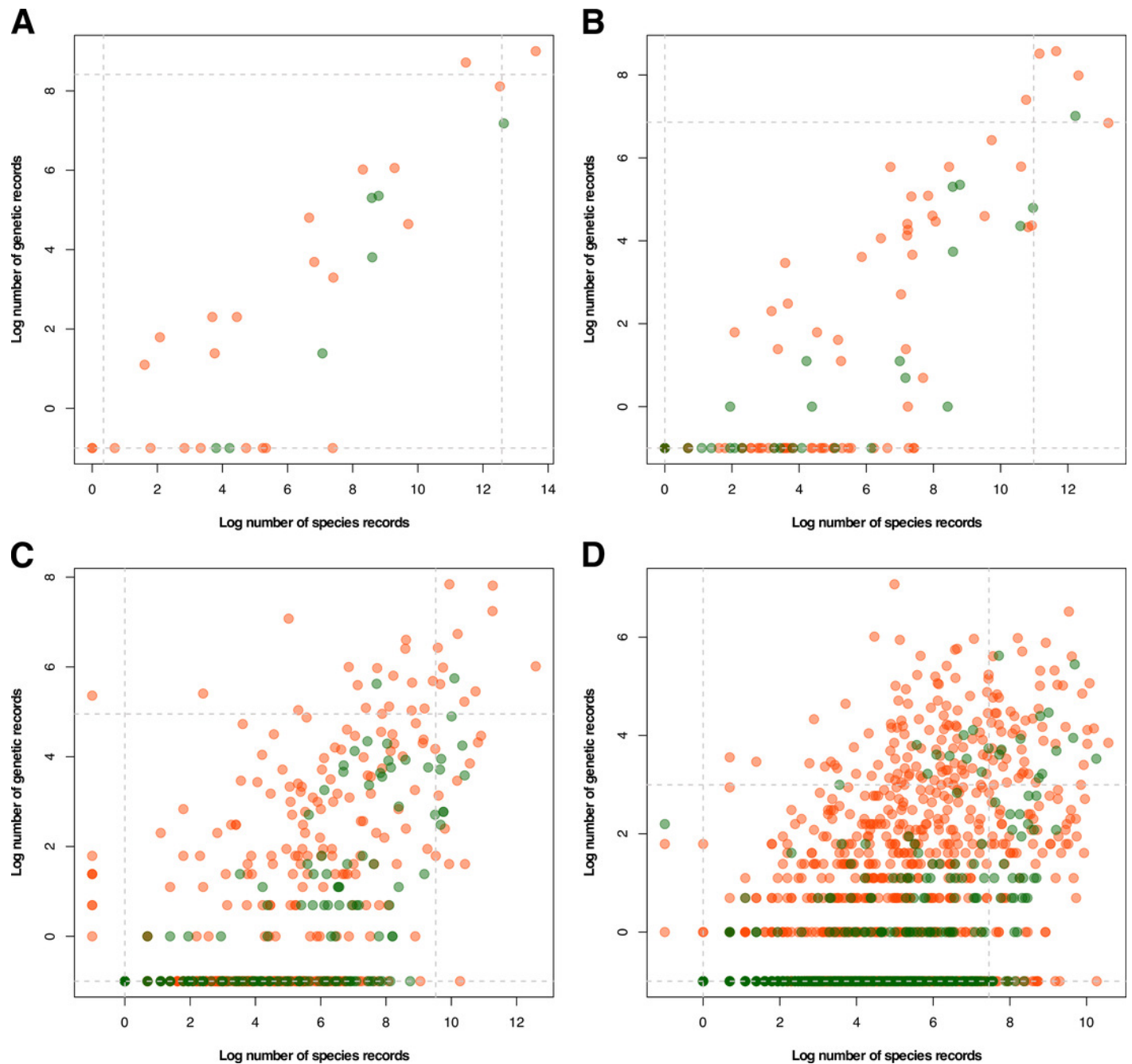
# Figure 3

Figure 3

Relationship between the amount of genetic and species data associated with each known animal and plant taxonomic group represented in the Philippine biodiversity data at different taxonomic levels
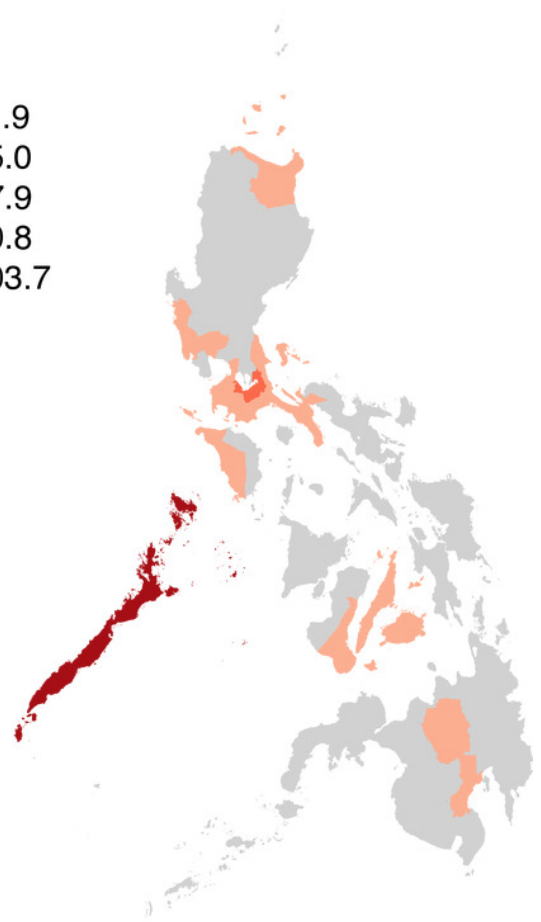
# Figure 4

