

## Large-bodied birds are over-represented in opportunistic 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 implicit and explicit  
4 biases in citizen science datasets remains an important first step before these data are used to  
5 address ecological questions and monitor biodiversity. One source of potential bias among  
6 datasets is the difference between those citizen science programs that collect opportunistic  
7 observations and those that have semi-structured or structured protocols for submitting  
8 observations. To quantify biases in an unstructured citizen science platform, we contrasted bird  
9 observations from the iNaturalist platform with that from a semi-structured citizen science  
10 platform — eBird — for the continental United States. We tested whether four traits of species  
11 (color, group size, body size, and commonness) predicted whether a species was over-  
12 represented in the opportunistic dataset. We found strong evidence that large-bodied birds were  
13 over-represented in the opportunistic citizen science dataset; moderate evidence that common  
14 species were over-represented in the opportunistic data; moderate evidence that species in large  
15 groups were over-represented; and no evidence that colorful species were over-represented in  
16 opportunistic citizen science data. Our results suggest that biases exist in opportunistic citizen  
17 science datasets, likely as a result of the detectability of a species and the inherent recording  
18 process. Future research in this space should continue to focus on quantifying and documenting  
19 biases in citizen science data, and understanding how these biases differ among unstructured,  
20 semi-structured, and structured citizen science platforms.

21

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

## 24 INTRODUCTION

25 Citizen science, or community science, — the involvement of volunteers in scientific endeavors  
26 — is increasingly seen as a cost-effective method for biodiversity monitoring and research.

27 Accordingly, the quantity and diversity of citizen science projects in the ecological and  
28 environmental sciences is increasing <sup>1</sup>. Such projects are quickly accumulating hundreds of  
29 millions of biodiversity observations around the world annually <sup>2,3</sup> expanding the spatial and  
30 temporal scope of our understanding in ecology, conservation, and natural resource management  
31 <sup>4,5</sup>. Citizen science projects vary widely in their scope, design, and intent <sup>6,7,8</sup>. Projects can range  
32 from unstructured (e.g., little training needed to participate and contribute  
33 opportunistic/incidental observations) to semi-structured (e.g., with minimal workflows and  
34 guidelines, but additional data collected with each observation can be included) to structured  
35 (e.g., prescribed sampling in space and time by mostly trained and experienced volunteers). The  
36 level of structure consequently influences the overall data quality of a particular project <sup>9,10</sup>.

37  
38 Data quality from citizen science projects has been questioned <sup>11,12</sup>, and such concerns can act as  
39 a barrier to the widespread use of citizen science data in ecology and conservation <sup>13</sup>. These  
40 concerns arise because citizen science data can be biased temporally, spatially, and/or  
41 taxonomically. Temporally, many citizen science datasets are biased because participants  
42 frequently sample on weekends <sup>14</sup> or disproportionately during specific times of the year such as  
43 spring migration for birds <sup>15</sup>. Spatially, there is often a disproportionate number of sightings from  
44 areas with large human populations <sup>16</sup>, areas with more accessibility <sup>17</sup>, regions with high  
45 biodiversity that attract observers <sup>18</sup>, and regions of the world with higher socioeconomic status  
46 <sup>19</sup>. Taxonomic biases also exist as some taxa receive orders of magnitude more citizen science

47 observations than other taxa, evidenced by the fact that birds represent a disproportionate amount  
48 of data in the Global Biodiversity Information Facility <sup>2</sup>. Even within citizen science projects  
49 focused on specific taxa, there can be strong taxonomic biases towards particularly charismatic  
50 species or those that are readily identified <sup>20,21</sup>. Such biases are not restricted to citizen science  
51 datasets, however, and many of the same biases are also present in professionally-collected data  
52 <sup>22</sup>, such as those associated with museum specimens <sup>23</sup>.

53

54 Despite potential biases in citizen science datasets, contrasts of data from volunteer participants  
55 to those contributed by professional scientists have shown that citizen science programs can  
56 provide reliable data <sup>12,24</sup>. For example, mark-recapture models of whale sharks are reliable  
57 whether using sightings reported by the public or by experienced researchers <sup>25</sup>, and volunteers  
58 perform comparably with professionals in identifying and monitoring invasive plant species <sup>26</sup>.  
59 Moreover, recent research has demonstrated the validity of using citizen science data for  
60 ecological questions such as estimating species distributions <sup>27,28,29</sup>, managing habitat for  
61 conservation <sup>30</sup>, estimating species richness <sup>31</sup>, monitoring pollination services <sup>32</sup>, and quantifying  
62 population trends <sup>33,34</sup>. These approaches are improved when using statistical solutions to  
63 account for known biases and noise in citizen science data <sup>3,35,36</sup>.

64

65 In addition to being an excellent resource for professional scientists to better understand  
66 ecological questions, citizen science projects are beneficial for society by encouraging increased  
67 engagement of the general public with science <sup>37,38</sup>. Many citizen science projects provide  
68 learning opportunities for their volunteers. For example, participants in citizen science projects  
69 have increased their knowledge about invasive weeds <sup>39,40,41</sup>, increased their knowledge of bird

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

77

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

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

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

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

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

97 turn, the likelihood of a species being submitted to an opportunistic citizen science database.

