

## Large-bodied birds are over-represented in unstructured citizen science data

Corey T. Callaghan<sup>1,2,3,\*</sup>, Alistair G. B. Poore<sup>2</sup>, Max Hofmann<sup>1,3</sup>, Christopher Roberts<sup>2</sup>, Henrique M. Pereira<sup>1,3</sup>

<sup>1</sup>German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Puschstr. 4, 04103 Leipzig, Germany

<sup>2</sup>Ecology & Evolution Research Centre; School of Biological, Earth and Environmental Sciences; UNSW Sydney; Sydney, NSW

<sup>3</sup>Institute of Biology, Martin Luther University Halle-Wittenberg, Am Kirchtor 1, 06108 Halle (Saale), Germany

\*Corresponding author:

email: [corey.callaghan@idiv.de](mailto:corey.callaghan@idiv.de)

**Note: This is a pre-print and has not undergone full peer-review.**

## 1 ABSTRACT

2 Citizen science platforms are quickly accumulating hundreds of millions of biodiversity  
3 observations around the world annually. Quantifying and correcting for the biases in citizen  
4 science datasets remains an important first step before these data are used to address ecological  
5 questions and monitor biodiversity. One source of potential bias among datasets is the difference  
6 between those citizen science programs that have unstructured protocols and those that have  
7 semi-structured or structured protocols for submitting observations. To quantify biases in an  
8 unstructured citizen science platform, we contrasted bird observations from the iNaturalist  
9 platform with that from a semi-structured citizen science platform — eBird — for the continental  
10 United States. We tested whether four traits of species (color, flock size, body size, and  
11 commonness) predicted if a species was under- or over-represented in the unstructured dataset  
12 compared with the semi-structured dataset. We found strong evidence that large-bodied birds  
13 were over-represented in the unstructured citizen science dataset; moderate evidence that  
14 common species were over-represented in the unstructured dataset; moderate evidence that  
15 species in large groups were over-represented; and no evidence that colorful species were over-  
16 represented in unstructured citizen science data. Our results suggest that biases exist in  
17 unstructured citizen science data when compared with semi-structured data, likely as a result of  
18 the detectability of a species and the inherent recording process. Importantly, in programs like  
19 iNaturalist the detectability process is two-fold — first, an individual needs to be detected, and  
20 second, it needs to be photographed, which is likely easier for many large-bodied species. Our  
21 results indicate that caution is warranted when using unstructured citizen science data in  
22 ecological modelling, and highlight body size as a fundamental trait that can be used as a  
23 covariate for modelling opportunistic species occurrence records, representing the detectability

24 or identifiability in unstructured citizen science datasets. Future research in this space should  
25 continue to focus on quantifying and documenting biases in citizen science data, and expand our  
26 research by including structured citizen science data to understand how biases differ among  
27 unstructured, semi-structured, and structured citizen science platforms.

28

29 *Keywords:* citizen science; biases; opportunistic data; presence-only data, species occurrence  
30 data, eBird; iNaturalist; species traits; detectability

## 31 INTRODUCTION

32 Citizen science, or community science, — the involvement of volunteers in scientific endeavors  
33 — is increasingly seen as a cost-effective method for biodiversity monitoring and research.  
34 Accordingly, the quantity and diversity of citizen science projects in the ecological and  
35 environmental sciences is increasing <sup>1</sup>. Such projects are quickly accumulating hundreds of  
36 millions of biodiversity observations around the world annually <sup>2,3</sup> expanding the spatial and  
37 temporal scope of our understanding in ecology, conservation, and natural resource management  
38 <sup>4,5</sup>. Citizen science projects vary widely in their scope, design, and intent <sup>6,7,8</sup>. Projects can range  
39 from unstructured (e.g., little training needed to participate and contribute  
40 opportunistic/incidental observations) to semi-structured (e.g., with minimal workflows and  
41 guidelines, but additional data collected with each observation can be included) to structured  
42 (e.g., prescribed sampling in space and time by mostly trained and experienced volunteers). The  
43 level of structure consequently influences the overall data quality of a particular project <sup>9,10</sup>.  
44  
45 Data quality from citizen science projects has been questioned <sup>11,12</sup>, and such concerns can act as  
46 a barrier to the widespread use of citizen science data in ecology and conservation <sup>13</sup>. These  
47 concerns arise because citizen science data can be biased temporally, spatially, and/or  
48 taxonomically. Temporally, many citizen science datasets are biased because participants  
49 frequently sample on weekends <sup>14</sup> or disproportionately during specific times of the year such as  
50 spring migration for birds <sup>15</sup>, or during specific times of day, such as the morning period when  
51 birds are most active. Spatially, there is often a disproportionate number of sightings from areas  
52 with large human populations <sup>16</sup>, areas with more accessibility <sup>17</sup>, regions with high biodiversity  
53 that attract observers <sup>18</sup>, and regions of the world with higher socioeconomic status <sup>19</sup>.

54 Taxonomic biases also exist as some taxa receive orders of magnitude more citizen science  
55 observations than other taxa, evidenced by the fact that birds represent a disproportionate amount  
56 of data in the Global Biodiversity Information Facility <sup>2</sup>. Even within citizen science projects  
57 focused on specific taxa, there can be strong taxonomic biases towards particularly charismatic  
58 species or those that are readily identified <sup>20, 21, 22, 23</sup>.

59

60 Despite potential biases in citizen science datasets, contrasts of data from volunteer participants  
61 to those contributed by more structured datasets have shown that citizen science programs can  
62 provide reliable data <sup>12, 24</sup>. For example, one case study found that mark-recapture models of  
63 whale sharks are reliable whether using sightings reported by the public or by experienced  
64 researchers <sup>25</sup>, and another case study found that unstructured data performs comparably with  
65 structured data in identifying and monitoring invasive plant species <sup>26</sup>. When analyzed  
66 appropriately, citizen science data can further our understanding of many facets of biodiversity,  
67 including estimating species distributions <sup>27, 28, 29</sup>, managing habitat for conservation <sup>30</sup>,  
68 estimating species richness <sup>31</sup>, monitoring pollination services <sup>32</sup>, and quantifying population  
69 trends <sup>33, 34</sup>. In the above examples, highlighting the potential uses of citizen science data,  
70 statistical solutions to account for known biases and noise in citizen science data are used <sup>3, 35, 36</sup>.

71

72 In addition to being an excellent resource for scientists to better understand ecological questions,  
73 citizen science projects are beneficial for society by encouraging increased engagement of the  
74 general public with science <sup>37, 38</sup>. Many citizen science projects provide learning opportunities for  
75 their volunteers. For example, participants in citizen science projects have increased their  
76 knowledge about invasive weeds <sup>39, 40, 41</sup>, increased their knowledge of bird biology and behavior

77 <sup>42</sup>, and even enhanced their conservation awareness and sense of place <sup>42,43</sup>. The ecological  
78 advances derived from citizen science data, combined with the important role it plays in  
79 community engagement with science, suggests that citizen science data will continue to play an  
80 important role in ecological and conservation research in the future <sup>2,4,38,44</sup>. However, what  
81 motivates volunteers to participate in science, and contribute observations, has important  
82 implications for the quality of the data obtained <sup>45</sup>, particularly if there are biases towards certain  
83 species, places, or times of sampling.

84

85 To ensure the continued and expanded use of citizen science data in ecology and conservation, it  
86 is important to document and understand the different biases present in citizen science datasets.  
87 Importantly, the degree of bias in a particular dataset will be influenced by the level of structure  
88 of that citizen science project. For example, unstructured projects (e.g., iNaturalist,  
89 [www.inaturalist.org](http://www.inaturalist.org)) or semi-structured projects (e.g., eBird, [www.ebird.org](http://www.ebird.org)) will generally be  
90 more spatially biased than structured projects that have pre-defined spatial sampling locations  
91 (e.g., Breeding Bird Surveys). Or, a citizen science project that collects incidental presence-only  
92 data, such as iNaturalist, is likely more susceptible to individual observer preferences compared  
93 with a semi-structured or structured project that requires all species encountered to be recorded  
94 by the observers. Charismatic species <sup>21</sup> can be over-represented in citizen science data because  
95 observers are more likely to record species that they, or society, consider more interesting or  
96 relevant <sup>46</sup>. Similarly, rare species are more likely to be the subject of conservation monitoring or  
97 more likely to be actively searched for by amateur naturalists <sup>47,48</sup> and thus can be over-  
98 represented in biodiversity datasets. In contrast, in some citizen science projects, abundant  
99 species can form a disproportionate number of records (e.g., <sup>49</sup>) because species' abundance and

100 ease of identification can lead to an increase in the number of records by casual observers <sup>50</sup>.

101 Inherently linked with observer preferences are issues of differences in species detectability <sup>50</sup>,

102 and the ease of making the observations. Therefore, species traits (e.g., body size, color, flock

103 size) may have an additive effect, influencing both the detectability of a species <sup>51, 52, 53</sup>, and in

104 turn, the likelihood of a species being submitted to an unstructured citizen science database.

105

106 Quantifying biases in citizen science datasets can help (1) researchers using these data to better

107 account for biases when drawing ecological conclusions, (2) the design and implementation of

108 future citizen science projects, and (3) understand what species or regions may need data

109 collection from professional scientists by understanding the ‘limits’ of citizen science projects <sup>19</sup>.

