

First submission

Guidance from your Editor

Please submit by **28 Mar 2020** for the benefit of the authors (and your \$200 publishing discount) .



Structure and Criteria

Please read the 'Structure and Criteria' page for general guidance.



Raw data check

Review the raw data.



Image check

Check that figures and images have not been inappropriately manipulated.

Privacy reminder: If uploading an annotated PDF, remove identifiable information to remain anonymous.

Files

Download and review all files from the [materials page](#).

4 Figure file(s)

1 Latex file(s)

3 Other file(s)



Structure and Criteria

Structure your review

The review form is divided into 5 sections. Please consider these when composing your review:

1. BASIC REPORTING
2. EXPERIMENTAL DESIGN
3. VALIDITY OF THE FINDINGS
4. General comments
5. Confidential notes to the editor

You can also annotate this PDF and upload it as part of your review

When ready [submit online](#).

Editorial Criteria

Use these criteria points to structure your review. The full detailed editorial criteria is on your [guidance page](#).

BASIC REPORTING

- Clear, unambiguous, professional English language used throughout.
- Intro & background to show context. Literature well referenced & relevant.
- Structure conforms to [Peerj standards](#), discipline norm, or improved for clarity.
- Figures are relevant, high quality, well labelled & described.
- Raw data supplied (see [Peerj policy](#)).

EXPERIMENTAL DESIGN

- Original primary research within [Scope of the journal](#).
- Research question well defined, relevant & meaningful. It is stated how the research fills an identified knowledge gap.
- Rigorous investigation performed to a high technical & ethical standard.
- Methods described with sufficient detail & information to replicate.

VALIDITY OF THE FINDINGS

- Impact and novelty not assessed. Negative/inconclusive results accepted. *Meaningful* replication encouraged where rationale & benefit to literature is clearly stated.
- All underlying data have been provided; they are robust, statistically sound, & controlled.
- Speculation is welcome, but should be identified as such.
- Conclusions are well stated, linked to original research question & limited to supporting results.



The best reviewers use these techniques

Tip

Example

Support criticisms with evidence from the text or from other sources

Smith et al (J of Methodology, 2005, V3, pp 123) have shown that the analysis you use in Lines 241-250 is not the most appropriate for this situation. Please explain why you used this method.

Give specific suggestions on how to improve the manuscript

Your introduction needs more detail. I suggest that you improve the description at lines 57- 86 to provide more justification for your study (specifically, you should expand upon the knowledge gap being filled).

Comment on language and grammar issues

The English language should be improved to ensure that an international audience can clearly understand your text. Some examples where the language could be improved include lines 23, 77, 121, 128 - the current phrasing makes comprehension difficult.

Organize by importance of the issues, and number your points

- 1. Your most important issue*
- 2. The next most important item*
- 3. ...*
- 4. The least important points*

Please provide constructive criticism, and avoid personal opinions

I thank you for providing the raw data, however your supplemental files need more descriptive metadata identifiers to be useful to future readers. Although your results are compelling, the data analysis should be improved in the following ways: AA, BB, CC

Comment on strengths (as well as weaknesses) of the manuscript

I commend the authors for their extensive data set, compiled over many years of detailed fieldwork. In addition, the manuscript is clearly written in professional, unambiguous language. If there is a weakness, it is in the statistical analysis (as I have noted above) which should be improved upon before Acceptance.

No one-size-fits-all solution to clean GBIF

Alexander Zizka^{Corresp., 1, 2}, **Fernanda Antunes Carvalho**³, **Alice Calvente**⁴, **Mabel Rocio Baez-Lizarazo**⁵, **Andressa Cabral**⁶, **Jéssica Fernanda Ramos Coelho**⁴, **Matheus Colli-Silva**⁶, **Mariana Ramos Fantinati**⁴, **Moabe Ferreira Fernandes**⁷, **Thais Ferreira-Araújo**⁴, **Fernanda Gondim Lambert Moreira**⁴, **Nathália Michelly da Cunha Santos**⁴, **Tiago Andrade Borges Santos**⁷, **Renata Clícia dos Santos-Costa**⁴, **Filipe Cabreirinha Serrano**⁸, **Ana Paula Alves da Silva**⁴, **Arthur de Souza Soares**⁴, **Paolla Gabryelle Cavalcante de Souza**⁴, **Eduardo Calisto Tomaz**⁴, **Valéria Fonseca Vale**⁴, **Tiago Luiz Vieira**⁷, **Alexandre Antonelli**^{9,10,11}

¹ German Centre for Integrative Biodiversity Research Halle Jena Leipzig (iDiv), Leipzig, Germany

² Naturalis Biodiversity Center, Leiden, The Netherlands

³ Departamento de Genética, Ecologia e Evolução, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

⁴ Departamento de Botânica e Zoologia, Universidade Federal do Rio Grande do Norte, Natal, Brazil

⁵ Departamento de Botânica, Universidade Federal do Rio Grande do Sul, Porto Alegre, Brazil

⁶ Departamento de Botânica, Universidade de São Paulo, São Paulo, Brazil

⁷ Departamento de Ciências Biológicas, Universidade Estadual de Feira de Santana, Feira de Santana, Brazil

⁸ Departamento de Ecologia, Universidade de São Paulo, São Paulo, Brazil

⁹ Gothenburg Global Biodiversity Centre, University of Gothenburg, Gothenburg, Sweden

¹⁰ Department for Biological and Environmental Sciences, University of Gothenburg, Gothenburg, Sweden

¹¹ Royal Botanic Gardens Kew, Richmond, United Kingdom

Corresponding Author: Alexander Zizka
Email address: alexander.zizka@idiv.de

Species occurrence records provide the basis for many biodiversity studies. They derive from geo-referenced specimens deposited in natural history collections and visual observations, such as those obtained through various mobile applications. Given the rapid increase in availability of such data, the control of quality and accuracy constitutes a particular concern. Automatic flagging and filtering are a scalable and reproducible means to identify potentially problematic records in datasets from public databases such as the Global Biodiversity Information Facility (GBIF; www.gbif.org). However, it is unclear how much data may be lost by filtering, whether the same tests should be applied across all taxonomic groups, and what is the effect of filtering for common downstream analyses. Here, we evaluate the effect of 13 recently proposed filters on the inference of species richness patterns and automated conservation assessments for 18 Neotropical taxa including animals, fungi, and plants, terrestrial and marine, downloaded from GBIF. We find that 29-90% of the records are potentially erroneous, with large variation across taxonomic groups. Tests for duplicated information, collection year, basis of record as well as urban areas and coordinates for terrestrial taxa in the sea or marine taxa on land have the greatest effect. While many flagged records might not be de facto erroneous, they could be overly imprecise and increase uncertainty in downstream analyses. Automated

flagging can help in identifying problematic records, but requires customization of which tests and thresholds should be applied to the taxonomic group and geographic area under focus. Our results stress the importance of thorough exploration of the meta-data associated with species records for biodiversity research.

No one-size-fits-all solution to clean GBIF

Alexander Zizka^{1,2}, Fernanda Antunes Carvalho³, Alice Calvente⁴, Mabel Rocio Baez-Lizarazo⁵, Andressa Cabral⁶, Jéssica Fernanda Ramos Coelho⁴, Matheus Colli-Silva⁶, Mariana Ramos Fantinati⁴, Moabe Ferreira Fernandes⁷, Thais Ferreira-Araújo⁴, Fernanda Gondim Lambert Moreira⁴, Nathália Michelly da Cunha Santos⁴, Tiago Andrade Borges Santos⁷, Renata Clicia dos Santos-Costa⁴, Filipe Cabreirinha Serrano⁸, Ana Paula Alves da Silva⁴, Arthur de Souza Soares⁴, Paolla Gabryelle Cavalcante de Souza⁴, Eduardo Calisto Tomaz⁴, Valéria Fonseca Vale⁴, Tiago Luiz Vieira⁷, and Alexandre Antonelli^{9,10,11}

¹sDiv, German Center for Integrative Biodiversity Research Halle-Jena-Leipzig (iDiv), Leipzig, Germany

²Naturalis Biodiversity Center, Leiden, The Netherlands

³Departamento de Genética, Ecologia e Evolução, Universidade Federal de Minas Gerais, Belo Horizonte, Brazil

⁴Departamento de Botânica e Zoologia, Universidade Federal do Rio Grande do Norte, Natal, Brazil

⁵Departamento de Botânica, Universidade Federal do Rio Grande do Sul, Porto Alegre, Brazil

⁶Departamento de Botânica, Universidade de São Paulo, São Paulo, Brazil

⁷Departamento de Ciências Biológicas, Universidade Estadual de Feira de Santana, Feira de Santana, Brazil