98

99 Quantifying the implicit and explicit biases in citizen science datasets can help (1) researchers

100 using these data to better account for biases when drawing ecological conclusions, (2) the design

101 and implementation of future citizen science projects, and (3) understand what species or regions

102 may need data collection from professional scientists by understanding the ‘limits’ of citizen

103 science projects <sup>19</sup>. Here, we quantify biases in bird observation data from an unstructured,

104 opportunistic citizen science project — iNaturalist — with that from a semi-structured one —

105 eBird. We restricted our comparison to birds because (1) birds are among the most popular taxa

106 with the non-scientific public, ensuring large sample sizes in both citizen science projects, and

107 (2) data on the species traits that may influence the likelihood of opportunistic observations are

108 readily available for birds. We assessed the over-representation or under-representation of bird

109 species’ observations in the unstructured opportunistic citizen science project compared to the

110 semi-structured project (see Figure 1). We then tested the following predictions: that (1) more

111 colorful species; (2) larger species; (3) species with the ‘least concern’ IUCN status; and (4)

112 more gregarious species (i.e., with larger group sizes) are over-represented in the opportunistic

113 citizen science dataset (iNaturalist) in contrast to the semi-structured citizen science dataset

114 (eBird). Our analysis highlights the importance of considering species’ traits when using citizen

115 science data in ecological research.

116

## 117 METHODS

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

121

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

132

133 *eBird citizen science data.* eBird is one of the most successful citizen science projects in the  
134 world, with almost 1 billion bird observations globally. It was launched in 2002 by the Cornell  
135 Lab of Ornithology and focuses on collecting reliable data on the distributions and relative  
136 abundance of birds throughout the world<sup>54</sup>. It is a semi-structured project where volunteers  
137 submit ‘checklists’ of all species seen and/or heard on birding outings. These checklists provide  
138 the ability to infer absences in the dataset for any species not recorded. Observers can submit

139 checklists at any time and place of their choosing, with no set protocols in place such as how  
140 long or how far to search. However, observers are asked to indicate the duration of and distance  
141 travelled during the birding outing when submitting their checklist. Filters are set — based on  
142 spatiotemporal coordinates — which restrict the species and their associated counts that can be  
143 submitted without approval from a regional expert reviewer<sup>56</sup>. We used the eBird basic dataset  
144 (version ebd\_May-2019) and restricted our analysis to data from the contiguous United States for  
145 the period from January 2010 to May 2019. We also restricted the data used to those of the best  
146 ‘quality’ by excluding incomplete checklists, checklists that were incidental or historical, which  
147 travelled >5 km, lasted <5 min, and lasted >240 min, minimizing the leverage of outliers on  
148 analyses<sup>57, 58</sup>.

149

150 *Filtering and aggregating the citizen science datasets.* Although both datasets are global in  
151 scope, we restricted our analysis to the contiguous United States as both of these citizen science  
152 projects initiated in the United States, and thus the data are most numerous from there. For  
153 comparisons, we aggregated data at the state level. This was done to account for differences that  
154 may exist throughout the entirety of the United States including differences in user behavior and  
155 the species pools that differ geographically. For each state, the eBird and iNaturalist data were  
156 summarized, providing a list of species for each state, including the percent of total eBird  
157 checklists that a species occurred on or the percent of total observations a species accounted for,  
158 respectively. In addition, the total number of observations in that state were summarized for both  
159 eBird and iNaturalist.

160

161 We used the eBird Clements taxonomy (version 2018) and all species from iNaturalist were  
162 matched with this taxonomy. A total of 1,030 species was initially collated from the eBird  
163 checklists, but many of these only occurred one or a few times — possibly representing  
164 misidentifications that had not yet been fixed by local reviewers or escaped and exotic birds  
165 which are incorporated in the eBird dataset but not considered part of the local avifauna or of  
166 interest to our analysis here. To account for these biases, we removed species that were on <1%  
167 of eBird checklists for a given state; trimming the eBird observations to the ‘core’ suite of  
168 species that occur in a state (*sensu*<sup>57</sup>). After trimming the species and harmonizing the taxonomy  
169 with iNaturalist, there were 507 species present and considered in our main analyses presented  
170 throughout the results. Although our results here are presented using the 1% cutoff level, we  
171 tested the sensitivity of this cutoff level and found comparable results across 0, 0.5, 1, and 1.5%  
172 cutoffs.