110 Here, we quantify biases in bird observation data from an unstructured, citizen science project —

111 iNaturalist — with that from a semi-structured one — eBird. We restricted our comparison to

112 birds because (1) birds are among the most popular taxa with the non-scientific public, ensuring

113 large sample sizes in both citizen science projects, and (2) data on the species traits that may

114 influence the likelihood of unstructured observations are readily available for birds. We assessed

115 the over-representation or under-representation of bird species’ observations in the unstructured

116 citizen science project compared to the semi-structured project (see Figure 1). We then tested the

117 following predictions: that (1) more colorful species; (2) larger species; (3) species with the

118 ‘least concern’ IUCN status; and (4) more gregarious species (i.e., with larger flock sizes) are

119 over-represented in the unstructured citizen science dataset (iNaturalist) in contrast to the semi-

120 structured citizen science dataset (eBird). Our analysis highlights the importance of considering

121 species’ traits when using citizen science data in ecological research.

122

## 123 METHODS

124 We made comparisons between iNaturalist ([www.inaturalist.org](http://www.inaturalist.org)) — an unstructured citizen  
125 science project — and eBird ([www.ebird.org](http://www.ebird.org)) — a semi-structured citizen science project<sup>15, 54</sup>.

126

127 *iNaturalist citizen science data.* iNaturalist is a multi-taxon citizen science project hosted by the  
128 California Academy of Sciences. It is an unstructured citizen science project where volunteers  
129 contribute opportunistic photos or sound recordings through a smart-phone or web-portal. Photos  
130 are then identified to the lowest possible taxonomic resolution using a community identification  
131 process, and once two users, or more than two-thirds, confirm the species-level identification of  
132 an organism it is considered “research grade”. Observations that are research grade are then  
133 uploaded to the Global Biodiversity Information Facility. We downloaded iNaturalist  
134 observations from the Global Biodiversity Information Facility for the contiguous United States  
135<sup>55</sup> for the period from January 2010 to May 2019, on December 3<sup>rd</sup>, 2019. For more details on the  
136 iNaturalist methodology, see here: <https://www.inaturalist.org/pages/getting+started>.

137

138 *eBird citizen science data.* eBird is one of the most successful citizen science projects in the  
139 world, with > 1 billion bird observations globally. It was launched in 2002 by the Cornell Lab of  
140 Ornithology and focuses on collecting reliable data on the distributions and relative abundance of  
141 birds throughout the world<sup>54</sup>. It is a semi-structured project where volunteers submit ‘checklists’  
142 of birds seen and/or heard on birding outings, following different protocols (e.g., stationary,  
143 incidental, or travelling). An important component of eBird that differentiates it from  
144 unstructured data collection is that users are required to indicate whether the checklist is  
145 ‘complete’ – meaning they included all species they were able to identify during that birding



146 outing. When using complete checklists only in an analysis, a user can infer non-detections in the  
147 dataset for any species not recorded. Observers can submit checklists at any time and place of  
148 their choosing, and for any duration and distance travelled. Most non-incidental checklists  
149 additionally include the duration and distance travelled while birding. Filters are set — based on  
150 spatiotemporal coordinates — which restrict the species and their associated counts that can be  
151 submitted without approval from a regional expert reviewer<sup>56</sup>. We used the eBird basic dataset  
152 (version ebd\_May-2019) and restricted our analysis to data from the contiguous United States for  
153 the period from January 2010 to May 2019. We also restricted the data used to those of the best  
154 ‘quality’ by excluding incomplete checklists, checklists that were incidental or historical, which  
155 travelled >5 km, lasted <5 min, and lasted >240 min, minimizing the leverage of outliers on  
156 analyses<sup>57, 58</sup>.

157

158 *Filtering and aggregating the citizen science datasets.* Although both datasets are global in  
159 scope, we restricted our analysis to the contiguous United States as both of these citizen science  
160 projects initiated in the United States, and thus the data are most numerous from there. For  
161 comparisons, we aggregated data at the state level. This was done to account for differences that  
162 may exist throughout the entirety of the United States including differences in user behavior and  
163 the species pools that differ geographically. We used the eBird Clements taxonomy (version  
164 2018) and all species from iNaturalist were matched with this taxonomy. A total of 1,030 species  
165 was initially collated from the eBird checklists, but many of these only occurred once or a few  
166 times — possibly representing misidentifications that had not yet been fixed by local reviewers  
167 or escaped and exotic birds which are incorporated in the eBird dataset but not considered part of  
168 the local avifauna or of interest to our analysis here. Although, these could represent scarce and

169 uncommon species in a state as well, albeit these are rarely sampled by iNaturalist. To account  
170 for these biases, we removed species that were on <1% of eBird checklists for a given state;  
171 trimming the eBird observations to the ‘core’ suite of species that occur in a state (*sensu*<sup>57</sup>).  
172 After trimming the species and harmonizing the taxonomy with iNaturalist, there were 507  
173 species remaining which were considered in our main analyses presented throughout the results.  
174 Although our results here are presented using the 1% cutoff level, we tested the sensitivity of this  
175 cutoff level and found comparable results across 0, 0.5, 1, and 1.5% cutoffs. For each state, the  
176 eBird and iNaturalist data were summarized by calculating the total number of observations in  
177 that state for every species. Using these aggregated data, we conducted preliminary comparisons  
178 of the unstructured and semi-structured datasets by quantifying the relationship between the  
179 number of eBird and iNaturalist observations at the state level, and at the species level.

180

#### 181 *Species-specific over- or under-representation in iNaturalist*

182 Our first analytical step was to model the log-log linear relationship between the total number of  
183 observations in iNaturalist and total number of observations in eBird for a species (Figure 1).  
184 This linear model was repeated separately for each state, where the response variable was log-  
185 transformed number of iNaturalist observations and the predictor variable was log-transformed  
186 number of eBird observations. Each model fitted was stratified by state, to account for inherent  
187 differences among states that were not of interest in our particular analysis, such as (1) the  
188 number of observers in a state, (2) the different relative abundance of a species throughout the  
189 United States, and (3) any other intrinsic differences that might exist among states that was not  
190 of interest in our analysis. A species with a high (i.e., positive) residual would be over-  
191 represented in iNaturalist relative to eBird, whereas a species with a low (i.e., negative) residual

192 would be under-represented in iNaturalist relative to eBird (Figure 1). Then we took the residuals  
193 from these models and used these as the response variables in our subsequent analyses of species  
194 characteristics (see below).

195

#### 196 *Species-specific trait data*

197 We tested whether four predictor variables (see Figure 1) would explain the over- or under-  
198 representation of bird species in the unstructured citizen science data. For each species, we used  
199 a proxy for their commonness/abundance, categorized according to IUCN status, taken from  
200 HBW BirdLife international checklist version 3 (<http://datazone.birdlife.org/species/taxonomy>).  
201 This variable was treated as an ordinal variable in our models (see below) and encompassed  
202 Least Concern, Vulnerable, and Near Threatened species. The three species recorded as  
203 endangered were removed from this analysis due to a lack of power at this level with so few  
204 observations. For each species we used the continuous predictor variables of (1) body size; (2)  
205 color; and (3) average flock size. Body sizes (adult body mass in grams) were taken from the  
206 amniote life history database compiled by Myhrvold et al. <sup>59</sup> and were log-transformed to meet  
207 normality assumptions. Color was taken from Dale et al. 2015 <sup>60</sup> and was extracted as RGB  
208 values for six patches per species (upper breast, lower breast, crown, forehead, nape, throat). To  
209 define a continuum of color where the brightest/most colorful (and likely most detectable species  
210 based on color) had the highest value we combined both the ‘distance from brown’ and the  
211 ‘brightness’ of a species for the data from Dale et al. 2015 <sup>60</sup>. Distance from brown was defined  
212 as the maximum Euclidian distance in the cubic RGB color space from brown (R = 102, B = 68,  
213 G = 0) from any of the six patches on a species, regardless of sex. Brightness was defined as the  
214 maximum relative luminance (i.e.,  $0.2126R + 0.7152G + 0.0722B$ ) from any of the six patches

215 on a species, regardless of sex. These two variables were combined and scaled from 0 to 1 for all  
216 species in Dale et al. 2015<sup>60</sup> and this value was used as our measure of color. Calculations were  
217 done in “Lab” space, an approximately perceptually uniform color space standardized by the  
218 Commission Internationale d'Eclairage. Exploratory analyses showed similar results with HSV  
219 color space. Flock size — an approximation of the gregariousness of a species — was taken from  
220 eBird as the average number of reported individuals among all checklists where a species was  
221 reported, across all data. We acknowledge that the number of a species reported on an eBird  
222 checklist likely encompasses both the gregariousness of a species as well as the density of a  
223 species in an area, as birders can travel through multiple territories.