⁸Departamento de Ecologia, Universidade de São Paulo, São Paulo, Brazil

⁹Gothenburg Global Biodiversity Centre, University of Gothenburg, Gothenburg, Sweden

¹⁰Department for Biological and Environmental Sciences, University of Gothenburg, Gothenburg, Sweden

¹¹Royal Botanic Gardens Kew, Richmond, United Kingdom

Corresponding author:

Alexander Zizka^{1,2}

Email address: alexander.zizka@idiv.de

ABSTRACT

Species occurrence records provide the basis for many biodiversity studies. They derive from geo-referenced specimens deposited in natural history collections and visual observations, such as those obtained through various mobile applications. Given the rapid increase in availability of such data, the control of quality and accuracy constitutes a particular concern. Automatic flagging and filtering are a scalable and reproducible means to identify potentially problematic records in datasets from public databases such as the Global Biodiversity Information Facility (GBIF; www.gbif.org). However, it is unclear how much data may be lost by filtering, whether the same tests should be applied across all taxonomic groups, and what is the effect of filtering for common downstream analyses. Here, we evaluate the effect of 13 recently proposed filters on the inference of species richness patterns and automated conservation assessments for 18 Neotropical taxa including animals, fungi, and plants, terrestrial and marine, downloaded from GBIF. We find that 29-90% of the records are potentially erroneous, with large variation across taxonomic groups. Tests for duplicated information, collection year, basis of record as well as urban areas and coordinates for terrestrial taxa in the sea or marine taxa on land have the greatest effect. While many flagged records might not be *de facto* erroneous, they could be overly imprecise and increase uncertainty in downstream analyses. Automated flagging can help in identifying problematic records, but requires customization of which tests and thresholds should be applied to the taxonomic group and geographic area under focus. Our results stress the importance of thorough exploration of the meta-data associated with species records for biodiversity research.

INTRODUCTION

Publicly available species distribution data have become a crucial resource in biodiversity research, including studies in ecology, biogeography, systematics and conservation biology. In particular, the availability of digitized collections from museums and herbaria, and citizen science observations has increased drastically over the last few years. As of today, the largest public aggregator for geo-referenced species occurrences data, the Global Biodiversity Information Facility (www.gbif.org), provides access to more than 1.3 billion geo-referenced occurrence records for species from across the globe and the tree of life.

A central challenge to the use of these publicly available species occurrence data in research are erroneous geographic coordinates (Anderson et al. 2016). Errors mostly arise because public databases integrate records collected with different methodologies in different places, at different times; often without centralized curation and only rudimentary meta-data. For instance, erroneous coordinates caused by data-entry errors or automated geo-referencing from vague locality descriptions are common (Maldonado et al. 2015; Yesson et al. 2007) and cause recurrent problems such as records of terrestrial species in the sea, records with coordinates assigned to the centroids of political entities, or records of species in captivity (Zizka, Silvestro, et al. 2019).

Manual data cleaning based on expert knowledge can detect these issues, but it is only feasible on small taxonomic or geographic scales, and it is time-consuming and difficult to reproduce. As an alternative, automated flagging methods to identify potentially problematic records have been proposed as a scalable option, as they are able to deal with datasets containing up to millions of records and many different taxa. Those methods are usually based on geographic gazetteers (e.g., Chamberlain 2016; Zizka, Silvestro, et al. 2019; Jin and Yang 2020) or on additional data, such as environmental variables (Robertson, Visser, and Hui 2016). Additionally, filtering procedures based on record meta-data, such as collection year, record type, and coordinate precisions have been proposed to improve the suitability of publicly available occurrence records for biodiversity research (Zizka, Silvestro, et al. 2019).

Erroneous records are especially problematic in conservation, where stakes are high. Recently proposed methods for automated conservation assessments could support the formal assessment procedures for the global Red List of the International Union for the Conservation of Nature (IUCN) (Dauby et al. 2017; Bachman et al. 2011; Pelletier et al. 2018). These methods approximate species' range size, namely the Extent of Occurrence (EOO, which is the area of a convex hull polygon comprising all records of a species), the Area of Occupancy (AOO, which is the sum of the area actually occupied by a species, calculated based on a small-scale regular grid), and the number of locations for a preliminary conservation assessment following IUCN Criterion B ("Geographic range"). These methods have been used to propose preliminary global (Stévant et al. 2019; Zizka, Azevedo, et al. 2019) and regional (Schmidt et al. 2017; Cosiaux et al. 2018) Red List assessments. However, all metrics but especially EOO, are sensitive to individual records with erroneous coordinates. Automated conservation assessments may therefore be biased, particularly if the number of records are few, as it is the case for many tropical species.

While newly proposed methods for automated cleaning of species records hold great promise for biodiversity research, their use across taxonomic groups and datasets remains poorly explored. Here, we test the effect of automated filtering of species geographic occurrence records on the number of records available in different groups of animals, fungi, and plants. Furthermore, we test the impact of automated filtering procedures for the accuracy of preliminary automated conservation assessments compared to full IUCN assessments. Specifically, we evaluate a pipeline of 13 automated filters to flag possibly erroneous records by using 1) record meta-data and 2) geographic gazetteers. We address three questions:

1. Which filters lead to the biggest loss of data when applied?
2. Does the importance of individual filters differ among taxonomic groups?
3. Does automated filtering improve the accuracy of automated conservation assessments?

MATERIAL AND METHODS

Choice of study taxa

This study is the outcome of a workshop held at the Federal University of Rio Grande do Norte in Natal, Brazil in October 2018 gathering students and researchers working with different taxonomic groups of animals, fungi, and plants across the Neotropics. Each participant analysed geographic occurrence data from their taxonomic group of interest and commented on the results for their group. Hence, we include groups based on the expertise of the participants rather than following some arbitrary choice of taxa and taxonomic ranks. We acknowledge a varying degree in documented expertise and number of years working on each group. We obtained public occurrence records for 18 taxa, including one

96 plant family, nine plant genera, one genus of fungi, three families and one genus of terrestrial arthropods, one family of
97 snakes, one family of skates, and one genus of bony fish (Table 1).

98 **Species occurrence data**

99 We downloaded occurrence information for all study groups from www.gbif.org using the `rgbif` v1.4.0 package
100 (Chamberlain 2017) in R (GBIF.org, 2019a-p,2020a,b). We downloaded only records with geographic coordinates
101 and limited the study area to latitudes smaller than 33° N and longitudes smaller than 35° W and larger than 120° W
102 reflecting the Neotropics (Morrone 2014), our main area of expertise. The natural distribution of all included taxa are
103 confined to the Neotropics except Arhynchobatidae, Diogenidae, Dipsadidae, Entomobryidae, *Gaylussacia*, Iridaceae,
104 Neanuridae, and *Tillandsia* for which we only obtained the Neotropical occurrences.

105 **Automated cleaning**

106 We followed the cleaning pipeline outlined by Zizka, Silvestro, et al. (2019) and first filtered the raw data downloaded
107 from GBIF (“raw”, hereafter) using meta-data for those records for which they were available, removing: 1) records
108 with a coordinate precision below 100 km (as this represent the grain size of many macro-ecological analyses); 2) fossil
109 records and records of unknown source; 3) records collected before 1945 (before the end of the Second World War,
110 since coordinates of old records are often imprecise); and 4) records with an individual count of less than one and more
111 than 99. Furthermore, we rounded the geographic coordinates to four decimal precision and retained only one record per
112 species per location (i.e., test for duplicated records). In a second step, we used the `clean_coordinates` function
113 of the `CoordinateCleaner` v2.0-11 package (Zizka, Silvestro, et al. 2019) with default options to flag errors
114 that are common to biological data sets (“filtered”, hereafter). These include: coordinates in the sea for terrestrial taxa
115 and on land for marine taxa, coordinates containing only zeros, coordinates assigned to country and province centroids,
116 coordinates within urban areas, and coordinates assigned to biodiversity institutions. See Table 2 for a summary of all
117 flagging tests we used.

118 **Downstream analyses**

119 We first generated species richness maps using 100x100 km grid cells for the raw and filtered datasets respectively, using
120 the package `speciesgeocodeR` v2.0-10 (Töpel et al. 2016). We then performed an automated conservation
121 assessment for all study groups based on both datasets using the `ConR` v1.2.4 package (Dauby et al. 2017). `ConR`
122 estimates the EOO, AOO, and the number of locations, and then suggests a preliminary conservation status based on
123 Criterion B of the global IUCN Red List. While these assessments are preliminary (see IUCN Standards and Petitions
124 Subcommittee 2017), they can be a proxy used by the IUCN to speed up full assessments. We then benchmarked the
125 preliminary conservation assessments against the global IUCN Red List assessments as far as available for the same
126 taxa, which we obtained from www.iucn.org via the `rredlist` v.0.5.0 package (Chamberlain 2018).