173

#### 174 *Species-specific trait data*

175 We tested whether four predictor variables (see Figure 1) would explain the over- or under-  
176 representation of bird species in the opportunistic citizen science data. For each species, we used  
177 a proxy for their commonness/abundance, categorized according to IUCN status, taken from  
178 HBW BirdLife international checklist version 3 (<http://datazone.birdlife.org/species/taxonomy>).  
179 This variable was treated as an ordinal variable in our models (see below) and encompassed  
180 Least Concern, Vulnerable, and Near Threatened species. The three species recorded as  
181 endangered were removed from this analysis due to a lack of power at this level with so few  
182 observations. For each species we used the continuous predictor variables of (1) body size; (2)  
183 color; and (3) average group size. Body sizes (adult body mass in grams) were taken from the

184 amniote life history database compiled by <sup>59</sup> and were log-transformed to meet normality  
185 assumptions. Color was taken from <sup>60</sup> and was extracted as RGB values for six patches per  
186 species. To define a continuum of color where the brightest/most colorful (and likely most  
187 detectable species) had the highest value we combined both the ‘distance from brown’ and the  
188 ‘brightness’ of a species for the data from <sup>60</sup>. Distance from brown was defined as the maximum  
189 Euclidian distance in the cubic RGB color space from brown (R = 102, B = 68, G = 0) from  
190 either the upper or lower breast patch of a species. Brightness was defined as the maximum  
191 relative luminance (i.e.,  $0.2126R + 0.7152G + 0.0722B$ ) from either the upper or lower breast  
192 patch of a species. These two variables were combined and scaled from 0 to 1 for all species in <sup>60</sup>  
193 and this value was used as our measure of color. Group size — an approximation of the  
194 gregariousness of a species — was taken from eBird as the average number of reported  
195 individuals among all checklists where a species was reported, across all data.

196

### 197 *Statistical analysis*

198 For the species traits we ran (1) separate models for every trait and (2) a global model with all  
199 traits included. This was done because there was much missing data for species’ traits and in  
200 order to obtain maximum power for each trait, we wanted to fit individual models. For example,  
201 of the original 1,030 species detected from eBird we had body size data for 84% of species  
202 whereas for color we only had data for 44% of species. This approach allowed us to test both the  
203 independent relationships (i.e., each predictor separately against the response variable) and the  
204 relationship of a predictor given the other predictor variables (i.e., all predictors against the  
205 response variable simultaneously).

206

207 In all instances, the response variable was the residual from a log-log linear model fit between  
208 the eBird observations and the iNaturalist observations for a given species. In this instance, a  
209 species with a high (i.e., positive) residual would be over-represented in iNaturalist relative to  
210 eBird, whereas a species with a low (i.e., negative) residual would be under-represented in  
211 iNaturalist (Figure 1) relative to eBird. Each model fitted was stratified by state, accounting for  
212 differences in (1) the number of observers in a state, (2) the different relative abundance of a  
213 species throughout the United States, and (3) any other intrinsic differences that might exist  
214 among states that was not of inherent interest in our analysis. Table 1 summarizes the average  
215 sample size for the respective models fit among predictor variables. To confirm the robustness of  
216 our results at an individual state level, we ran a linear mixed effect model where the response  
217 variable was the residuals from a log-log linear model fit between the eBird observations and the  
218 iNaturalist observations for a given species, the predictor variables were the respective traits, and  
219 the random effect was state. Again, the models varied in sample size among predictor variables  
220 (see Table 1).

221

### 222 *Data analysis and availability*

223 All analyses were carried out in R software<sup>61</sup> and relied heavily on the tidyverse workflow<sup>62</sup>.  
224 Mixed-effects models were fitted using the lme4 package<sup>63</sup> and p-values were extracted using  
225 the lmerTest package<sup>64</sup>. Data and code to reproduce these analyses are available in a GitHub  
226 repository ([https://github.com/coreytcallaghan/inaturalist\\_preferences](https://github.com/coreytcallaghan/inaturalist_preferences)) and will be permanently  
227 archived in a Zenodo repository upon acceptance of this article.

228

## 229 RESULTS

230 A total of 507 species across the United States was included in our analysis. These species  
231 comprised a total of 255,727,592 eBird and 1,107,224 iNaturalist observations. At the state level,  
232 the number of eBird checklists and the number of iNaturalist observations were strongly  
233 correlated (Figure 2a;  $R^2 = 0.58$ ,  $p$ -value  $< 0.001$ ). Similarly, at the species level, the total  
234 number of iNaturalist observations and eBird observations for a given species was strongly  
235 correlated (Figure 2b;  $R^2 = 0.9$ ), and for both datasets the number of observations per species  
236 was positively-skewed (Figure S1). We also found that the percent of eBird checklists a species  
237 was found on and the percent of total iNaturalist observations a species comprised was strongly  
238 correlated among states (Figure S2), suggesting that species are sampled to a similar extent in  
239 opportunistic and semi-structured citizen science projects.

240  
241 Our analyses showed that larger species were more likely to be over-represented in the  
242 opportunistic citizen science dataset, with the residuals from the contrast between datasets  
243 strongly associated with body size (Figure 3, estimate = 0.11,  $t = 31.59$ ,  $p < 0.001$ ). We found no  
244 evidence that more colorful birds were over-represented in opportunistic citizen science data  
245 (estimate = -0.01,  $t = -0.413$ ,  $p = 0.68$ ) and moderate evidence that gregarious species were over-  
246 represented in opportunistic citizen science data (estimate=0.033,  $t=6.118$ ,  $p < 0.001$ ). There was  
247 some evidence that species which are of least concern (with IUCN status treated as an ordinal  
248 variable) were more commonly found in the opportunistic citizen science data (Figure 3; Figure  
249 S3, estimate = 0.078,  $t = 7.73$ ,  $p < 0.001$ ). The results from these individual linear models ran at  
250 the state level were confirmed by linear mixed models with state as a random effect. When  
251 considering all traits simultaneously in a linear mixed-effects model, the above patterns remained  
252 broadly similar (Figure 4).