224

#### 225 *Statistical analysis*

226 We used mixed effects models to examine the effects of species traits on the relative bias  
227 between our unstructured and semi-structured citizen science datasets. The response variable was  
228 the residuals from a log-log linear model fit between the eBird observations and the iNaturalist  
229 observations for a given species (described above), the predictor variables were the respective  
230 traits, and the random effect was state. First, we ran a global model where all traits were included  
231 as predictor variables: log<sub>10</sub>-transformed body size, log<sub>10</sub>-transformed flock size, color, and  
232 IUCN status treated as an ordinal variable. Second, to confirm the results of this global model,  
233 we ran four separate models – one for each trait as listed above – because there was much  
234 missing data for species’ traits. This approach allowed us to test the relationship of a predictor  
235 given the other predictor variables (i.e., all predictors against the response variable  
236 simultaneously) as well as the independent relationships (i.e., each predictor separately against  
237 the response variable).

238

239 *Data analysis and availability*240 All analyses were carried out in R software<sup>61</sup> and relied heavily on the tidyverse workflow<sup>62</sup>.241 Mixed-effects models were fitted using the lme4 package<sup>63</sup> and p-values were extracted using242 the lmerTest package<sup>64</sup>. Data and code to reproduce these analyses are available in a GitHub243 repository ([https://github.com/coreytcallaghan/inaturalist\\_preferences](https://github.com/coreytcallaghan/inaturalist_preferences)) and will be permanently

244 archived in a Zenodo repository upon acceptance of this article.

245

## 246 RESULTS

247 A total of 255,727,592 eBird and 1,107,224 iNaturalist observations were used in our analysis.

248 At the state level, the number of eBird checklists and the number of iNaturalist observations

249 were strongly correlated (Figure 2a;  $R^2 = 0.58$ , p-value < 0.001). Similarly, at the species level,

250 the total number of iNaturalist observations and eBird observations for a given species was

251 strongly correlated (Figure 2b;  $R^2 = 0.9$ ), and for both datasets the number of observations per

252 species was positively-skewed (Figure S1). We also found that the percent of eBird checklists a

253 species was found on and the percent of total iNaturalist observations a species comprised was

254 strongly correlated among states (Figure S2), suggesting that species are sampled to a

255 proportionally similar extent in unstructured and semi-structured citizen science projects.

256

257 Across the 507 species included in our analyses (Table S1), we showed that larger species were

258 more likely to be over-represented in the unstructured citizen science dataset, and this was true in

259 most states, as illustrated by the empirical comparison (Figure 3a). The empirical comparison

260 also showed over-representation of flock size in the unstructured dataset, although some states

261 showed a negative relationship indicating the possibility that this trait varies in space (Figure 3b).  
262 There was no discernible pattern in the relationship between color and over- or under-  
263 representation in iNaturalist data (Figure 3c), and there was some evidence that Least Concern  
264 species were over-represented in the iNaturalist data (Figure 3d).

265

266 In contrast to our empirical comparisons (Figure 3), our mixed effects multiple regression linear  
267 model (N=3986) with state as a random effect (Figure 4) found strong evidence that body size  
268 (parameter estimate=0.049; 95% CI=0.023, 0.073) and flock size (parameter estimate=0.051;  
269 95% CI=0.034, 0.069) were over-represented in iNaturalist compared with eBird; moderate  
270 evidence that common species were over-represented (parameter estimate=0.027; 95% CI=-  
271 0.003, 0.058); and no evidence that color influenced the over- or under-representation of a  
272 species in iNaturalist (parameter estimate=-0.008; 95% CI=-0.064, 0.048). The patterns found in  
273 the multiple regression model qualitatively matched that of the individual trait models, where  
274 more observations were included in some instances (see Table S2).

275

## 276 DISCUSSION

277 We compared two popular citizen science platforms throughout the continental United States and  
278 found that there was strong agreement between the relative number of observations of a species  
279 in iNaturalist and eBird, albeit there were about 200 times more observations in eBird than  
280 iNaturalist. This suggests that species are observed at similar rates in both citizen science  
281 projects — i.e., the inherent processes driving observation in both unstructured and semi-  
282 structured citizen science projects are similar. Nevertheless, in support of our predictions (Figure  
283 1) we found strong evidence that large-bodied birds are over-represented in the unstructured

284 citizen science dataset compared with the semi-structured dataset. We also found moderate  
285 evidence that common species were over-represented in the unstructured data, and weak  
286 evidence that species in large flocks were over-represented. In contrast to our prediction,  
287 however, we found no evidence that brightly-colored species were over-represented in  
288 unstructured citizen science data.

289  
290 Our finding that large-bodied birds were over-represented in an unstructured citizen science  
291 dataset is probably because larger-bodied birds are more detectable<sup>53,65</sup>. This confirms previous  
292 research which has shown that smaller-bodied taxa are under-represented in citizen science data  
293<sup>66,67,68</sup>, but this may not be the case for some taxa such as mammals<sup>69</sup>. However, it is difficult to  
294 know whether this is an inherent preference shown by users of the unstructured citizen science  
295 data, or if this comes about as part of the recording process (e.g., species' detectability;<sup>50</sup>).  
296 Species detectability is complex and can be linked to species' mobility or habitat preferences of  
297 the species themselves; for example, large-bodied wading birds generally occurring in open  
298 wetlands are more easily detected than small-bodied songbirds generally occurring in dense  
299 forest understory.

300  
301 Related to detectability, an important distinction between iNaturalist and eBird is how  
302 identifications are made. For an observer to make a record in iNaturalist, usually a photograph is  
303 uploaded (although sound recordings are also accepted). Because a photograph is needed, the  
304 detectability process is two-fold — first, it needs to be detected, and second, it needs to be  
305 photographed, which is likely easier for many large-bodied birds. Longer lenses, often restricted  
306 to serious photographers, may be needed to photograph smaller-bodied birds whereas

307 smartphones can usually capture a sufficient image of a larger-bodied bird. In contrast to  
308 iNaturalist, in eBird, a lot of identifications are made acoustically, and identification can  
309 sometimes also use contextual clues such as behavior or habitat of the bird — often difficult to  
310 capture in a photograph. Most traits analyzed here are related to visual encounter/identification,  
311 thus likely explaining the differences found between the unstructured iNaturalist and the semi-  
312 structured eBird data. To illustrate this difference, in New York state, the most under-represented  
313 species in iNaturalist (i.e., with the lowest residuals) are Marsh Wren, American Crow, Warbling  
314 Vireo, Least Flycatcher, Willow Flycatcher – all species that are identified largely acoustically.  
315 In contrast, the most over-represented species in iNaturalist (i.e., with the highest residuals) are  
316 House Sparrow, American Robin, Palm Warbler, Northern Mockingbird – all species that are  
317 easy to visually see and thus detect and photograph (Table S1). Therefore, the bias towards  
318 large-bodied birds in the unstructured data is probably a result of detectability and the ability to  
319 capture a photograph<sup>53</sup>. Photographs can also be uploaded to eBird, and a further test of this  
320 hypothesis could interrogate the species in eBird which have photographs uploaded. This process  
321 is similar in insects, for example, which are generally small, but larger insects (e.g., butterflies)  
322 are both easier to observe, photograph, and identify — making it likely that the biases we found  
323 in birds generalize to insects as well. Indeed, a study of bugs and beetles found that smaller  
324 species are typically less represented in citizen science data<sup>68</sup>. Importantly, because this body  
325 size bias is systematic, it is likely easier to model as we know that this data is not missing at  
326 random (e.g.,<sup>70</sup>) and thus body size can be included in various modelling processes when using  
327 unstructured citizen science data (e.g.,<sup>67</sup>).

328



329 Similar to body size, we found that birds which occur in larger groups (i.e., flocks) are over-  
330 represented in the unstructured dataset. This, again, may be inherently linked to the recording  
331 process, rather than a specific bias or preference of the observers themselves. This is because  
332 common birds, that occur in large flocks, are more likely to be seen and thus submitted to the  
333 unstructured citizen science data <sup>65</sup>. A larger flock will likely also provide more opportunities to  
334 capture a photograph than when observing a single individual, as has been illustrated in the  
335 detectability of animals from aerial surveys by professionals <sup>71</sup>.

336  
337 One explanation for the least concern birds being somewhat over-represented in iNaturalist is  
338 user behavior — eBird data are more likely to be derived from avid birdwatchers (e.g., those that  
339 search out uncommon birds and keep serious lists) compared with iNaturalist data which may be  
340 derived from more recreational birdwatchers that focus on ‘backyard’ species. The types of  
341 participants, and their motivations, of iNaturalist and eBird are therefore likely very different as  
342 has generally been shown among citizen science projects (e.g., <sup>72</sup>). Participants submitting  
343 observations to eBird are likely better at identifying birds than those submitting to iNaturalist and  
344 can also rely on acoustic and contextual clues to make identifications, as discussed above.  
345 Importantly, our analysis focused on only unstructured versus semi-structured data, but future  
346 work should expand this comparison to include structured datasets (e.g., breeding bird surveys)  
347 to understand if the biases found here also exist when compared with more structured datasets.  
348 For example, there may be a skew in eBird data towards rare birds when compared to  
349 standardized surveys (e.g., breeding bird surveys) resulting from birders preferentially adding  
350 rare and uncommon species. Such a finding would further highlight the divergence in behavior  
351 between the users of iNaturalist and eBird.