127 **Evaluation of results**

128 Each author provided an informed comment on the performance of the raw and cleaned datasets, concerning the number
129 of flagged records and the accuracy of the overall species richness maps. We then compared the agreement between
130 automated conservation assessments based on raw and filtered occurrences with the global IUCN Red List for those taxa
131 where IUCN assessments were available (www.iucn.org).

132 We did all analyses in the R computing environment (R Core Team 2019), using standard libraries for data handling
133 and visualization (Wickham 2018; Garnier 2018; Ooms 2014, 2019; Hijmans 2019). All scripts are available from a
134 zenodo repository ([doi:10.5281/zenodo.3695102](https://doi.org/10.5281/zenodo.3695102)).

135 **RESULTS**

136 We retrieved a total of 218,899 species occurrence records, with a median of 2,844 records per study group and 10
137 records per species (Table 3, Appendix 1). We obtained most records for Dipsadidae (64,249) and fewest for *Thozetella*
138 (51). The species with most records was *Harengula jaguana* (19,878).

139 Our automated tests flagged a total of 97,018 records (Fig. 2), with a median of 45% per group. Overall, the most
140 important test was for duplicated records (on average 35.5% per taxonomic group). The filtering steps based on record
141 meta-data that flagged the largest number of records were the basis of records (5.9%) and the collection year (3.4%).
142 The most important automated tests were for urban area (8.6%) and the occurrence from records of terrestrial taxa in the
143 sea and marine taxa on land (4.3%, see Table 3 and Appendix 1 in the electronic supplement for further details and

144 the absolute numbers). Only a few records were flagged by the coordinate precision, zero coordinates and biodiversity
145 institution tests (Fig. 3).

146 Entomobryidae, Diogenidae, and Neanuridae had the highest fraction of flagged records (Table 3). In general, the
147 different filters we tested were of similar importance for different study groups. There were few outstanding exceptions,
148 including the particularly high proportions of records flagged on the “basis of record test” for *Tityus* (7.0%), Dipsadidae
149 (5.6%), *Prosthechea* (5.0%) and *Tillandsia* (4.9%), by the collection year for Dipsadidae (11.3%), by the taxonomic
150 identification level for *Tityus* (1.6%), by the capital coordinates for *Oocephalus* (6.1%) and *Gaylussacia* (3.2%), by the
151 seas/land test for Diogenidae and *Thozetella*, and by the urban areas test for *Oocephalus* (13.3%) and Iridaceae (12.3%).
152 Furthermore, Entomobryidae differed considerably from all other study taxa with exceptionally high numbers of records
153 flagged by the “basis of record,” level of identification’ and urban areas” tests.

154 Geographically, the records flagged by the “basis of record” and “individual count” tests were concentrated in
155 Central America and southern North America, and a relatively high number of records flagged due to their proximity to
156 the centroids of political entities were located on Caribbean islands (Fig. 3). See Appendix 2 for species richness maps
157 using the raw and cleaned data for all study groups.

158 We found IUCN assessments for 579 species that were also included in our distribution data from 11 of our study
159 groups (Table 4, Appendix 3). The fraction of species evaluated varied among the study group, with a maximum of
160 100% for *Harengula* and *Lepismium* and a minimum of 2.3% for Iridaceae (note that the number of total species varied
161 considerably among groups). The median percentage of species per study group with an IUCN assessment was 15%. A
162 total of 102 species were listed as Threatened by the IUCN global Red List (CR = 19, EN = 40, VU = 43) and 477 as
163 Not Threatened.

164 We obtained automated conservation assessments for 2,181 species in the filtered dataset. Based on the filtered data,
165 the automated conservation assessment evaluated 1,382 species as possibly threatened (63.4%, CR = 495, EN = 577, VU
166 = 310, see Appendix 3 for assessments of all species). The automated assessment based on the filtered dataset agreed
167 with the IUCN assessment for identifying species as possibly threatened (CR, EN, VU) for 358 species (64%; Table 4).
168 Filtering reduced the EOO by -18.4% and the AOO by -9.9% on median per group. For the raw dataset the agreement
169 with IUCN was higher at 381 species (65.7%).

170 DISCUSSION

171 Automated flagging based on meta-data and automatic tests flagged on average 45% of the records per taxonomic
172 group. The filters for basis of record, duplicates, collection year, and urban areas flagged the highest fraction of records
173 (Question 1). The importance of different tests was similar across taxonomic groups, with exceptions especially for
174 the tests on basis of record, collection year, capital coordinates, and urban areas (Question 2). The results for species
175 richness were similar between the raw and filtered data with some improvements by the filters. We found little impact of
176 filtering on the accuracy of the automated conservation assessments (Question 3).

177 The relevance of individual filters

178 The aim of automated filtering is to identify possibly problematic records that are unsuitable for particular downstream
179 analyses. Hence, our tests may identify true errors (e.g., wrong coordinates such as those with exclusively zeros, equal
180 latitude and longitude and records of terrestrial species in the sea or marine species on land) as well as records that are
181 not necessarily erroneous, but have a low precision (e.g., collection year, coordinate precision, centroids of political
182 units) or are redundant for certain analyses (duplicates). The importance of each test and the linked thresholds must
183 be judged based on the specific downstream analyses. For instance, records assigned to country centroids might be
184 acceptable for inter-continental comparisons, but can be considered erroneous for species distribution modelling on a
185 local scale. As our results show, adapting automated tests to the geographic study area or the taxonomic study group can
186 be warranted. For instance, the high number of records flagged for centroids on the Lesser Antilles (Fig. 3) might be
187 overly strict, although we chose a conservative distance for the Political centroid test (1 km).

188 Several reasons may explain the high proportion of records flagged as duplicates. First, the deposition of duplicates
189 from the same specimen hosted at different collections is common practice, especially for plants, where a specimen
190 duplication is easily feasible. Second, independent collections at similar localities may occur, in particular for local
191 endemics. Third, low coordinate precision, for instance based on automated geo-referencing from locality descriptions,
192 may lump records from nearby localities. Duplicated records do not represent errors *per se*, but may constitute an
193 unnecessary computational burden for large-scale biodiversity analyses (Antonelli et al. 2018), which is often a
194 bottleneck for analyses of datasets with millions of records.

Similarities and differences among taxa

The number of records flagged by individual tests was similar across study groups, suggesting that similar error types might be relevant for collections of plants and animals. Therefore, the same tests can be used across taxonomic groups. Some notable exceptions stress the need to adapt each test to the taxonomic study group to balance data quality and data availability. The high fraction of records flagged by the “basis of record test” for *Tityus*, Dipsadidae, *Prosthechea* and *Tillandsia*, were mostly caused by a high number of records in these groups based on unknown collection methods, which might be caused by the contribution of specimen datasets lacking this information for these groups. The high fraction of records flagged by the collection year test for Dispadidae was caused by a high collection effort already in the late 1880s and early 1900s, as can be expected for a charismatic group of reptiles, but also by 500 records dated to the year 1700. The latter records likely represent a data entry error: they are all contributed to GBIF from the same institution, and the institutions code for unavailable collection dates is 1700-01-01 - 2014-01-01, which has likely erroneously been converted to 1700. The high number of species flagged at capital coordinates and within urban areas for the plant groups Iridaceae and *Oocephalus* might be related to horticulture since at least some species in those groups are commonly cultivated as ornamentals.

The general agreement between the species richness maps based on raw and filtered data was encouraging, for the use of this data for large-scale biogeographic research (Fig. 4, Appendix 2). The filter based on political centroids had an important impact on species richness patterns, which is congruent with the results from a previous study in the coffee family (Maldonado et al. 2015). Records assigned to country or province centroids are often old records, which are geo-referenced at a later point based on vague locality descriptions. These records are at the same time more likely to represent dubious species names, since they might be old synonyms or type specimens of species that have only been collected and described once, which are erroneously increasing species numbers.

We overall consider the effect of the automated filters as positive since they identified the above-mentioned issues and increased the data precision and reduced computational burden (Table 3, Appendix 2). However, in some cases filters failed to remove major issues, often due to incomplete metadata. For instance, for Diogenidae we found at least two records of an extinct species (*Paguristes mexicanus*) which slipped the “basis of record” test because they were marked as “preserved specimen” rather than “fossil specimen”. Furthermore, for Entomobryidae we found that for 1,996 records the meta-data on taxonomic rank was “UNRANKED” despite all of them being identified to species level, leading to a high fraction of records removed by the “Identification level” filter. Additionally automated filters might be overly strict or unsuitable for certain taxa. For instance, in Entomobryidae, 2,004 samples were marked as material samples, and therefore removed by our global filter retaining only specimen and observation data, which in this case was overly strict.