Figure 4

Maps of the sampling distribution of barcode and species occurrence data on animal and plant taxa across the Philippines and the relationship between the two datasets in terms of province
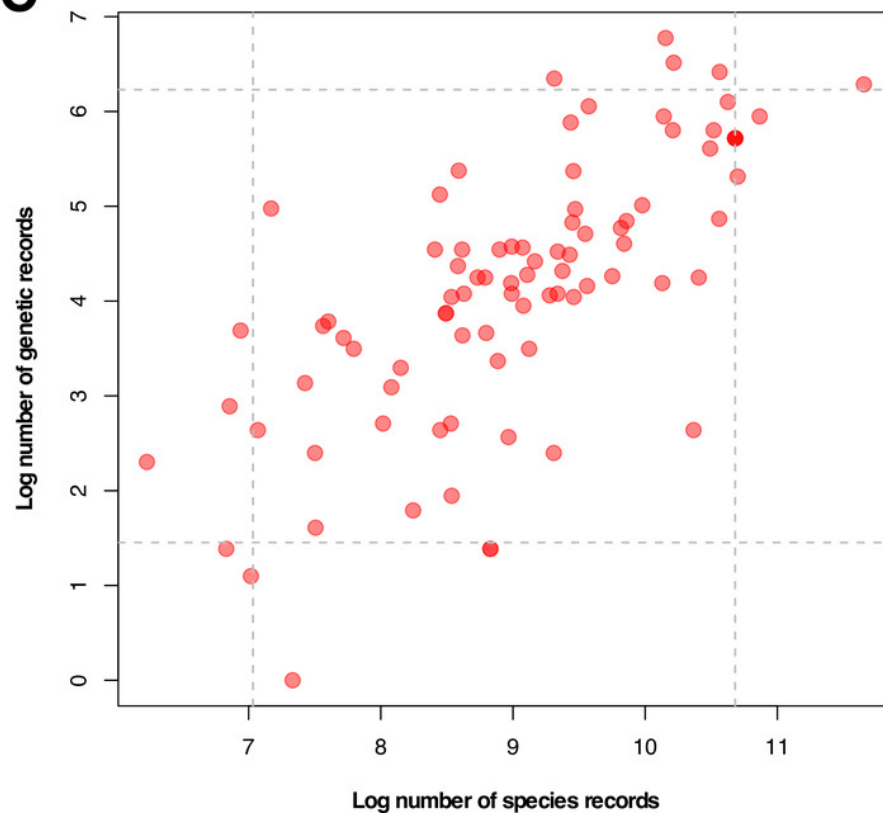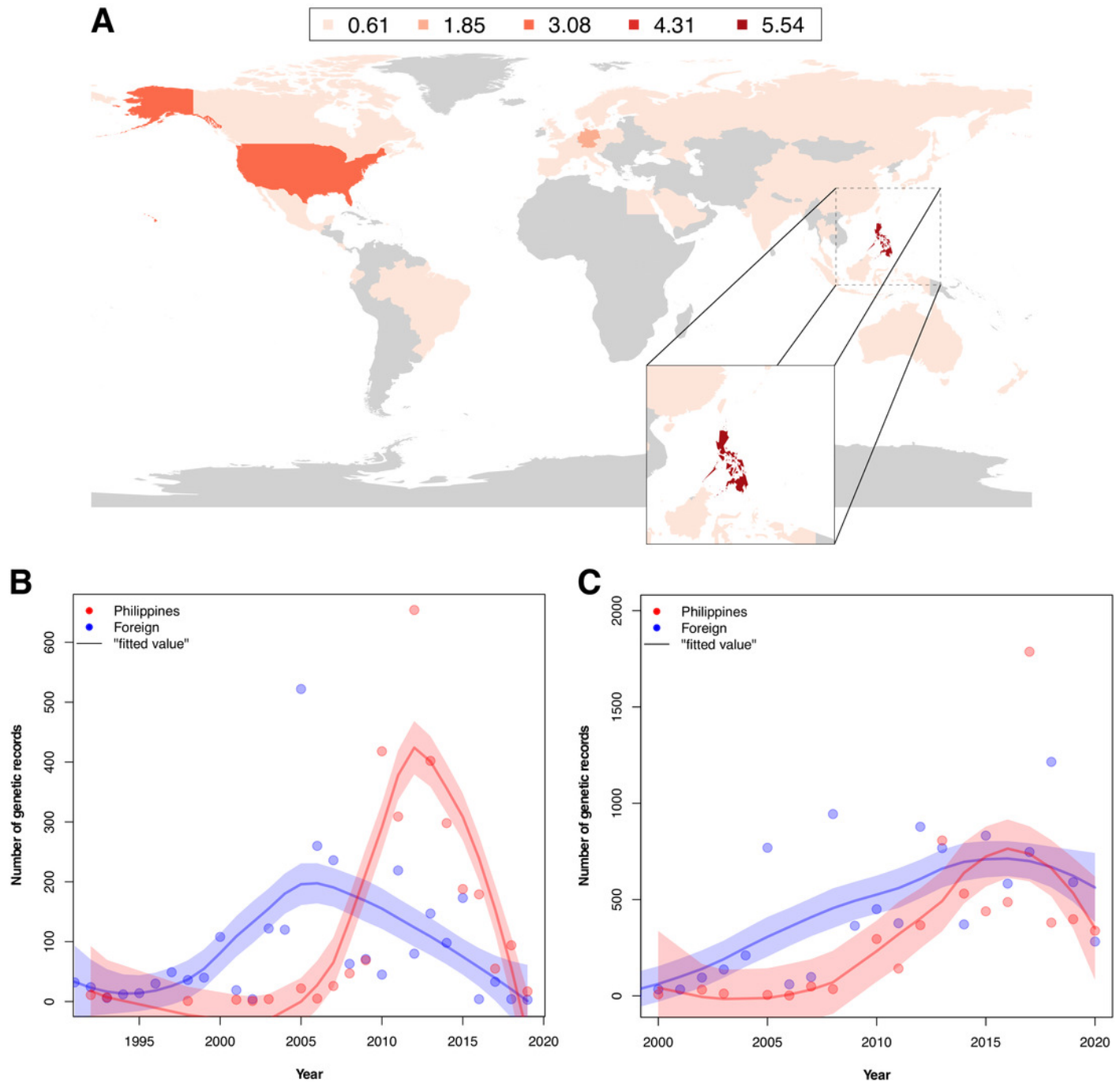
# Figure 5

Figure 5

Map of the distribution of barcode data on Philippine animal and plant biodiversity contributed by different countries across the world and their contribution to documenting efforts across the years

# Figure 6

Figure 6

Heatmap matrix showcasing the relationship between the number of barcode records associated with regions that have been sampled and the regions of local institutions that contributed the data