352  
353 The lack of signal of the colorfulness of a species in predicting over-representation in iNaturalist  
354 could suggest that iNaturalist users are not limited by ‘attractiveness/aesthetics’ but mostly by  
355 detectability, as discussed above (Figure 4). Alternatively, the lack of a signal here could be a  
356 result of the comparison being between a semi-structured and an unstructured dataset – i.e., both  
357 eBird and iNaturalist are skewed towards more colorful species, and a comparison with a  
358 structured dataset will help test this hypothesis. Quantifying the influence of color on  
359 detectability remains a challenge (e.g., <sup>73</sup>). In contrast to our results, others have demonstrated a  
360 clear preference of ‘color’ by the general public in online google searches of birds <sup>74</sup>. However,  
361 the role of aesthetics, or color, by the public may be complex as illustrated by one study which  
362 found that only blue and yellow were significant in determining bird ‘beauty’ <sup>75</sup>. In other taxa,  
363 more colorful insect species are more commonly reported <sup>68</sup>, as well as more patterned and  
364 morphologically interesting species. This may suggest, at least in the case of insects, that  
365 contributors are selecting subjects based on their visual aesthetics, not just their detectability.  
366 The discrepancies between our results and that of <sup>68</sup> suggest that the influence of traits may vary  
367 between different taxa, making it important to explore these relationships for a range of  
368 organisms rather than extrapolating the results of birds, or bugs and beetles, to other groups.

369  
370 While citizen science data are undoubtedly valuable for ecology and conservation <sup>4, 76, 77</sup>, there  
371 remain limits to the use of citizen science datasets <sup>13, 78</sup>. The ability to sample remote regions, for  
372 example, will likely remain a limitation in citizen science data, and this has been well-recognized  
373 <sup>17</sup>. Quantifying the limits of citizen science datasets for use in ecology and conservation remains  
374 an important step for the future widespread use of citizen science data in ecology and

375 conservation. Data-integration — where noisy citizen science data are integrated with  
376 professionally-curated datasets — will likely be increasingly important in the future use of  
377 citizen science data <sup>79, 80</sup>. By knowing the biases present in citizen science data, experts can  
378 preferentially generate data that maximize the integration process, for example by collecting data  
379 from remote regions. Further, professional scientists should use limited funding to target species  
380 that are likely to be under-represented in citizen science datasets — i.e., rare, small-bodied,  
381 species.

382

383 Ultimately, citizen science data will continue to perform, at least in part, a substantial role in the  
384 future of ecology and conservation research <sup>44</sup>. Understanding, documenting, and quantifying the  
385 biases associated with these data remains an important first step before the widespread use of  
386 these data in answering ecological questions and biodiversity monitoring <sup>5</sup>. Our results highlight  
387 that for birds, semi-structured eBird out-samples unstructured iNaturalist data, but the number of  
388 observations recorded per species are strongly correlated between the two platforms. When  
389 looking at the differences in this relationship, it is clear that biases exist, likely due to the biases  
390 in the unstructured iNaturalist data. We note that we compared the unstructured dataset to a  
391 semi-structured dataset, and the semi-structured dataset does not necessarily represent the  
392 “truth”. The biases found here, could also be present when comparing a semi-structured dataset  
393 to true density or abundance of birds in the landscape. To better understand these differences,  
394 future research in this space should continue to focus on quantifying and documenting biases in  
395 citizen science data, and understanding how these biases change from unstructured to semi-  
396 structured to structured citizen science platforms. Nevertheless, our results demonstrate the

397 importance of using species-specific traits, when modelling citizen science datasets <sup>27, 29, 52, 81, 82,</sup>  
398 83,84.

399

#### 400 ACKNOWLEDGEMENTS

401 We thank the countless contributors to both eBird and iNaturalist who contribute their  
402 observations as well as the Cornell Lab of Ornithology and the California Academy of Sciences  
403 to their commitment of making citizen science data open access. CTC, HMP, and MH  
404 acknowledge funding of iDiv via the German Research Foundation (DFG FZT 118). CTC was  
405 supported by a Marie Skłodowska-Curie Individual Fellowship (No 891052).

406

#### 407 AUTHOR CONTRIBUTIONS

408 CTC conceived and led the study with input from all authors. CTC performed the analyses with  
409 input from all authors. CTC wrote the first version of the manuscript and all authors contributed  
410 to editing the manuscript.

411

#### 412 COMPETING INTERESTS

413 The author(s) declare no competing interests.

## 414 REFERENCES

- 415 1. Pocock, M. J., Tweddle, J. C., Savage, J., Robinson, L. D. & Roy, H. E. The diversity and  
416 evolution of ecological and environmental citizen science. *PLoS One* **12**, e0172579  
417 (2017).
- 418 2. Chandler, M. *et al.* Contribution of citizen science towards international biodiversity  
419 monitoring. *Biological Conservation* **213**, 280–294 (2017).
- 420 3. Chandler, M. *et al.* Involving citizen scientists in biodiversity observation. in *The GEO*  
421 *handbook on biodiversity observation networks* 211–237 (Springer, Cham, 2017).
- 422 4. McKinley, D. C. *et al.* Citizen science can improve conservation science, natural resource  
423 management, and environmental protection. *Biological Conservation* **208**, 15–28 (2017).
- 424 5. Pereira, H. M. *et al.* Monitoring essential biodiversity variables at the species level. in  
425 *The GEO handbook on biodiversity observation networks* 79–105 (Springer, Cham,  
426 2017).
- 427 6. Wiggins, A. & Crowston, K. From conservation to crowdsourcing: A typology of citizen  
428 science. in *2011 44th hawaii international conference on system sciences* 1–10 (IEEE,  
429 2011).
- 430 7. Haklay, M. Citizen science and volunteered geographic information: Overview and  
431 typology of participation. in *Crowdsourcing geographic knowledge* 105–122 (Springer,  
432 2013).
- 433 8. Kelling, S. *et al.* Using semistructured surveys to improve citizen science data for  
434 monitoring biodiversity. *BioScience* **69**, 170–179 (2019).

- 435 9. Welvaert, M. & Caley, P. Citizen surveillance for environmental monitoring: Combining  
436 the efforts of citizen science and crowdsourcing in a quantitative data framework.  
437 *SpringerPlus* **5**, 1890 (2016).
- 438 10. Callaghan, C. T., Rowley, J. J., Cornwell, W. K., Poore, A. G. & Major, R. E. Improving  
439 big citizen science data: Moving beyond haphazard sampling. *PLoS Biology* **17**,  
440 e3000357 (2019).
- 441 11. Bonter, D. N. & Cooper, C. B. Data validation in citizen science: A case study from  
442 project FeederWatch. *Frontiers in Ecology and the Environment* **10**, 305–307 (2012).
- 443 12. Kosmala, M., Wiggins, A., Swanson, A. & Simmons, B. Assessing data quality in citizen  
444 science. *Frontiers in Ecology and the Environment* **14**, 551–560 (2016).
- 445 13. Burgess, H. K. *et al.* The science of citizen science: Exploring barriers to use as a primary  
446 research tool. *Biological Conservation* **208**, 113–120 (2017).
- 447 14. Courter, J. R., Johnson, R. J., Stuyck, C. M., Lang, B. A. & Kaiser, E. W. Weekend bias  
448 in citizen science data reporting: Implications for phenology studies. *International*  
449 *journal of biometeorology* **57**, 715–720 (2013).
- 450 15. Sullivan, B. L. *et al.* The eBird enterprise: An integrated approach to development and  
451 application of citizen science. *Biological Conservation* **169**, 31–40 (2014).
- 452 16. Kelling, S. *et al.* Can observation skills of citizen scientists be estimated using species  
453 accumulation curves? *PloS one* **10**, e0139600 (2015).
- 454 17. Tiago, P., Ceia-Hasse, A., Marques, T. A., Capinha, C. & Pereira, H. M. Spatial  
455 distribution of citizen science casuistic observations for different taxonomic groups.  
456 *Scientific reports* **7**, 1–9 (2017).

- 457 18. Geldmann, J. *et al.* What determines spatial bias in citizen science? Exploring four  
458 recording schemes with different proficiency requirements. *Diversity and Distributions*  
459 **22**, 1139–1149 (2016).
- 460 19. Callaghan, C. T. *et al.* Three frontiers for the future of biodiversity research using citizen  
461 science data. *BioScience* **71**, 55–63 (2021).
- 462 20. Ward, D. F. Understanding sampling and taxonomic biases recorded by citizen scientists.  
463 *Journal of insect conservation* **18**, 753–756 (2014).
- 464 21. Troudet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R. & Legendre, F. Taxonomic bias in  
465 biodiversity data and societal preferences. *Scientific reports* **7**, 1–14 (2017).
- 466 22. Martín-López, B., Montes, C., Ramírez, L. & Benayas, J. What drives policy decision-  
467 making related to species conservation? *Biological Conservation* **142**, 1370–1380 (2009).
- 468 23. Boakes, E. H. *et al.* Distorted views of biodiversity: Spatial and temporal bias in species  
469 occurrence data. *PLoS Biol* **8**, e1000385 (2010).
- 470 24. Aceves-Bueno, E. *et al.* The accuracy of citizen science data: A quantitative review.  
471 *Bulletin of the Ecological Society of America* **98**, 278–290 (2017).
- 472 25. Davies, T. K., Stevens, G., Meekan, M. G., Struve, J. & Rowcliffe, J. M. Can citizen  
473 science monitor whale-shark aggregations? Investigating bias in mark–recapture  
474 modelling using identification photographs sourced from the public. *Wildlife Research*  
475 **39**, 696–704 (2013).
- 476 26. Crall, A. W. *et al.* Assessing citizen science data quality: An invasive species case study.  
477 *Conservation Letters* **4**, 433–442 (2011).