It should be noted, that while in this study we focused on meta-data and geographic filtering, taxonomic cleaning—the resolution of synonymies and identification of accepted names—is another important part of data curation, but depends on taxon-specific taxonomic backbones and synonymy lists which are not readily available for many groups.

The impact of filtering on the accuracy of automated conservation assessments

The accuracy of the automated conservation assessment was in the same range as found by previous studies (Nic Lughadha et al. 2019; Zizka, Azevedo, et al. 2019). The similar accuracy of the raw and filtered dataset for the automated conservation assessment was surprising, in particular given the EOO and AOO reduction observed in the filtered dataset (Table 4) and the impact of errors on spatial analyses observed in previous studies (Gueta and Carmel 2016). The robustness of the automated assessment was likely due to the fact that the EOO for most species was large, even after the considerable reduction caused by filtering. This might be caused by the structure of our comparison, which only included species that were evaluated by the IUCN Red List (and not considered as Data Deficient) and at the same time had occurrences recorded in GBIF. Those inclusion criteria are likely to have biased the datasets towards species with large ranges, since generally more data for them are available. The robustness of automated conservation assessments to data quality is encouraging, although these methods are only an approximation (and not replacements) of full IUCN Red List assessments, especially for species with few collection records (Rivers et al. 2011).

CONCLUSIONS

Our results suggest that between one to two thirds of the occurrence records obtained from public databases might be unsuitable for downstream biodiversity analyses. While the majority of these records might not be erroneous *per se*, they might be overly imprecise and thereby increase uncertainty of downstream results or add computational burden on big data analyses.

246 Our results stress the importance of (meta-)data exploration prior to most biodiversity analyses. Automated flagging
247 can help to identify problematic records, but also highlight the necessity to customize tests and thresholds to the specific
248 taxonomic groups and geographic area of interest. The putative errors we encountered point to the importance to train
249 researchers and students to curate species occurrence datasets and to allow users to provide data aggregators such as
250 GBIF with feedback for particular records, thus contributing to overall accuracy and precision of public biodiversity
251 databases.

252 **ACKNOWLEDGEMENTS**

253 This study enrolled participants of the workshop “Biodiversity data: from field to yield” led by Alice Calvente, Fernanda
254 Carvalho, Alexander Zizka, and Alexandre Antonelli through the Programa de Pós Graduação em Sistemática e Evolução
255 of the Universidade Federal do Rio Grande do Norte (UFRN) and promoted by the 6th Conference on Comparative
256 Biology of Monocotyledons - Monocots VI. We thank the Pró-reitoria de Pesquisa and the Pró-reitoria de Pós-graduação
257 of UFRN for financial support (edital 02/2016 - internacionalização). AZ is funded by iDiv via the German Research
258 Foundation (DFG FZT 118), specifically through sDiv, the Synthesis Centre of iDiv. AA is supported by the Swedish
259 Research Council, the Knut and Alice Wallenberg Foundation, the Swedish Foundation for Strategic Research and the
260 Royal Botanic Gardens, Kew. FS was financed by the Coordenação de Aperfeiçoamento de Pessoal de Nível Superior -
261 Brasil (CAPES) - Finance Code 001 and Fundação de Amparo à Pesquisa do estado de São Paulo (FAPESP) (FAPESP,
262 process 2015/20215-7).

263 **SUPPLEMENTARY MATERIAL**

- 264 • Appendix 1 - Absolute number of flagged records per taxonomic group and test
- 265 • Appendix 2 - Taxon specific richness maps and comments
- 266 • Appendix 3 - Full results of the conservation assessment

267 **TABLES**

Table 1. The study groups and their taxonomy. This study includes three marine and 15 terrestrial taxa, seven of them animals, one group of fungi and ten plants, belonging to 16 different orders.

Taxon	Taxon rank	Realm	Common name	'Phylum'	Order	Family
Diogenidae	Family	Marine	Hermit crabs	Arthropoda	Decapoda	Diogenidae
Entomobryidae	Family	Terrestrial	Springtails	Arthropoda	Entomobryomorpha	Entomobryidae
Neanuridae	Family	Terrestrial	Springtails	Arthropoda	Poduromorpha	Neanuridae
<i>Tityus</i>	Genus	Terrestrial	Scorpions	Arthropoda	Scorpiones	Buthidae
Arhynchobatidae	Family	Marine	Skates	Chordata	Rajiformes	Arhynchobatidae
Dipsadidae	Family	Terrestrial	Snakes	Chordata	Squamata	Dipsadidae
<i>Harengula</i>	Genus	Marine	Herrings	Chordata	Clupeiformes	Clupeidae
<i>Thozetella</i>	Genus	Terrestrial	Sac fungi	Ascomycota	Chaetosphaeriales	Chaetosphaeriaceae
<i>Conchocarpus</i>	Genus	Terrestrial	NA	Angiosperms	Sapindales	Rutaceae
<i>Gaylussacia</i>	Genus	Terrestrial	Huckleberries	Angiosperms	Ericales	Ericaceae
<i>Harpalyce</i>	Genus	Terrestrial	NA	Angiosperms	Fabales	Fabaceae
Iridaceae	Family	Terrestrial	NA	Angiosperms	Asparagales	Iridaceae
<i>Lepismium</i>	Genus	Terrestrial	Cacti	Angiosperms	Caryophyllales	Cactaceae
<i>Ocephalus</i>	Genus	Terrestrial	NA	Angiosperms	Lamiales	Lamiaceae
<i>Pilosocereus</i>	Genus	Terrestrial	NA	Angiosperms	Caryophyllales	Cactaceae
<i>Prosthechea</i>	Genus	Terrestrial	Orchids	Angiosperms	Asparagales	Orchidaceae
<i>Tillandsia</i>	Genus	Terrestrial	Bromeliads	Angiosperms	Poales	Bromeliaceae
<i>Tocoyena</i>	Genus	Terrestrial	NA	Angiosperms	Gentianales	Rubiaceae

Table 2. The automated filters used in this study.

Test	Type	Rationale
Basis of record	Meta-data	Records might be unsuitable or unreliable for certain analyses dependent on their source, e.g. "fossil" or "unknown"
Collection year	Meta-data	Coordinates from old records are more likely to be imprecise or erroneous coordinates since they are derived from geo-referencing based on the locality description. This is more problematic for older records, since names or borders of places may change
Coordinate precision	Meta-data	Records may be unsuitable for a study if their precision is lower than the study analysis scale
Identification level	Meta-data	Records may be unsuitable if they are not identified to species level.
Individual count	Meta-data	Records may be unsuitable if the number of recorded individuals is 0 (record of absence) or if the count is too high, as this is often related to records from barcoding or indicative of data entry problems.
Biodiversity institutions	Gazetteer-based	Records may have coordinates at the location of biodiversity institutions, e.g. because they were erroneously entered with the physical location of the specimen or because they represent individuals from captivity or horticulture
Duplicates	Gazetteer-based	Duplicated records may add unnecessary computational burden, in particular for large scale biodiversity analyses and distribution modelling for many species
Equal lat/lon	Gazetteer-based	Coordinates with equal latitude and longitude are usually indicative of data entry errors
Capitals	Gazetteer-based	Records may be assigned to the coordinates of country capitals based on a vague locality description
Political centroids	Gazetteer-based	Records may be assigned to the coordinates of the centroids of political entities based on a vague locality description
Sea	Gazetteer-based	Coordinates from terrestrial organisms in the sea are usually indicative of data entry errors, e.g. swapped latitude and longitude
Urban areas	Gazetteer-based	Records from urban areas are not necessarily errors, but often represent imprecise records automatically geo-referenced from vague locality descriptions or old records from different land-use types
Zeros	Gazetteer-based	Coordinates with plain zeros are often indicative of data entry errors

Table 3. The impact of automated filtering on occurrence records for 18 Neotropical taxa downloaded from www.gbif.org. From column four onwards the numbers show the percentage of records flagged by the respective test. Only tests that flagged at least 0.1% of the records in any group are shown. Individual records can be flagged by multiple tests, therefore the sum of percentage from all tests can supersede the total percentage.