253

## 254 DISCUSSION

255 We compared two popular citizen science platforms throughout the continental United States and  
256 found that there was strong agreement between the relative number of observations of a species  
257 in iNaturalist and eBird, albeit there were about 200 times more observations in eBird than  
258 iNaturalist. This suggests that species are observed at similar rates in both citizen science  
259 projects — i.e., the inherent processes driving observation in both opportunistic and semi-  
260 structured citizen science projects are similar. Nevertheless, in support of our predictions (Figure  
261 1) we found strong evidence that large-bodied birds are over-represented in the opportunistic  
262 citizen science dataset compared with the semi-structured dataset. We also found moderate  
263 evidence that common species were over-represented in the opportunistic data, and weak  
264 evidence that species in large flocks were over-represented. In contrast to our prediction,  
265 however, we found no evidence that brightly-colored species were over-represented in  
266 opportunistic citizen science data.

267

268 Our finding that large-bodied birds were over-represented in an opportunistic citizen science  
269 dataset is probably because larger-bodied birds are more detectable<sup>53, 65</sup>. Thus, smaller-bodied  
270 taxa are under-represented in citizen science data<sup>66, 67, 68</sup>, but this may not be the case for other  
271 taxa such as mammals<sup>69</sup>. However, it is difficult to know whether this is an inherent preference  
272 shown by users of the opportunistic citizen science data, or if this comes about as part of the  
273 recording process (e.g., species' detectability;<sup>50</sup>). Species detectability is complex and can be  
274 linked to a species' mobility or habitat preferences of the species themselves; for example, large-  
275 bodied wading birds generally occurring in open wetlands are more easily detected than small-

276 bodied songbirds generally occurring in dense forest understory. For amphibians and reptiles,  
277 climatic niches are not fully sampled by citizen science datasets due in part to life history and  
278 habitat sampling biases <sup>29</sup>. Moreover, in order for an observer to make a record in iNaturalist,  
279 usually a photo is uploaded (although sound recordings are also accepted). Because a photo is  
280 needed, the detectability process is two-fold — first, it needs to be detected, and second, it needs  
281 to be photographed, which is likely easier for many large-bodied birds. Longer lenses, often  
282 restricted to serious photographers, may be needed to photograph smaller-bodied birds whereas  
283 smartphones can usually capture a sufficient image of a larger-bodied bird. The bias towards  
284 large-bodied birds in the opportunistic data is probably a result of detectability and the ability to  
285 capture a photograph <sup>53</sup>. This process is similar in insects, for example, which are generally  
286 small, but larger insects (e.g., butterflies) are both easier to observe, photograph, and identify —  
287 making it likely that the biases we found in birds generalize to insects as well. Indeed, a study of  
288 bugs and beetles found that smaller species are typically less represented in citizen science data  
289 <sup>68</sup>. Importantly, because this represents a form of systematic bias, it is likely easier to model this  
290 bias as we know that this data is not missing at random (e.g., <sup>70</sup>) and thus body size should be  
291 included in various modelling processes when using opportunistic citizen science data (e.g., <sup>67</sup>).  
292

293 Similar to body size, we found that birds which occur in larger groups (i.e., flocks) and those that  
294 are of least concern are over-represented in the opportunistic dataset. This, again, may be  
295 inherently linked to the recording process, rather than a specific bias or preference of the  
296 observers themselves. This is because common birds, that occur in large flocks, are more likely  
297 to be seen and thus submitted to the opportunistic citizen science data <sup>65</sup>. A larger flock will  
298 likely also provide more opportunities to capture a photograph than when observing a single

299 individual, as has been illustrated in the detectability of animals from aerial surveys by  
300 professionals 71. One explanation for the least concern birds being over-represented in  
301 iNaturalist is user behavior — eBird data are more likely to be derived from avid birdwatchers  
302 (e.g., those that search out uncommon birds and keep serious lists) compared with iNaturalist  
303 data which may be derived from more recreational birdwatchers that focus on ‘backyard’  
304 species. Another important distinction between iNaturalist and eBird is how identifications are  
305 made. In eBird, most identifications are made acoustically, whereas a photo is generally required  
306 for iNaturalist. Most traits analyzed here are related to visual encounter/identification, thus  
307 potentially explaining the differences found between the opportunistic iNaturalist and the semi-  
308 structured eBird data.