- 478 27. Strien, A. J. van, Swaay, C. A. van & Termaat, T. Opportunistic citizen science data of  
479 animal species produce reliable estimates of distribution trends if analysed with  
480 occupancy models. *Journal of Applied Ecology* **50**, 1450–1458 (2013).
- 481 28. Johnston, A., Moran, N., Musgrove, A., Fink, D. & Baillie, S. R. Estimating species  
482 distributions from spatially biased citizen science data. *Ecological Modelling* **422**,  
483 108927 (2020).
- 484 29. Tiago, P., Pereira, H. M. & Capinha, C. Using citizen science data to estimate climatic  
485 niches and species distributions. *Basic and Applied Ecology* **20**, 75–85 (2017).
- 486 30. Sullivan, B. L. *et al.* Using open access observational data for conservation action: A case  
487 study for birds. *Biological Conservation* **208**, 5–14 (2017).
- 488 31. Callaghan, C. T. *et al.* Citizen science data accurately predicts expert-derived species  
489 richness at a continental scale when sampling thresholds are met. *Biodiversity and*  
490 *Conservation* **29**, 1323–1337 (2020).
- 491 32. Birkin, L. & Goulson, D. Using citizen science to monitor pollination services.  
492 *Ecological Entomology* **40**, 3–11 (2015).
- 493 33. Delaney, D. G., Sperling, C. D., Adams, C. S. & Leung, B. Marine invasive species:  
494 Validation of citizen science and implications for national monitoring networks.  
495 *Biological Invasions* **10**, 117–128 (2008).
- 496 34. Schultz, C. B., Brown, L. M., Pelton, E. & Crone, E. E. Citizen science monitoring  
497 demonstrates dramatic declines of monarch butterflies in western north america.  
498 *Biological Conservation* **214**, 343–346 (2017).
- 499 35. Bird, T. J. *et al.* Statistical solutions for error and bias in global citizen science datasets.  
500 *Biological Conservation* **173**, 144–154 (2014).



- 501 36. Isaac, N. J., Strien, A. J. van, August, T. A., Zeeuw, M. P. de & Roy, D. B. Statistics for  
502 citizen science: Extracting signals of change from noisy ecological data. *Methods in*  
503 *Ecology and Evolution* **5**, 1052–1060 (2014).
- 504 37. Dickinson, J. L. *et al.* The current state of citizen science as a tool for ecological research  
505 and public engagement. *Frontiers in Ecology and the Environment* **10**, 291–297 (2012).
- 506 38. Bonney, R. *et al.* Next steps for citizen science. *Science* **343**, 1436–1437 (2014).
- 507 39. Jordan, R. C., Gray, S. A., Howe, D. V., Brooks, W. R. & Ehrenfeld, J. G. Knowledge  
508 gain and behavioral change in citizen-science programs. *Conservation biology* **25**, 1148–  
509 1154 (2011).
- 510 40. Crall, A. W. *et al.* The impacts of an invasive species citizen science training program on  
511 participant attitudes, behavior, and science literacy. *Public Understanding of Science* **22**,  
512 745–764 (2013).
- 513 41. Jordan, R. C., Ballard, H. L. & Phillips, T. B. Key issues and new approaches for  
514 evaluating citizen-science learning outcomes. *Frontiers in Ecology and the Environment*  
515 **10**, 307–309 (2012).
- 516 42. Evans, C. *et al.* The neighborhood nestwatch program: Participant outcomes of a citizen-  
517 science ecological research project. *Conservation Biology* **19**, 589–594 (2005).
- 518 43. Haywood, B. K., Parrish, J. K. & Dolliver, J. Place-based and data-rich citizen science as  
519 a precursor for conservation action. *Conservation Biology* **30**, 476–486 (2016).
- 520 44. Pocock, M. J. *et al.* A vision for global biodiversity monitoring with citizen science. in  
521 *Advances in ecological research* vol. 59 169–223 (Elsevier, 2018).

- 522 45. Tiago, P., Gouveia, M. J., Capinha, C., Santos-Reis, M. & Pereira, H. M. The influence  
523 of motivational factors on the frequency of participation in citizen science activities.  
524 *Nature Conservation* **18**, 61 (2017).
- 525 46. Isaac, N. J. & Pocock, M. J. Bias and information in biological records. *Biological*  
526 *Journal of the Linnean Society* **115**, 522–531 (2015).
- 527 47. Angulo, E. & Courchamp, F. Rare species are valued big time. *PloS one* **4**, e5215 (2009).
- 528 48. Booth, J. E., Gaston, K. J., Evans, K. L. & Armsworth, P. R. The value of species rarity  
529 in biodiversity recreation: A birdwatching example. *Biological Conservation* **144**, 2728–  
530 2732 (2011).
- 531 49. Rowley, J. J. *et al.* FrogID: Citizen scientists provide validated biodiversity data on frogs  
532 of australia. *Herpetological Conservation and Biology* **14**, 155–70 (2019).
- 533 50. Boakes, E. H. *et al.* Patterns of contribution to citizen science biodiversity projects  
534 increase understanding of volunteers recording behaviour. *Scientific reports* **6**, 33051  
535 (2016).
- 536 51. Garrard, G. E., McCarthy, M. A., Williams, N. S., Bekessy, S. A. & Wintle, B. A. A  
537 general model of detectability using species traits. *Methods in Ecology and Evolution* **4**,  
538 45–52 (2013).
- 539 52. Denis, T. *et al.* Biological traits, rather than environment, shape detection curves of large  
540 vertebrates in neotropical rainforests. *Ecological Applications* **27**, 1564–1577 (2017).
- 541 53. Sólymos, P., Matsuoka, S. M., Stralberg, D., Barker, N. K. & Bayne, E. M. Phylogeny  
542 and species traits predict bird detectability. *Ecography* **41**, 1595–1603 (2018).
- 543 54. Wood, C., Sullivan, B., Iliff, M., Fink, D. & Kelling, S. eBird: Engaging birders in  
544 science and conservation. *PLoS Biol* **9**, e1001220 (2011).

- 545 55. GBIF.org (3<sup>rd</sup> December 2019). GBIF occurrence download.  
546 <https://www.doi.org/10.15468/dl.lpwczr>
- 547 56. Gilfedder, M. *et al.* Brokering trust in citizen science. *Society & natural resources* **32**,  
548 292–302 (2019).
- 549 57. Callaghan, C., Lyons, M., Martin, J., Major, R. & Kingsford, R. Assessing the reliability  
550 of avian biodiversity measures of urban greenspaces using eBird citizen science data.  
551 *Avian Conservation and Ecology* **12**, (2017).
- 552 58. Johnston, A. *et al.* Best practices for making reliable inferences from citizen science data:  
553 Case study using eBird to estimate species distributions. *BioRxiv* 574392 (2019).
- 554 59. Myhrvold, N. P. *et al.* An amniote life-history database to perform comparative analyses  
555 with birds, mammals, and reptiles: Ecological archives E096-269. *Ecology* **96**, 3109–  
556 3109 (2015).
- 557 60. Dale, J., Dey, C. J., Delhey, K., Kempnaers, B. & Valcu, M. The effects of life history  
558 and sexual selection on male and female plumage colouration. *Nature* **527**, 367–370  
559 (2015).
- 560 61. R Core Team. *R: A language and environment for statistical computing*. (R Foundation  
561 for Statistical Computing, 2020).
- 562 62. Wickham, H. *et al.* Welcome to the tidyverse. *Journal of Open Source Software* **4**, 1686  
563 (2019).
- 564 63. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting linear mixed-effects models  
565 using lme4. *Journal of Statistical Software* **67**, 1–48 (2015).
- 566 64. Kuznetsova, A., Brockhoff, P. B. & Christensen, R. H. B. lmerTest package: Tests in  
567 linear mixed effects models. *Journal of Statistical Software* **82**, 1–26 (2017).