Taxon	Summary			Meta data filters					Gazetteer-based filters						
	Total records	Records flagged	Fraction flagged [%]	Basis of record [%]	Collection year [%]	Coordinate precision [%]	Id-level [%]	Individual count [%]	Biodiversity Institutions [%]	Capitals [%]	Duplicates [%]	Political centroids [%]	Sea/land area [%]	Urban areas [%]	Zeros [%]
Diogenidae	13,840	9,508	68.7	1.7	2.5	0.0	0.0	0.0	0.0	0.7	33.8	0.2	44.3	1.3	0.0
Entomobryidae	2,767	2,498	90.3	72.9	2.0	0.0	72.1	0.0	0.1	0.1	85.5	0.0	0.0	70.1	0.0
Neanuridae	689	461	66.9	2.9	1.3	0.0	0.0	0.0	0.0	0.0	62.4	0.0	0.0	2.0	0.0
<i>Tityus</i>	1,018	562	55.2	7.0	0.4	1.8	1.6	0.0	0.5	1.2	43.5	0.1	0.0	6.9	0.0
Arhynchobatidae	14,633	5,635	38.5	1.7	1.3	0.0	0.9	0.0	0.0	0.0	35.4	0.0	3.8	1.9	0.0
Dipsadidae	64,249	37,058	57.7	5.6	11.3	0.8	0.0	0.1	0.3	1.8	46.3	0.4	0.0	8.5	0.0
<i>Harengula</i>	36,697	11,384	31.0	1.0	0.4	0.0	0.3	0.0	0.0	0.2	27.0	0.1	5.5	0.2	0.0
<i>Thozetella</i>	51	18	35.3	0.0	0.0	0.0	0.0	0.0	0.0	0.0	27.5	0.0	23.5	2.0	0.0
<i>Conchocarpus</i>	1,551	670	43.2	0.5	1.9	0.1	0.0	0.0	0.1	0.0	39.6	0.9	0.4	2.3	0.0
<i>Gaylussacia</i>	3,998	1,887	47.2	0.7	4.4	0.6	0.0	0.0	0.1	3.2	41.8	1.1	0.1	5.2	0.0
<i>Harpalyce</i>	870	288	33.1	0.5	5.5	0.7	0.0	0.9	0.0	1.0	26.0	1.3	0.0	3.8	0.0
Iridaceae	23,127	7,774	33.6	0.9	4.7	0.1	0.0	1.3	0.4	1.0	17.1	0.4	0.1	12.3	0.0
<i>Lepismium</i>	825	245	29.7	0.0	2.1	0.0	0.0	0.0	0.0	0.1	21.9	0.1	0.0	7.8	0.0
<i>Oocephalus</i>	883	435	49.3	0.0	0.7	0.3	0.0	0.1	0.0	6.1	41.9	0.8	0.0	13.3	0.0
<i>Pilosocereus</i>	1,940	501	25.8	1.8	7.0	0.0	0.0	0.9	0.2	0.5	16.8	0.5	0.0	2.1	0.0
<i>Prosthechea</i>	6,617	2,088	31.6	5.0	8.3	0.1	0.0	0.2	0.0	0.4	19.6	1.7	0.0	0.9	0.1
<i>Tillandsia</i>	42,222	14,908	35.3	4.9	5.1	0.1	0.0	1.0	0.3	0.7	19.8	0.7	0.0	9.2	0.0
<i>Tocoyena</i>	2,922	1,098	37.6	0.1	1.9	0.2	0.0	0.5	0.0	0.8	32.3	0.8	0.2	5.0	0.0
Total	218,899	97,018	44.3	4.2	5.6	0.3	1.0	0.4	0.2	1.0	32.3	0.4	4.0	7.1	0.0

Table 4. Conservation assessment for 11 Neotropical taxa of plants and animals based on three datasets. IUCN: global red list assessment obtained from www.iucn.org; GBIF Raw: Preliminary conservation assessment based on IUCN Criterion B using ConR and the raw dataset from GBIF; GBIF filtered: Preliminary conservation assessment based on IUCN Criterion B using ConR and the filtered dataset. Only taxa with at least one species evaluated by IUCN shown.

Taxon	IUCN			GBIF Raw			GBIF Filtered				
	n taxa	Evaluated [%]	Threatened [%]	n taxa	Threatened [%]	Match with IUCN [%]	n taxa	Threatened [%]	Match with IUCN [%]	EOO change compared to raw [%]	AOO change compared to raw [%]
<i>Arhynchobatidae</i>	37	51.3	17.9	39	35.9	45.0	39	41.0	40.0	-32.7	-18.5
<i>Dipsadidae</i>	520	68.0	8.8	638	58.3	63.0	598	59.9	61.2	-2.3	-15.6
<i>Harengula</i>	4	100.0	0.0	4	0.0	100.0	4	0.0	100.0	-38.0	-36.9
<i>Conchocarpus</i>	4	8.7	0.0	46	63.0	100.0	45	62.2	100.0	-15.3	-7.1
<i>Gaylussacia</i>	2	3.3	0.0	61	59.0	50.0	58	60.3	50.0	-22.5	-8.6
<i>Harpalyce</i>	3	15.0	5.0	20	65.0	66.7	17	58.8	50.0	-18.4	-16.5
<i>Iridaceae</i>	13	2.3	0.2	531	64.4	50.0	466	62.9	62.5	-18.2	-12.3
<i>Lepismium</i>	6	100.0	0.0	6	16.7	83.3	6	16.7	83.3	-33.9	-7.9
<i>Pilosocereus</i>	41	80.9	19.1	47	55.3	73.7	46	56.5	71.1	-8.5	-5.8
<i>Tillandsia</i>	54	11.6	6.0	464	61.4	85.2	453	62.7	83.3	-13.7	-9.9
<i>Tocoyena</i>	3	13.6	4.5	22	31.8	66.7	21	38.1	66.7	-23.0	-9.5