309

310 The lack of signal of the colorfulness of a species in predicting over-representation in iNaturalist  
311 could suggest that iNaturalist users are not limited by ‘attractiveness/aesthetics’ but mostly by  
312 detectability, as discussed above (Figure 4). Quantifying the influence of color on detectability  
313 remains a challenge (e.g., <sup>72</sup>). In contrast to our results, <sup>68</sup> found that more colorful insect species  
314 are more commonly reported, as well as more patterned and morphologically interesting species.  
315 This may suggest, at least in the case of insects, that contributors are selecting subjects based on  
316 their visual aesthetics, not just their detectability. The discrepancies between our results and that  
317 of <sup>68</sup> suggest that the influence of traits may vary between different taxa, making it important to  
318 explore these relationships for a range of organisms rather than extrapolating the results of birds,  
319 or bugs and beetles, to other groups.

320

321 While citizen science data are undoubtedly valuable for ecology and conservation <sup>4, 73, 74</sup>, there  
322 remain limits to the use of citizen science datasets <sup>13, 75</sup>. The ability to sample remote regions, for  
323 example, will likely remain a limitation in citizen science data, and this has been well-recognized  
324 <sup>17</sup>. Quantifying the limits of citizen science datasets for use in ecology and conservation remains  
325 an important step for the future widespread use of citizen science data in ecology and  
326 conservation. Data-integration — where noisy citizen science data are integrated with  
327 professionally-curated datasets — will likely be increasingly important in the future use of  
328 citizen science data <sup>76, 77</sup>. By knowing the biases present in citizen science data, experts can  
329 preferentially generate data that maximize the integration process, for example by collecting data  
330 from remote regions. Further, professional scientists should use limited funding to target species  
331 that are likely to be under-represented in citizen science datasets — i.e., rare, small-bodied,  
332 species.

333  
334 Ultimately, citizen science data will continue to perform, at least in part, a substantial role in the  
335 future of ecology and conservation research <sup>44</sup>. Understanding, documenting, and quantifying the  
336 biases associated with these data remains an important first step before the widespread use of  
337 these data in answering ecological questions and biodiversity monitoring <sup>5</sup>. Our results highlight  
338 that for birds, semi-structured eBird out-samples opportunistic iNaturalist data, but the number of  
339 observations recorded per species are strongly correlated between the two platforms. When  
340 looking at the differences in this relationship, it is clear that biases exist, likely due to the biases  
341 in the opportunistic iNaturalist data. We note that we compared the opportunistic dataset to a  
342 semi-structured dataset, and the semi-structured dataset does not necessarily represent the  
343 “truth”. The biases found here, could also be present when comparing a semi-structured dataset

344 to true density or abundance of birds in the landscape. To better understand these differences,  
345 future research in this space should continue to focus on quantifying and documenting biases in  
346 citizen science data, and understanding how these biases change from unstructured to semi-  
347 structured to structured citizen science platforms. Nevertheless, our results demonstrate the  
348 importance of using species-specific traits, when modelling citizen science datasets<sup>27, 29, 52, 78, 79,</sup>  
349 <sup>80</sup>.

350

### 351 ACKNOWLEDGEMENTS

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

357

### 358 AUTHOR CONTRIBUTIONS

359 CTC conceived and led the study with input from all authors. CTC performed the analyses with  
360 input from all authors. CTC wrote the first version of the manuscript and all authors contributed  
361 to editing the manuscript.

362

### 363 COMPETING INTERESTS

364 The author(s) declare no competing interests.

## 365 REFERENCES

- 366 1. Pocock, M. J., Tweddle, J. C., Savage, J., Robinson, L. D. & Roy, H. E. The diversity and  
367 evolution of ecological and environmental citizen science. *PLoS One* **12**, e0172579  
368 (2017).
- 369 2. Chandler, M. *et al.* Contribution of citizen science towards international biodiversity  
370 monitoring. *Biological Conservation* **213**, 280–294 (2017).
- 371 3. Chandler, M. *et al.* Involving citizen scientists in biodiversity observation. in *The GEO*  
372 *handbook on biodiversity observation networks* 211–237 (Springer, Cham, 2017).
- 373 4. McKinley, D. C. *et al.* Citizen science can improve conservation science, natural resource  
374 management, and environmental protection. *Biological Conservation* **208**, 15–28 (2017).
- 375 5. Pereira, H. M. *et al.* Monitoring essential biodiversity variables at the species level. in  
376 *The GEO handbook on biodiversity observation networks* 79–105 (Springer, Cham,  
377 2017).
- 378 6. Wiggins, A. & Crowston, K. From conservation to crowdsourcing: A typology of citizen  
379 science. in *2011 44th hawaii international conference on system sciences* 1–10 (IEEE,  
380 2011).
- 381 7. Haklay, M. Citizen science and volunteered geographic information: Overview and  
382 typology of participation. in *Crowdsourcing geographic knowledge* 105–122 (Springer,  
383 2013).
- 384 8. Kelling, S. *et al.* Using semistructured surveys to improve citizen science data for  
385 monitoring biodiversity. *BioScience* **69**, 170–179 (2019).