- 568 65. Johnston, A. *et al.* Species traits explain variation in detectability of UK birds. *Bird Study*  
569 **61**, 340–350 (2014).
- 570 66. Steen, V. A., Elphick, C. S. & Tingley, M. W. An evaluation of stringent filtering to  
571 improve species distribution models from citizen science data. *Diversity and*  
572 *Distributions* **25**, 1857–1869 (2019).
- 573 67. Henckel, L., Bradter, U., Jönsson, M., Isaac, N. J. & Snäll, T. Assessing the usefulness of  
574 citizen science data for habitat suitability modelling: Opportunistic reporting versus  
575 sampling based on a systematic protocol. *Diversity and Distributions* **26**, 1276–1290  
576 (2020).
- 577 68. Caley, P., Welvaert, M. & Barry, S. C. Crowd surveillance: Estimating citizen science  
578 reporting probabilities for insects of biosecurity concern. *Journal of Pest Science* **93**,  
579 543–550 (2020).
- 580 69. Périquet, S., Roxburgh, L., Roux, A. le & Collinson, W. J. Testing the value of citizen  
581 science for roadkill studies: A case study from south africa. *Frontiers in Ecology and*  
582 *Evolution* **6**, 15 (2018).
- 583 70. Nakagawa, S. & Freckleton, R. P. Model averaging, missing data and multiple  
584 imputation: A case study for behavioural ecology. *Behavioral Ecology and Sociobiology*  
585 **65**, 103–116 (2011).
- 586 71. Schlossberg, S., Chase, M. & Griffin, C. Using species traits to predict detectability of  
587 animals on aerial surveys. *Ecological Applications* **28**, 106–118 (2018).
- 588 72. Aristeidou, M., Scanlon, E. & Sharples, M. Profiles of engagement in online  
589 communities of citizen science participation. *Computers in Human Behavior* **74**, 246–256  
590 (2017).

- 591 73. Troschianko, J., Skelhorn, J. & Stevens, M. Quantifying camouflage: How to predict  
592 detectability from appearance. *BMC evolutionary biology* **17**, 1–13 (2017).
- 593 74. Schuetz, J. G. & Johnston, A. Characterizing the cultural niches of North American birds.  
594 *Proceedings of the National Academy of Sciences* **22**, 10868–10873 (2019).
- 595 75. Lišková, S. & Frynta, D. What determines bird beauty in human eyes? *Anthrozoös* **26**,  
596 27–41 (2013).
- 597 76. Tulloch, A. I., Possingham, H. P., Joseph, L. N., Szabo, J. & Martin, T. G. Realising the  
598 full potential of citizen science monitoring programs. *Biological Conservation* **165**, 128–  
599 138 (2013).
- 600 77. Kobori, H. *et al.* Citizen science: A new approach to advance ecology, education, and  
601 conservation. *Ecological research* **31**, 1–19 (2016).
- 602 78. Callaghan, C. T., Poore, A. G., Major, R. E., Rowley, J. J. & Cornwell, W. K. Optimizing  
603 future biodiversity sampling by citizen scientists. *Proceedings of the Royal Society B* **286**,  
604 20191487 (2019).
- 605 79. Pacifici, K. *et al.* Integrating multiple data sources in species distribution modeling: A  
606 framework for data fusion. *Ecology* **98**, 840–850 (2017).
- 607 80. Robinson, O. J. *et al.* Integrating citizen science data with expert surveys increases  
608 accuracy and spatial extent of species distribution models. *Diversity and Distributions* **26**,  
609 976–986 (2020).
- 610 81. Strien, A. J. van, Termaat, T., Groenendijk, D., Mensing, V. & Kery, M. Site-occupancy  
611 models may offer new opportunities for dragonfly monitoring based on daily species  
612 lists. *Basic and Applied Ecology* **11**, 495–503 (2010).

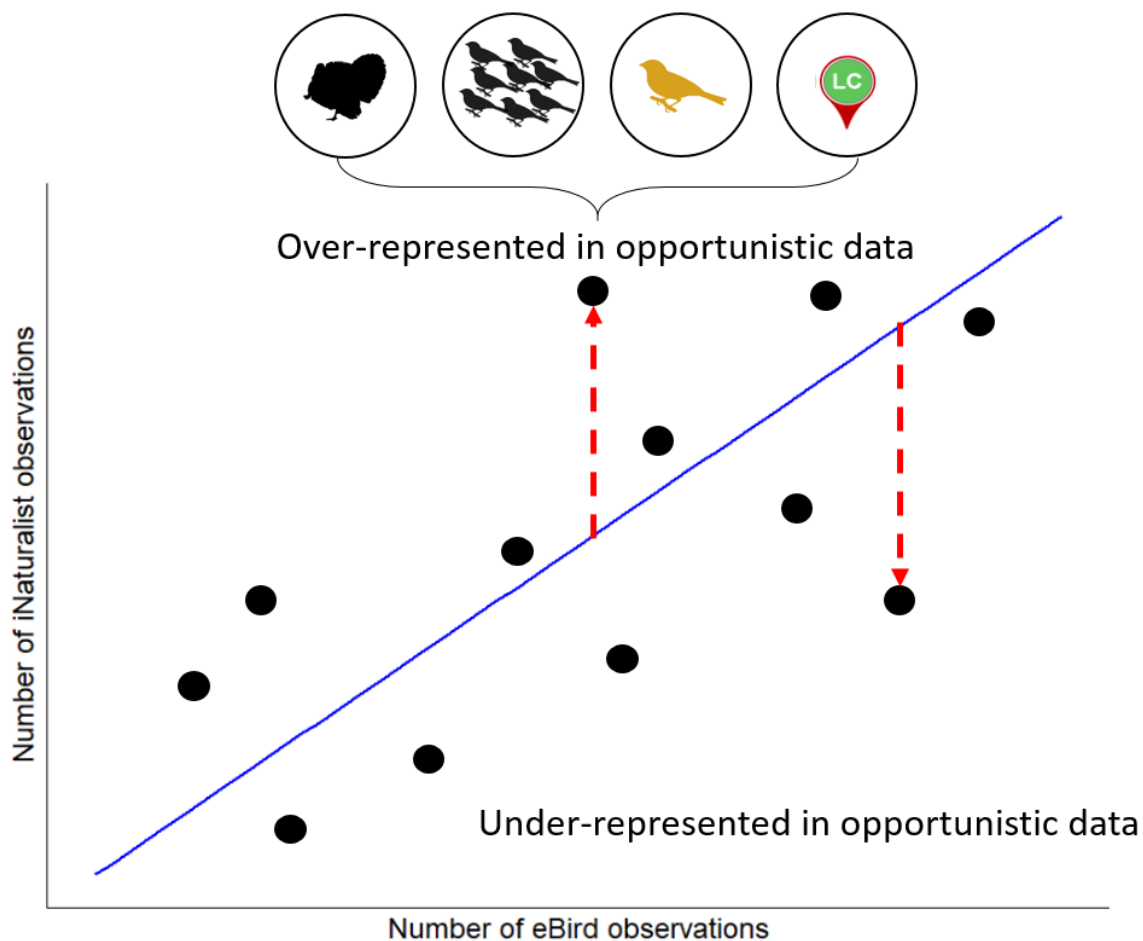
- 613 82. Van der Wal, R. *et al.* Mapping species distributions: A comparison of skilled naturalist  
614 and lay citizen science recording. *Ambio* **44**, 584–600 (2015).
- 615 83. Dennis, E. B., Morgan, B. J., Brereton, T. M., Roy, D. B. & Fox, R. Using citizen science  
616 butterfly counts to predict species population trends. *Conservation biology* **31**, 1350–  
617 1361 (2017).
- 618 84. Stoudt, S., Goldstein, B. R. & De Valpine, P. Identifying charismatic bird species and  
619 traits with community science data. bioRxiv. <https://doi.org/10.1101/2021.06.05.446577>