268 **FIGURES**

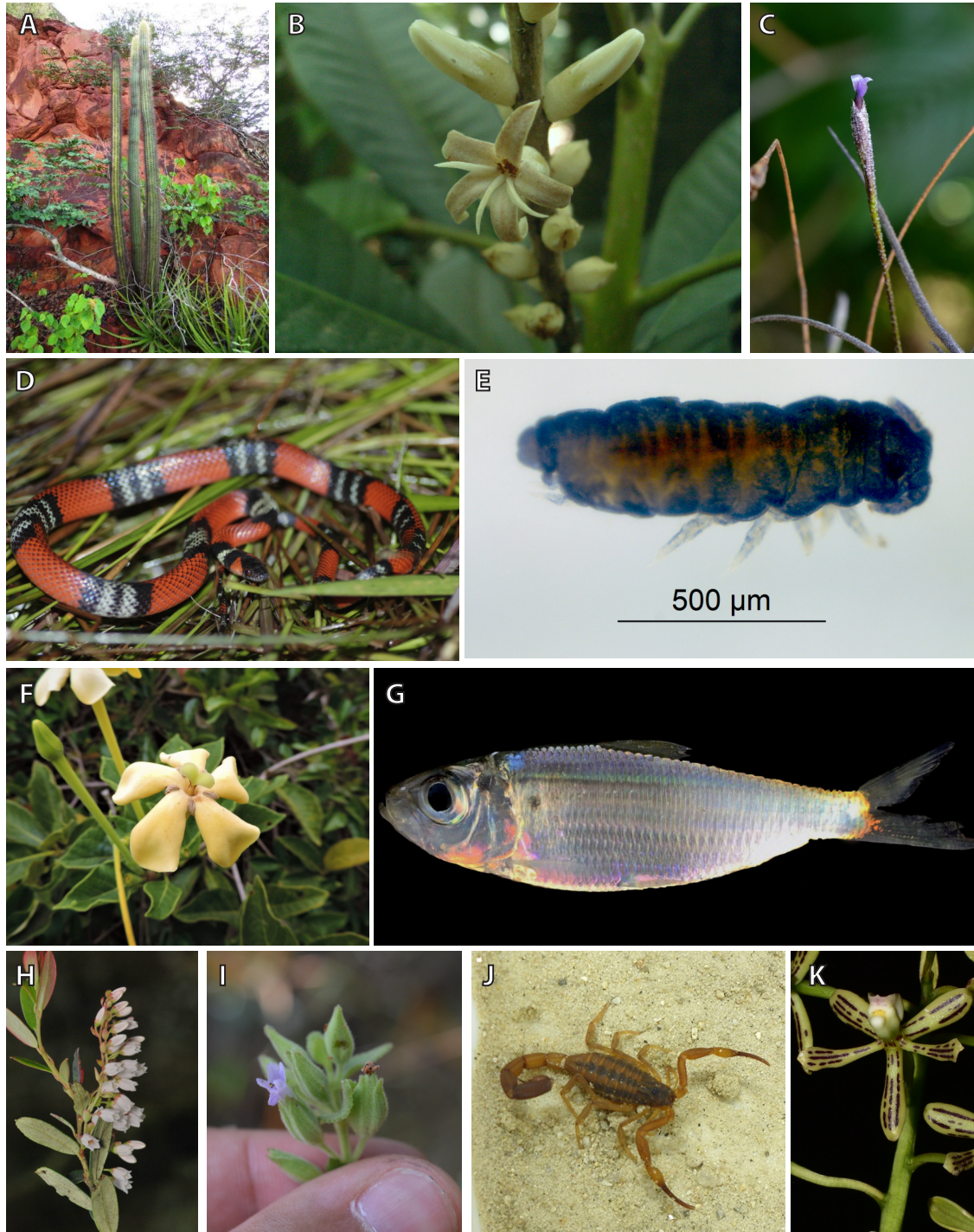


Figure 1. Examples for taxa included in this study. **A)** *Pilosocereus pusillibaccatus* (*Pilosocereus*), **B)** *Conchocarpus macrocarpus* (*Conchocarpus*); **C)** *Tillandsia recurva* (*Tillandsia*); **D)** *Oxyrhopus guibei* (*Dipsadidae*); **E)** *Aethiopella ricardoii* (*Neanuridae*); **F)** *Tocoyena formosa* (*Tocoyena*); **G)** *Harengula jaguana* (*Harengula*); **H)** *Gaylussacia decipiens* (*Gaylussacia*); **I)** *Oocephalus foliosus* (*Oocephalus*); **J)** *Tityus carvalhoi* (*Tityus*); **K)** *Prosthechea vespa* (*Prosthechea*), Image credits: A) Pamela Lavor, B) Juliana El-Ottra, C) Eduardo Tomaz, D) Filipe C Serrano, E) Raiane Vital da Paz, F) Fernanda GL Moreira, G) Thais Ferreira-Araujo, H) Luiz Menini Neto, I) Arthur Soares, J) Renata C Santos-Costa, K) Tiago Vieira

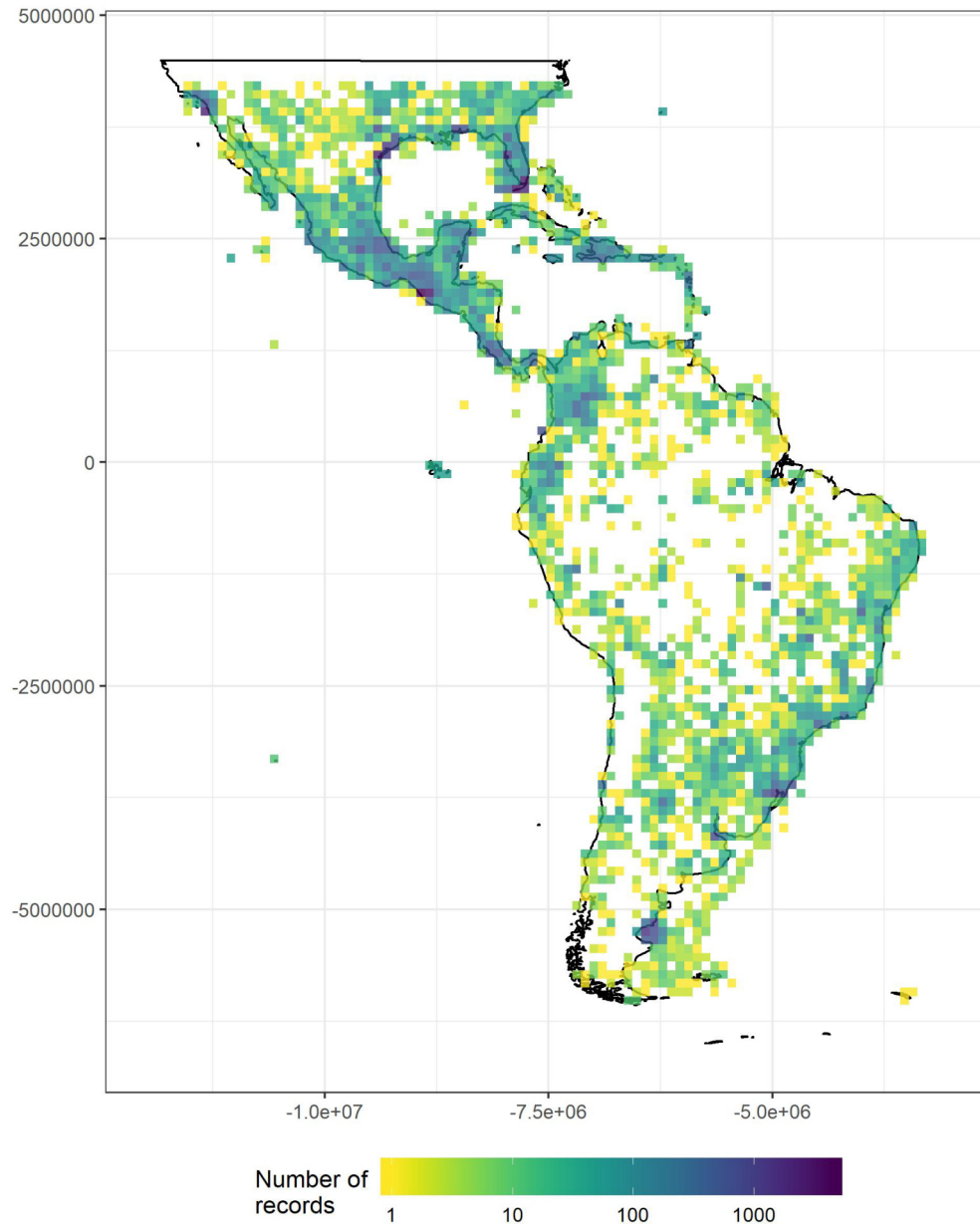


Figure 2. The difference in records available between the raw and filtered dataset of 18 Neotropical taxa including animals, fungi, and plants, plotted in a 100x100 km grid across the Neotropics.

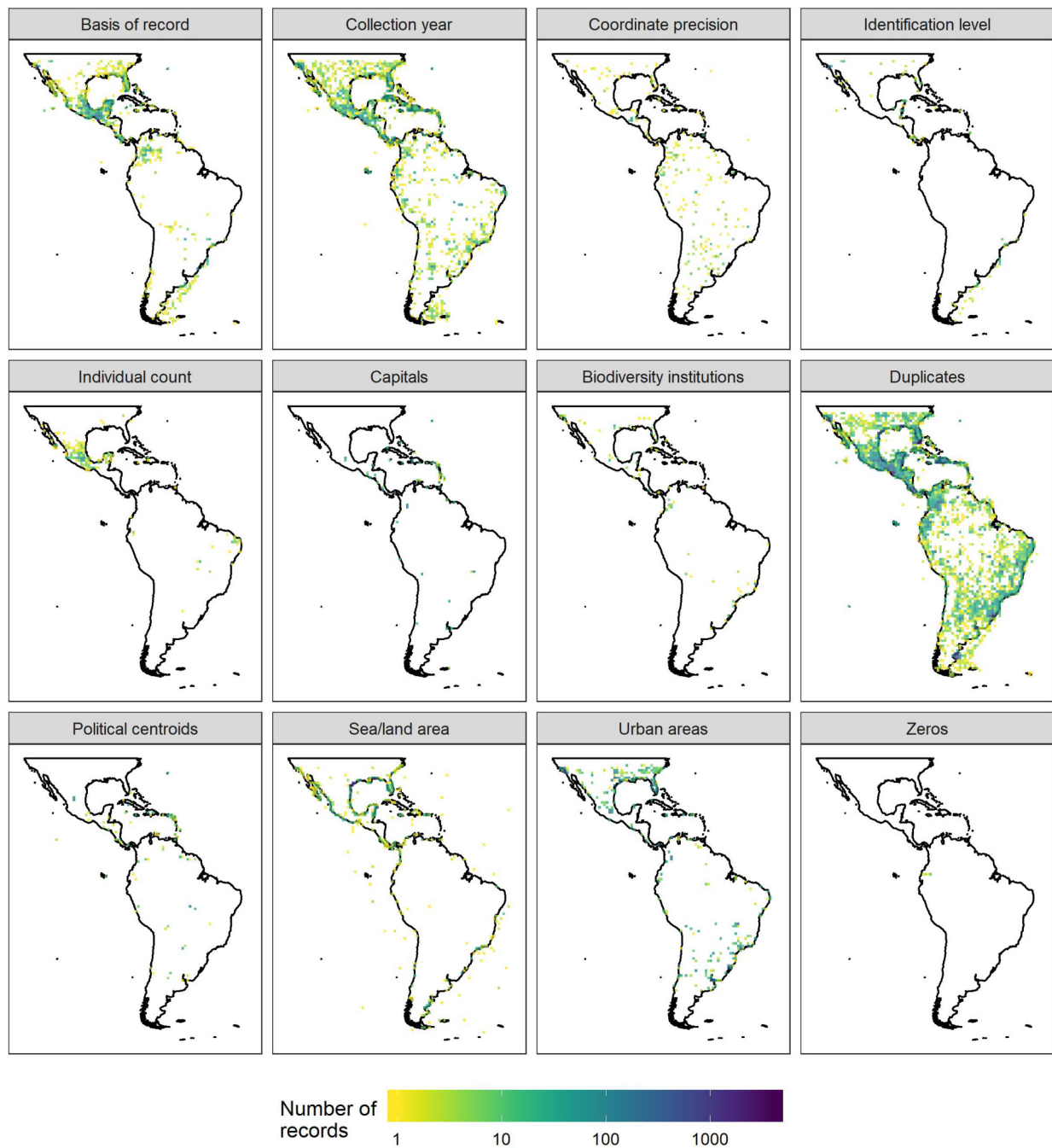


Figure 3. Geographic location of the occurrence records flagged by the automated tests applied in this study. Only filters that flagged at least 0.1% of records in any taxon are shown.