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

- 408 18. Geldmann, J. *et al.* What determines spatial bias in citizen science? Exploring four  
409 recording schemes with different proficiency requirements. *Diversity and Distributions*  
410 **22**, 1139–1149 (2016).
- 411 19. Callaghan, C. T. *et al.* Three frontiers for the future of biodiversity research using citizen  
412 science data. *BioScience* **71**, 55–63 (2021).
- 413 20. Ward, D. F. Understanding sampling and taxonomic biases recorded by citizen scientists.  
414 *Journal of insect conservation* **18**, 753–756 (2014).
- 415 21. Troudet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R. & Legendre, F. Taxonomic bias in  
416 biodiversity data and societal preferences. *Scientific reports* **7**, 1–14 (2017).
- 417 22. Martín-López, B., Montes, C., Ramírez, L. & Benayas, J. What drives policy decision-  
418 making related to species conservation? *Biological Conservation* **142**, 1370–1380 (2009).
- 419 23. Boakes, E. H. *et al.* Distorted views of biodiversity: Spatial and temporal bias in species  
420 occurrence data. *PLoS Biol* **8**, e1000385 (2010).
- 421 24. Aceves-Bueno, E. *et al.* The accuracy of citizen science data: A quantitative review.  
422 *Bulletin of the Ecological Society of America* **98**, 278–290 (2017).
- 423 25. Davies, T. K., Stevens, G., Meekan, M. G., Struve, J. & Rowcliffe, J. M. Can citizen  
424 science monitor whale-shark aggregations? Investigating bias in mark–recapture  
425 modelling using identification photographs sourced from the public. *Wildlife Research*  
426 **39**, 696–704 (2013).
- 427 26. Crall, A. W. *et al.* Assessing citizen science data quality: An invasive species case study.  
428 *Conservation Letters* **4**, 433–442 (2011).

- 429 27. Strien, A. J. van, Swaay, C. A. van & Termaat, T. Opportunistic citizen science data of  
430 animal species produce reliable estimates of distribution trends if analysed with  
431 occupancy models. *Journal of Applied Ecology* **50**, 1450–1458 (2013).
- 432 28. Johnston, A., Moran, N., Musgrove, A., Fink, D. & Baillie, S. R. Estimating species  
433 distributions from spatially biased citizen science data. *Ecological Modelling* **422**,  
434 108927 (2020).
- 435 29. Tiago, P., Pereira, H. M. & Capinha, C. Using citizen science data to estimate climatic  
436 niches and species distributions. *Basic and Applied Ecology* **20**, 75–85 (2017).
- 437 30. Sullivan, B. L. *et al.* Using open access observational data for conservation action: A case  
438 study for birds. *Biological Conservation* **208**, 5–14 (2017).
- 439 31. Callaghan, C. T. *et al.* Citizen science data accurately predicts expert-derived species  
440 richness at a continental scale when sampling thresholds are met. *Biodiversity and*  
441 *Conservation* **29**, 1323–1337 (2020).
- 442 32. Birkin, L. & Goulson, D. Using citizen science to monitor pollination services.  
443 *Ecological Entomology* **40**, 3–11 (2015).
- 444 33. Delaney, D. G., Sperling, C. D., Adams, C. S. & Leung, B. Marine invasive species:  
445 Validation of citizen science and implications for national monitoring networks.  
446 *Biological Invasions* **10**, 117–128 (2008).
- 447 34. Schultz, C. B., Brown, L. M., Pelton, E. & Crone, E. E. Citizen science monitoring  
448 demonstrates dramatic declines of monarch butterflies in western north america.  
449 *Biological Conservation* **214**, 343–346 (2017).
- 450 35. Bird, T. J. *et al.* Statistical solutions for error and bias in global citizen science datasets.  
451 *Biological Conservation* **173**, 144–154 (2014).

- 452 36. Isaac, N. J., Strien, A. J. van, August, T. A., Zeeuw, M. P. de & Roy, D. B. Statistics for  
453 citizen science: Extracting signals of change from noisy ecological data. *Methods in*  
454 *Ecology and Evolution* **5**, 1052–1060 (2014).
- 455 37. Dickinson, J. L. *et al.* The current state of citizen science as a tool for ecological research  
456 and public engagement. *Frontiers in Ecology and the Environment* **10**, 291–297 (2012).
- 457 38. Bonney, R. *et al.* Next steps for citizen science. *Science* **343**, 1436–1437 (2014).
- 458 39. Jordan, R. C., Gray, S. A., Howe, D. V., Brooks, W. R. & Ehrenfeld, J. G. Knowledge  
459 gain and behavioral change in citizen-science programs. *Conservation biology* **25**, 1148–  
460 1154 (2011).
- 461 40. Crall, A. W. *et al.* The impacts of an invasive species citizen science training program on  
462 participant attitudes, behavior, and science literacy. *Public Understanding of Science* **22**,  
463 745–764 (2013).
- 464 41. Jordan, R. C., Ballard, H. L. & Phillips, T. B. Key issues and new approaches for  
465 evaluating citizen-science learning outcomes. *Frontiers in Ecology and the Environment*  
466 **10**, 307–309 (2012).
- 467 42. Evans, C. *et al.* The neighborhood nestwatch program: Participant outcomes of a citizen-  
468 science ecological research project. *Conservation Biology* **19**, 589–594 (2005).
- 469 43. Haywood, B. K., Parrish, J. K. & Dolliver, J. Place-based and data-rich citizen science as  
470 a precursor for conservation action. *Conservation Biology* **30**, 476–486 (2016).
- 471 44. Pocock, M. J. *et al.* A vision for global biodiversity monitoring with citizen science. in  
472 *Advances in ecological research* vol. 59 169–223 (Elsevier, 2018).