620

621

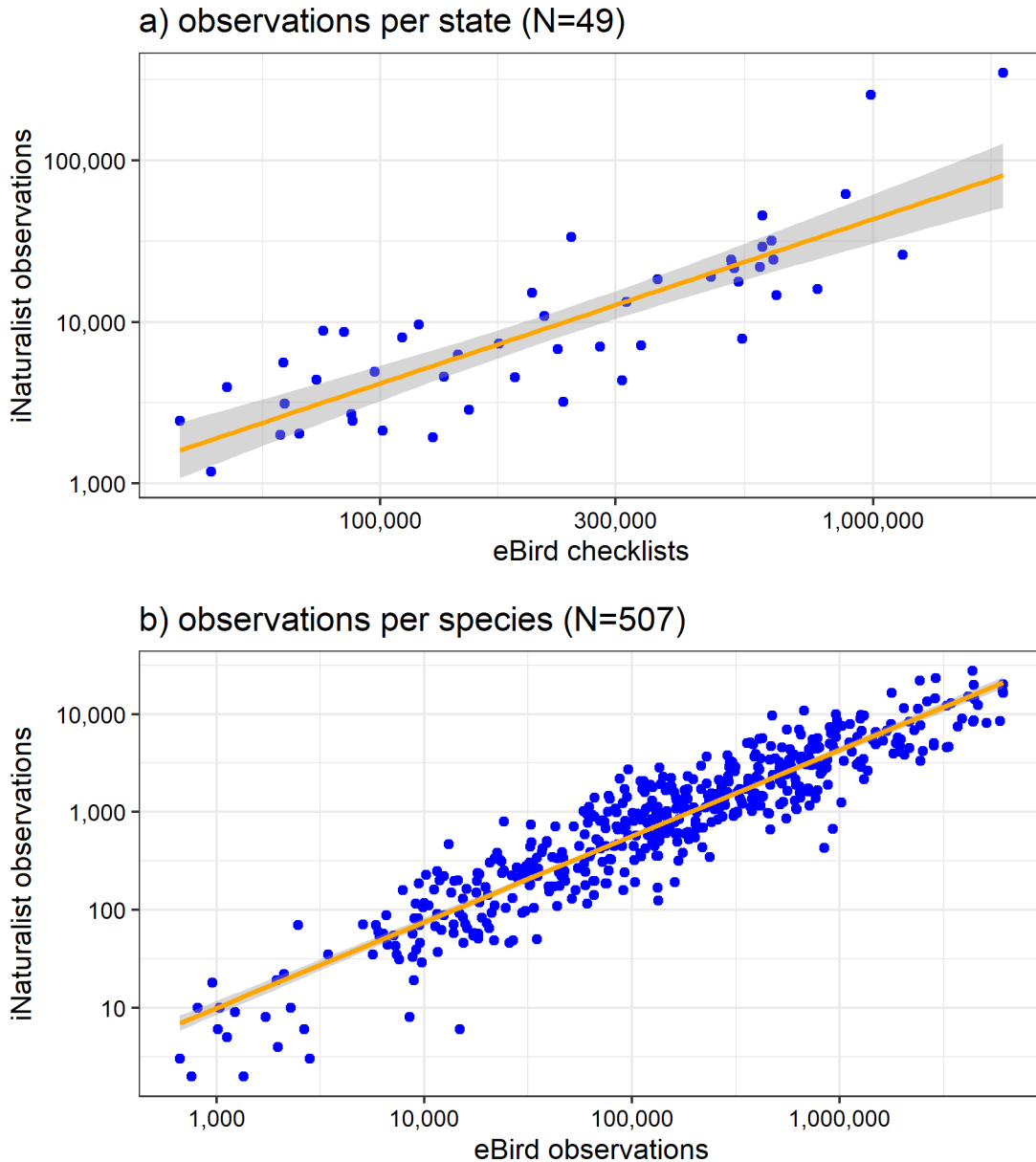
622 FIGURES

623



624

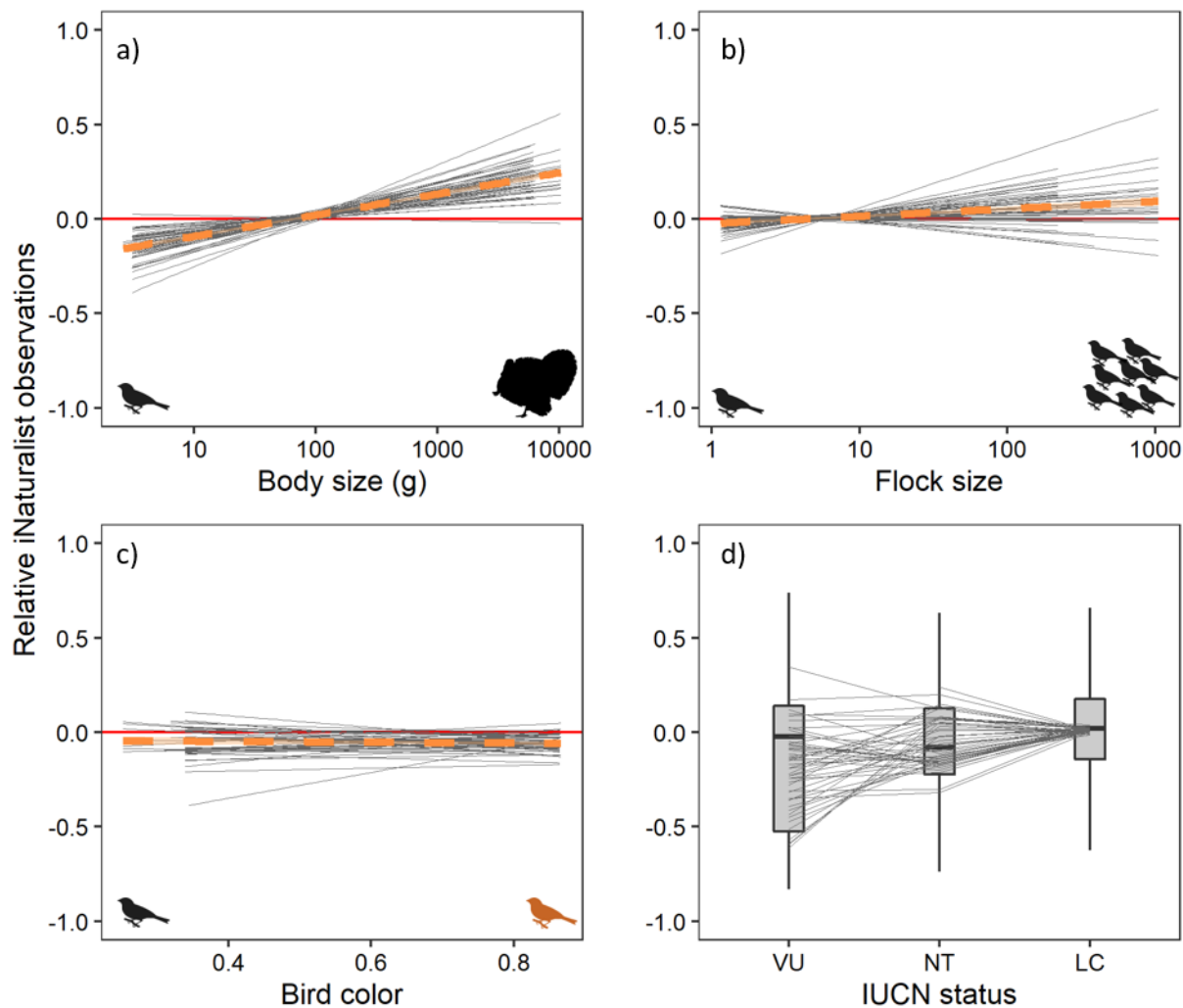
625 **Figure 1.** A conceptual figure depicting the methods used in our analysis. We used the residual  
 626 from the relationship between the number of eBird observations (i.e., semi-structured citizen  
 627 science observations) and iNaturalist observations (i.e., unstructured citizen science  
 628 observations) to quantify the over- or under-representation of a species in unstructured citizen  
 629 science data. We predicted that species which were over-represented in unstructured iNaturalist  
 630 data would be larger in size, occur more frequently in large flocks, be brighter in color, and be  
 631 categorized as Least Concern IUCN status (a proxy for commonness).



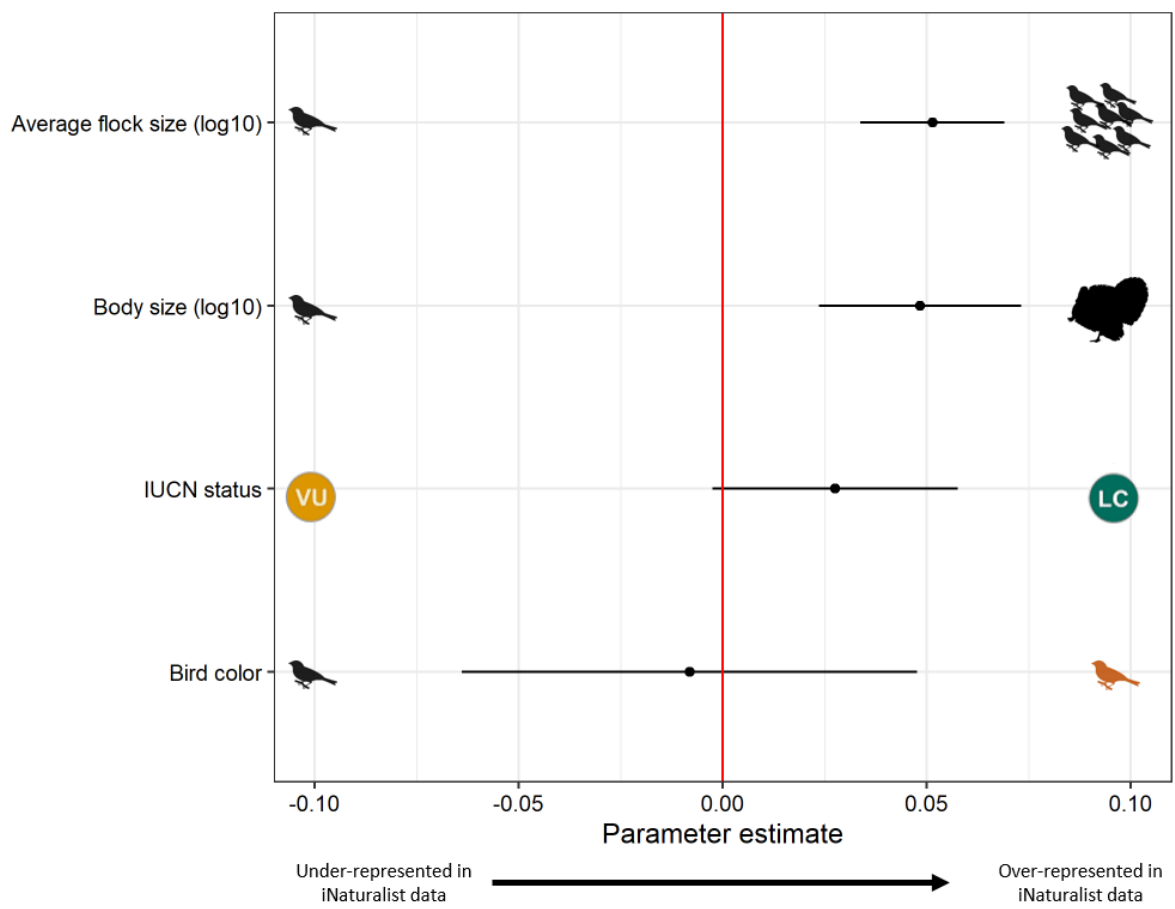
632

633 **Figure 2.** a) The relationship between the total number of eBird checklists and total number of  
 634 iNaturalist observations for 49 states, including the District of Columbia. There was strong  
 635 evidence that these variables were correlated ( $R^2=0.58$ ,  $p$ -value  $<0.001$ ) suggesting that sampling  
 636 between datasets is correlated among states. b) The relationship between the number of  
 637 observations for a species from eBird (x-axis) and the number of observations for a species from  
 638 iNaturalist (y-axis) for only eBird species which were found on  $>1\%$  of eBird checklists.