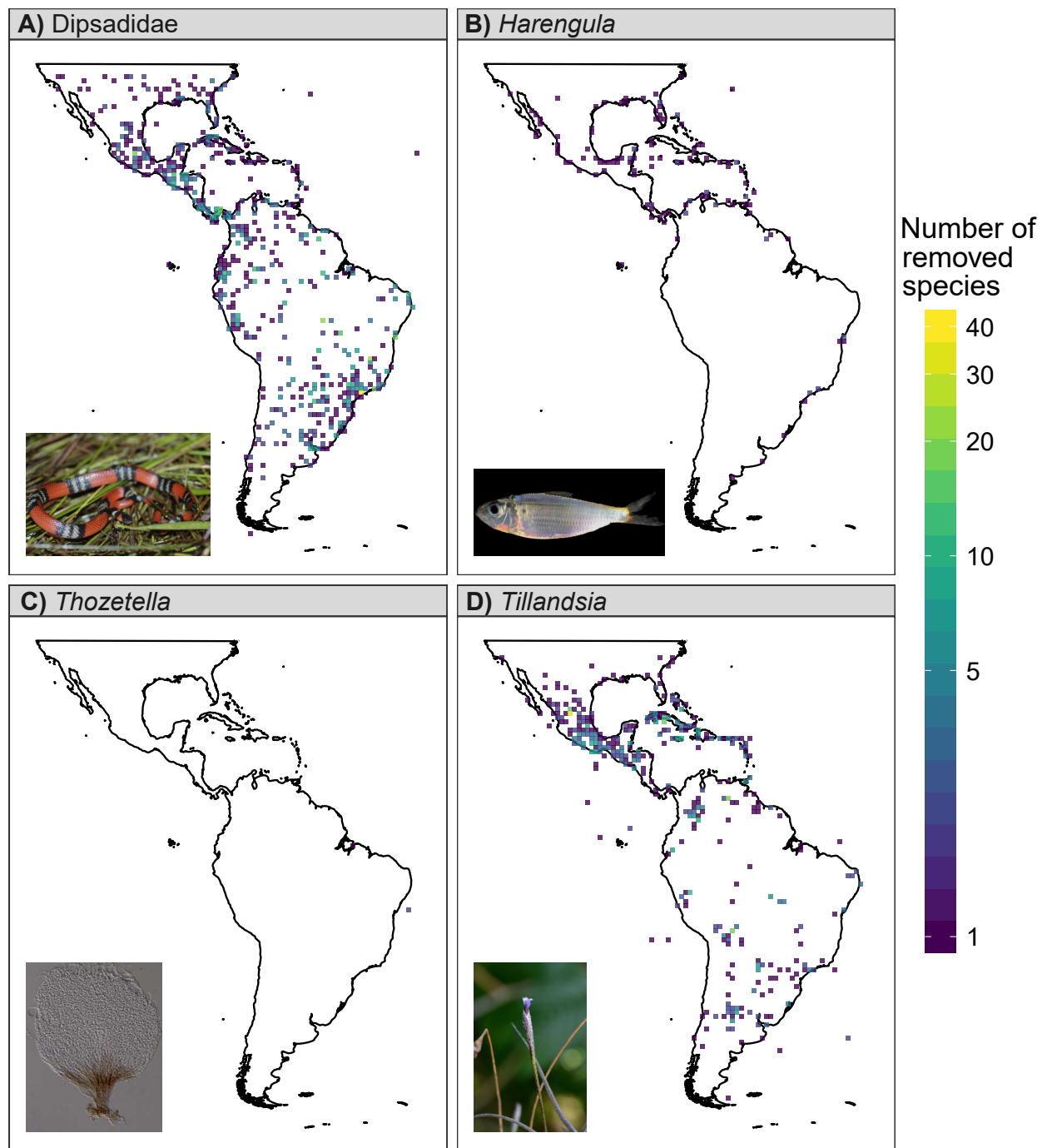


Figure 4. Illustrative examples for the difference in species richness between the raw and filtered dataset (raw - filtered) from four of the study taxa. Photo credits for C) by Tiago Andrade Borges Santos, otherwise as in Figure 1.

269 REFERENCES

- 270 Anderson, Robert P, Miguel Araújo, Antoine Guisan, Jorge M Lobo, Enrique Martínez-Meyer, Townsend Peterson,
271 and Jorge Soberón. 2016. “Final Report of the Task Group on GBIF Data Fitness for Use in Distribution Modelling -
272 Are species occurrence data in global online repositories fit for modeling species distributions? The case of the Global
273 Biodiversity Information Facility (GBIF).” Copenhagen, Denmark: GBIF.
- 274 Antonelli, Alexandre, Alexander Zizka, Fernanda Antunes Carvalho, Ruud Scharn, Christine D. Bacon, Daniele
275 Silvestro, and Fabien L Condamine. 2018. “Amazonia is the primary source of Neotropical biodiversity.” *Proceedings*
276 *of the National Academy of Sciences* 115 (23): 6034–9. <https://doi.org/10.1073/pnas.1713819115>.
- 277 Bachman, Steven P., Justin Moat, Andrew Hill, Javier de la Torre, and Ben Scott. 2011. “Supporting Red List
278 threat assessments with GeoCAT: geospatial conservation assessment tool.” *ZooKeys* 150 (November): 117–26.
279 <https://doi.org/10.3897/zookeys.150.2109>.
- 280 Chamberlain, Scott. 2016. “scrubr: Clean Biological Occurrence Records.” [https://cran.r-project.](https://cran.r-project.org/package=scrubr)
281 [org/package=scrubr](https://cran.r-project.org/package=scrubr).
- 282 ———. 2018. *rredlist: 'IUCN' Red List Client*. <https://cran.r-project.org/package=rredlist>.
- 283 Chamberlain, Scott A. 2017. “rgbif: Interface to the Global Biodiversity Information Facility API. R package version
284 0.9.9.” <https://github.com/ropensci/rgbif>.
- 285 Cosiaux, Ariane, Lauren M. Gardiner, Fred W. Stauffer, Steven P. Bachman, Bonaventure Sonké, William J.
286 Baker, and Thomas L. P. Couvreur. 2018. “Low extinction risk for an important plant resource: Conservation
287 assessments of continental African palms (Arecaceae/Palmae).” *Biological Conservation* 221 (May): 323–33. <https://doi.org/10.1016/j.biocon.2018.02.025>.
- 288 Dauby, Gilles, Tariq Stévant, Vincent Droissart, Ariane Cosiaux, Vincent Deblauwe, Murielle Simo-Droissart, Marc S.
289 M. Sosef, et al. 2017. “ConR: An R package to assist large-scale multispecies preliminary conservation assessments using
290 distribution data.” *Ecology and Evolution* 7 (24): 11292–11303. <https://doi.org/10.1002/ece3.3704>.
- 291 Garnier, Simon. 2018. *viridis: Default color maps from 'matplotlib'*. [https://cran.r-project.org/](https://cran.r-project.org/package=viridis)
292 [package=viridis](https://cran.r-project.org/package=viridis).
- 293 GBIF.org. 2019a. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.sojrfp)
294 [15468/dl.sojrfp](https://doi.org/doi.org/10.15468/dl.sojrfp).
- 295 ———. 2019b. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.bx0jjw)
296 [15468/dl.bx0jjw](https://doi.org/doi.org/10.15468/dl.bx0jjw).
- 297 ———. 2019c. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.zv6kuq)
298 [15468/dl.zv6kuq](https://doi.org/doi.org/10.15468/dl.zv6kuq).
- 299 ———. 2019d. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.uutyb6)
300 [15468/dl.uutyb6](https://doi.org/doi.org/10.15468/dl.uutyb6).
- 301 ———. 2019e. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.8hnzfo)
302 [15468/dl.8hnzfo](https://doi.org/doi.org/10.15468/dl.8hnzfo).
- 303 ———. 2019f. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.zznjbv)
304 [15468/dl.zznjbv](https://doi.org/doi.org/10.15468/dl.zznjbv).
- 305 ———. 2019g. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.zjjpmh)
306 [15468/dl.zjjpmh](https://doi.org/doi.org/10.15468/dl.zjjpmh).
- 307 ———. 2019h. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.4srw8a)
308 [15468/dl.4srw8a](https://doi.org/doi.org/10.15468/dl.4srw8a).
- 309 ———. 2019i. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.rpkjsh)
310 [15468/dl.rpkjsh](https://doi.org/doi.org/10.15468/dl.rpkjsh).
- 311 ———. 2019j. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.762543)
312 [15468/dl.762543](https://doi.org/doi.org/10.15468/dl.762543).
- 313 ———. 2019k. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.nmzgi9)
314 [15468/dl.nmzgi9](https://doi.org/doi.org/10.15468/dl.nmzgi9).
- 315 ———. 2019l. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.scmkx5)
316 [15468/dl.scmkx5](https://doi.org/doi.org/10.15468/dl.scmkx5).
- 317 ———. 2019m. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.d34gos)
318 [15468/dl.d34gos](https://doi.org/doi.org/10.15468/dl.d34gos).
- 319 ———. 2019n. “(29 December 2019) GBIF Occurrence Download.” [https://doi.org/doi.org/10.](https://doi.org/doi.org/10.15468/dl.6bzfz4)
320 [15468/dl.6bzfz4](https://doi.org/doi.org/10.15468/dl.6bzfz4).
- 321