- 473 45. Tiago, P., Gouveia, M. J., Capinha, C., Santos-Reis, M. & Pereira, H. M. The influence  
474 of motivational factors on the frequency of participation in citizen science activities.  
475 *Nature Conservation* **18**, 61 (2017).
- 476 46. Isaac, N. J. & Pocock, M. J. Bias and information in biological records. *Biological*  
477 *Journal of the Linnean Society* **115**, 522–531 (2015).
- 478 47. Angulo, E. & Courchamp, F. Rare species are valued big time. *PloS one* **4**, e5215 (2009).
- 479 48. Booth, J. E., Gaston, K. J., Evans, K. L. & Armsworth, P. R. The value of species rarity  
480 in biodiversity recreation: A birdwatching example. *Biological Conservation* **144**, 2728–  
481 2732 (2011).
- 482 49. Rowley, J. J. *et al.* FrogID: Citizen scientists provide validated biodiversity data on frogs  
483 of australia. *Herpetological Conservation and Biology* **14**, 155–70 (2019).
- 484 50. Boakes, E. H. *et al.* Patterns of contribution to citizen science biodiversity projects  
485 increase understanding of volunteers recording behaviour. *Scientific reports* **6**, 33051  
486 (2016).
- 487 51. Garrard, G. E., McCarthy, M. A., Williams, N. S., Bekessy, S. A. & Wintle, B. A. A  
488 general model of detectability using species traits. *Methods in Ecology and Evolution* **4**,  
489 45–52 (2013).
- 490 52. Denis, T. *et al.* Biological traits, rather than environment, shape detection curves of large  
491 vertebrates in neotropical rainforests. *Ecological Applications* **27**, 1564–1577 (2017).
- 492 53. Sólymos, P., Matsuoka, S. M., Stralberg, D., Barker, N. K. & Bayne, E. M. Phylogeny  
493 and species traits predict bird detectability. *Ecography* **41**, 1595–1603 (2018).
- 494 54. Wood, C., Sullivan, B., Iliff, M., Fink, D. & Kelling, S. eBird: Engaging birders in  
495 science and conservation. *PLoS Biol* **9**, e1001220 (2011).

- 496 55. GBIF.org (3<sup>rd</sup> December 2019). GBIF occurrence download.  
497 <https://www.doi.org/10.15468/dl.lpwczr>
- 498 56. Gilfedder, M. *et al.* Brokering trust in citizen science. *Society & natural resources* **32**,  
499 292–302 (2019).
- 500 57. Callaghan, C., Lyons, M., Martin, J., Major, R. & Kingsford, R. Assessing the reliability  
501 of avian biodiversity measures of urban greenspaces using eBird citizen science data.  
502 *Avian Conservation and Ecology* **12**, (2017).
- 503 58. Johnston, A. *et al.* Best practices for making reliable inferences from citizen science data:  
504 Case study using eBird to estimate species distributions. *BioRxiv* 574392 (2019).
- 505 59. Myhrvold, N. P. *et al.* An amniote life-history database to perform comparative analyses  
506 with birds, mammals, and reptiles: Ecological archives E096-269. *Ecology* **96**, 3109–  
507 3109 (2015).
- 508 60. Dale, J., Dey, C. J., Delhey, K., Kempnaers, B. & Valcu, M. The effects of life history  
509 and sexual selection on male and female plumage colouration. *Nature* **527**, 367–370  
510 (2015).
- 511 61. R Core Team. *R: A language and environment for statistical computing*. (R Foundation  
512 for Statistical Computing, 2020).
- 513 62. Wickham, H. *et al.* Welcome to the tidyverse. *Journal of Open Source Software* **4**, 1686  
514 (2019).
- 515 63. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting linear mixed-effects models  
516 using lme4. *Journal of Statistical Software* **67**, 1–48 (2015).
- 517 64. Kuznetsova, A., Brockhoff, P. B. & Christensen, R. H. B. lmerTest package: Tests in  
518 linear mixed effects models. *Journal of Statistical Software* **82**, 1–26 (2017).