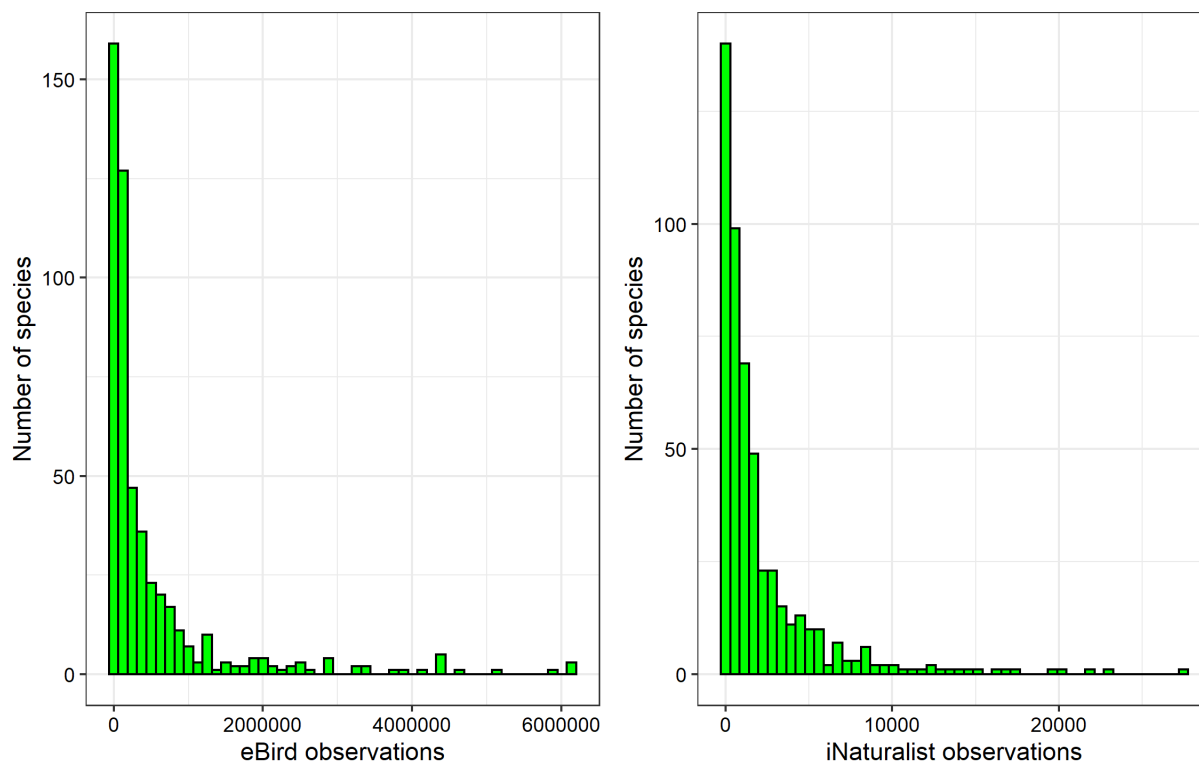
639  
 640 **Figure 3.** The relationship between a) body size of a species, b) flock size, c) color and d)  
 641 commonness and the residuals of a linear model fit between iNaturalist and eBird observations  
 642 (see Figure 1). These results demonstrate that there is a strong bias of body size in iNaturalist  
 643 compared with eBird. Positive values on the y-axis mean over-represented in iNaturalist and  
 644 negative values on the y-axis mean under-represented in iNaturalist. Body size and flock size are  
 645 represented on a log<sub>10</sub> scale. Each line represents a state (N=49). For a-c), the overall  
 646 relationship pooling states is represented by the orange fitted line and 95% confidence interval.  
 647



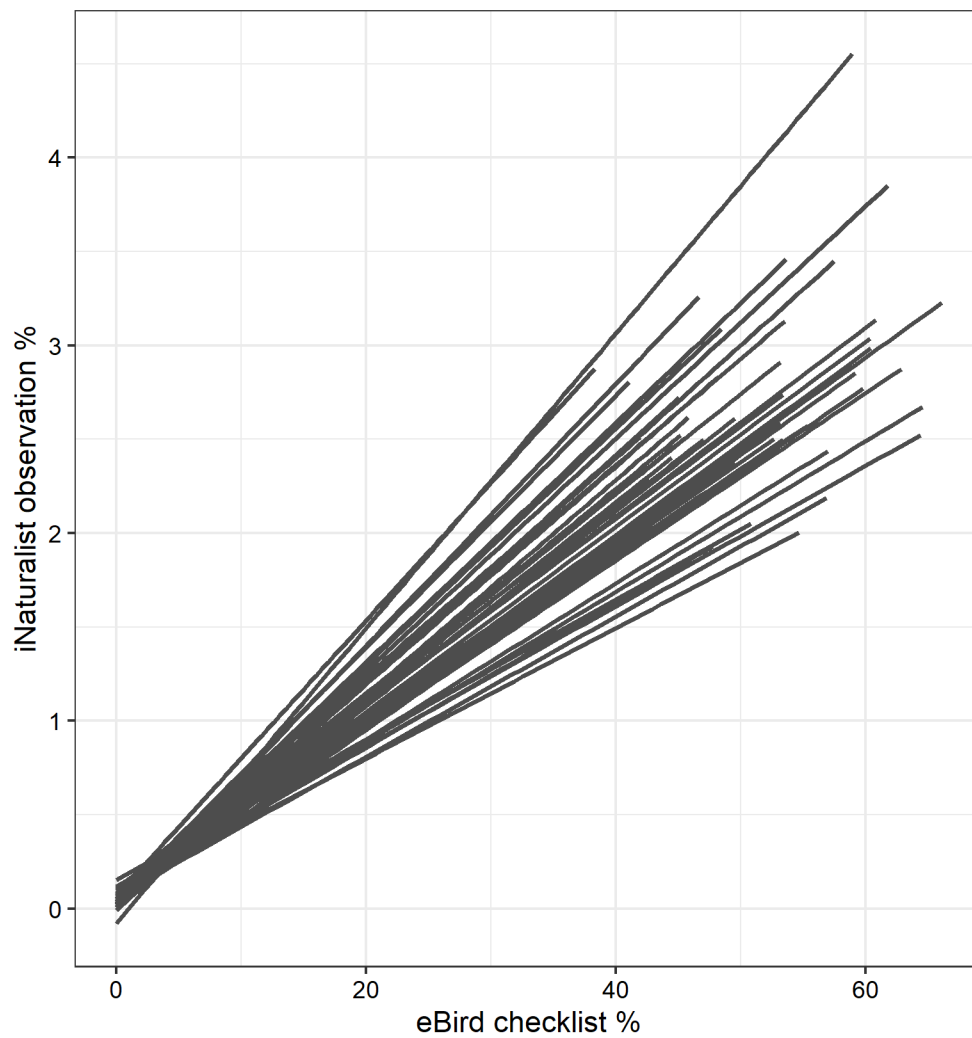
648  
649  
650  
651  
652  
653  
654  
655

**Figure 4.** Results of a linear mixed effect model where all four variables were considered simultaneously, and state was a random effect. Strong support was found for body size and flock size (their 95% confidence interval does not overlap 0), whereas moderate support was found for IUCN status, and no support was found for color.

## SUPPLEMENTARY FIGURES



**Figure S1.** Histograms of the number of observations for a species from both eBird and iNaturalist citizen science projects.



**Figure S2.** Among states (each line represents a state; N=49) we found that the percent of eBird checklists a species was found on and the percent of all iNaturalist observations a species comprised was strongly correlated.

**Table S1.** Uploaded separately. Raw data used for modelling, including the residual difference between iNaturalist and eBird, stratified by state.

**Table S2.** Results of single regression models, where each trait was treated separately, and consequently had different sample sizes in the model fit. Each model was fit with the residuals used as the response variable, the specific trait as the predictor variable, where body size and flock size were log10-transformed and IUCN was treated as an ordinal variable, and state was a random effect. This analysis was performed to confirm the results of the multiple regression mixed effects analysis presented in the main results (Figure 4).

	estimate	t	p-value	Number of obs
Body size	0.11	31.59	<0.001	7743
Color	-0.01	-0.413	0.68	4542
Flock size	0.033	6.118	<0.001	8702
IUCN status	0.078	7.73	<0.001	7629