- 322 ———. 2019o. “(29 December 2019) GBIF Occurrence Download.” <https://doi.org/doi.org/10.15468/dl.wkwque>.
- 323
- 324 ———. 2019p. “(29 December 2019) GBIF Occurrence Download.” <https://doi.org/doi.org/10.15468/dl.zj2cyj>.
- 325
- 326 ———. 2020a. “(25 February 2020) GBIF Occurrence Download.” <https://doi.org/doi.org/10.15468/dl.qazjh4>.
- 327
- 328 ———. 2020b. “(25 February 2020) GBIF Occurrence Download.” <https://doi.org/doi.org/10.15468/dl.ixq7wh>.
- 329
- 330 Gueta, Tomer, and Yohay Carmel. 2016. “Quantifying the value of user-level data cleaning for big data: A case study using mammal distribution models.” *Ecological Informatics* 34: 139–45. <https://doi.org/10.1016/j.ecoinf.2016.06.001>.
- 331
- 332
- 333 Hijmans, Robert J. 2019. “raster: Geographic data analysis and modeling.” <https://cran.r-project.org/package=raster>.
- 334
- 335 IUCN Standards and Petitions Subcommittee. 2017. “Guidelines for Using the IUCN Red List - Categories and Criteria. Version 13. Prepared by the Standards and Petitions Subcommittee. Downloadable from <http://www.iucnredlist.org/documents/RedListGuidelines.pdf>.”
- 336
- 337
- 338 Jin, Jing, and Jun Yang. 2020. “BDcleaner: A workflow for cleaning taxonomic and geographic errors in occurrence data archived in biodiversity databases.” *Global Ecology and Conservation* 21 (March): e00852. <https://doi.org/10.1016/j.gecco.2019.e00852>.
- 339
- 340
- 341 Maldonado, Carla, Carlos I. Molina, Alexander Zizka, Claes Persson, Charlotte M. Taylor, Joaquina Albán, Eder Chilquillo, Nina Rønsted, and Alexandre Antonelli. 2015. “Estimating species diversity and distribution in the era of Big Data: to what extent can we trust public databases?” *Global Ecology and Biogeography* 24 (8): 973–84. <https://doi.org/10.1111/geb.12326>.
- 342
- 343
- 344
- 345 Morrone, Juan J. 2014. “Biogeographical regionalisation of the Neotropical region.” *Zootaxa* 3782 (1): 1. <https://doi.org/10.11646/zootaxa.3782.1.1>.
- 346
- 347 Nic Lughadha, Eimear, Barnaby E. Walker, Cátia Canteiro, Helen Chadburn, Aaron P. Davis, Serene Hargreaves, Eve J. Lucas, et al. 2019. “The use and misuse of herbarium specimens in evaluating plant extinction risks.” *Philosophical Transactions of the Royal Society B: Biological Sciences* 374 (1763): 20170402. <https://doi.org/10.1098/rstb.2017.0402>.
- 348
- 349
- 350
- 351 Ooms, Jeroen. 2014. “The jsonlite Package: A Practical and Consistent Mapping Between JSON Data and R Objects.” *arXiv*. <https://arxiv.org/abs/1403.2805>.
- 352
- 353 ———. 2019. *writexl: Export Data Frames to Excel 'xlsx' Format*. <https://cran.r-project.org/package=writexl>.
- 354
- 355 Pelletier, Tara A., Bryan C. Carstens, David C. Tank, Jack Sullivan, and Anahí Espíndola. 2018. “Predicting plant conservation priorities on a global scale.” *Proceedings of the National Academy of Sciences* 115 (51): 13027–32. <https://doi.org/10.1073/pnas.1804098115>.
- 356
- 357
- 358 R Core Team. 2019. “R: A language and environment for statistical computing.” Austria, Vienna: R Foundation for Statistical Computing. <https://www.r-project.org/>.
- 359
- 360 Rivers, Malin C., Lin Taylor, Neil a. Brummitt, Thomas R. Meagher, David L. Roberts, and Eimear Nic Lughadha. 2011. “How many herbarium specimens are needed to detect threatened species?” *Biological Conservation* 144 (10): 2541–7. <https://doi.org/10.1016/j.biocon.2011.07.014>.
- 361
- 362
- 363 Robertson, Mark P, Vernon Visser, and Cang Hui. 2016. “Biogeo: an R package for assessing and improving data quality of occurrence record datasets.” *Ecography* 39: 394–401. <https://doi.org/10.1111/ecog.02118>.
- 364
- 365 Schmidt, Marco, Alexander Zizka, Salifou Traoré, Mandingo Ataholo, Cyrille Chatelain, Philippe Daget, Stefan Dressler, et al. 2017. “Diversity, distribution and preliminary conservation status of the flora of Burkina Faso.” *Phytotaxa Monographs* 304 (1): 1–215.
- 366
- 367
- 368 Stévant, T., G. Dauby, P. P. Lowry, A. Blach-Overgaard, V. Droissart, D. J. Harris, B. A. Mackinder, et al. 2019. “A third of the tropical African flora is potentially threatened with extinction.” *Science Advances* 5 (11): eaax9444. <https://doi.org/10.1126/sciadv.aax9444>.
- 369
- 370
- 371 Töpel, Mats, Alexander Zizka, Maria Fernanda Calió, Ruud Scharn, Daniele Silvestro, and Alexandre Antonelli. 2016. “SpeciesGeoCoder: Fast Categorization of Species Occurrences for Analyses of Biodiversity, Biogeography, Ecology, and Evolution.” *Systematic Biology* 66 (2): syw064. <https://doi.org/10.1093/sysbio/syw064>.
- 372
- 373
- 374 Wickham, Hadley. 2018. “tidyverse: Easily install and load the 'Tidyverse'.” <https://cran.r-project.org/package=tidyverse>.
- 375

376 Yesson, Chris, Peter W Brewer, Tim Sutton, Neil Caithness, Jaspreet S Pahwa, Mikhaila Burgess, W Alec Gray, et
377 al. 2007. “How Global Is the Global Biodiversity Information Facility?” Edited by James Beach. *PLoS ONE* 2 (11):
378 e1124. <https://doi.org/10.1371/journal.pone.0001124>.

379 Zizka, Alexander, Josue Azevedo, Elton Leme, Beatriz Neves, Andrea Ferreira, Daniel Caceres, and Georg
380 Zizka. 2019. “Biogeography and conservation status of the pineapple family (Bromeliaceae),” no. June: 1–13.
381 <https://doi.org/10.1111/ddi.13004>.

382 Zizka, Alexander, Daniele Silvestro, Tobias Andermann, Josué Azevedo, Camila Duarte Ritter, Daniel Edler, Harith
383 Farooq, et al. 2019. “CoordinateCleaner: Standardized cleaning of occurrence records from biological collection
384 databases.” Edited by Tiago Quental. *Methods in Ecology and Evolution* 10 (5): 744–51. [https://doi.org/10.](https://doi.org/10.1111/2041-210X.13152)
385 [1111/2041-210X.13152](https://doi.org/10.1111/2041-210X.13152).