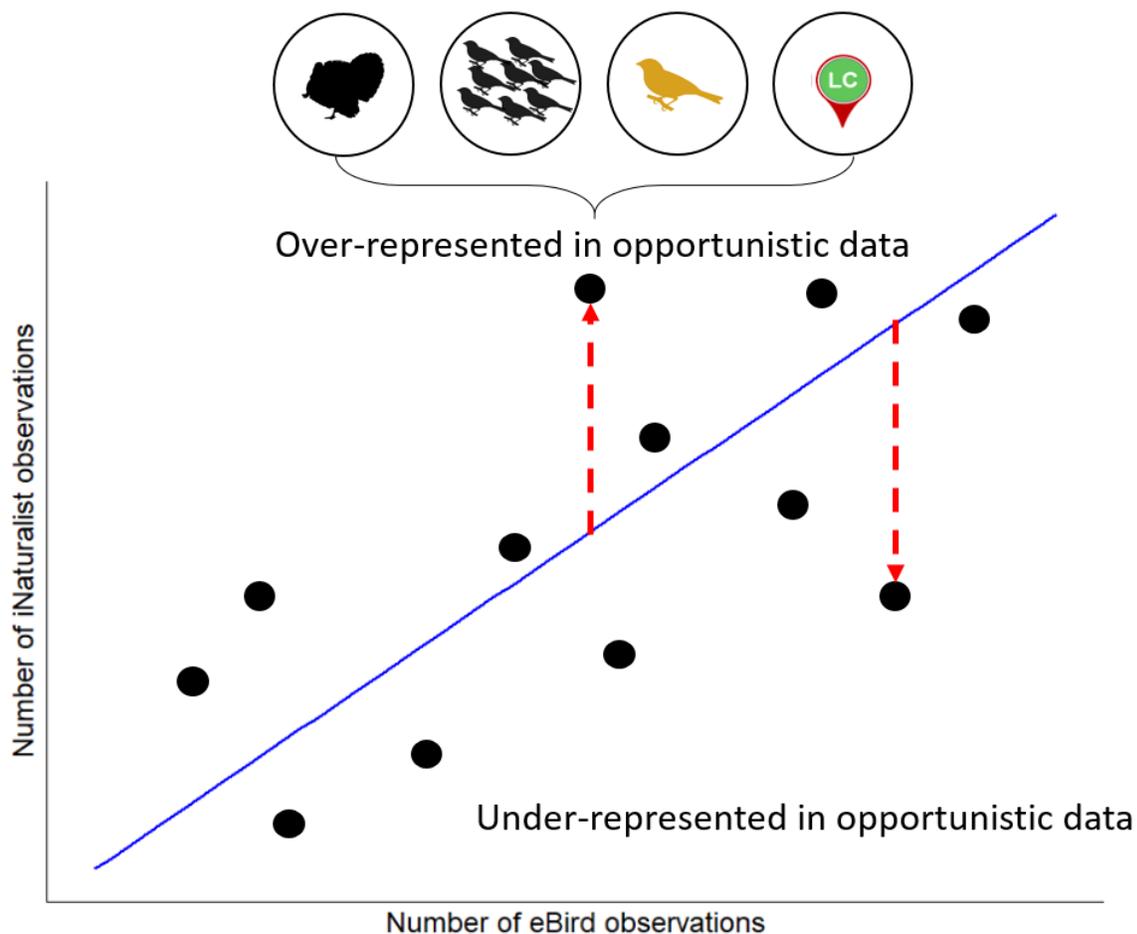
- 519 65. Johnston, A. *et al.* Species traits explain variation in detectability of UK birds. *Bird Study*  
520 **61**, 340–350 (2014).
- 521 66. Steen, V. A., Elphick, C. S. & Tingley, M. W. An evaluation of stringent filtering to  
522 improve species distribution models from citizen science data. *Diversity and*  
523 *Distributions* **25**, 1857–1869 (2019).
- 524 67. Henckel, L., Bradter, U., Jönsson, M., Isaac, N. J. & Snäll, T. Assessing the usefulness of  
525 citizen science data for habitat suitability modelling: Opportunistic reporting versus  
526 sampling based on a systematic protocol. *Diversity and Distributions* **26**, 1276–1290  
527 (2020).
- 528 68. Caley, P., Welvaert, M. & Barry, S. C. Crowd surveillance: Estimating citizen science  
529 reporting probabilities for insects of biosecurity concern. *Journal of Pest Science* **93**,  
530 543–550 (2020).
- 531 69. Périquet, S., Roxburgh, L., Roux, A. le & Collinson, W. J. Testing the value of citizen  
532 science for roadkill studies: A case study from south africa. *Frontiers in Ecology and*  
533 *Evolution* **6**, 15 (2018).
- 534 70. Nakagawa, S. & Freckleton, R. P. Model averaging, missing data and multiple  
535 imputation: A case study for behavioural ecology. *Behavioral Ecology and Sociobiology*  
536 **65**, 103–116 (2011).
- 537 71. Schlossberg, S., Chase, M. & Griffin, C. Using species traits to predict detectability of  
538 animals on aerial surveys. *Ecological Applications* **28**, 106–118 (2018).
- 539 72. Troscianko, J., Skelhorn, J. & Stevens, M. Quantifying camouflage: How to predict  
540 detectability from appearance. *BMC evolutionary biology* **17**, 1–13 (2017).

- 541 73. Tulloch, A. I., Possingham, H. P., Joseph, L. N., Szabo, J. & Martin, T. G. Realising the  
542 full potential of citizen science monitoring programs. *Biological Conservation* **165**, 128–  
543 138 (2013).
- 544 74. Kobori, H. *et al.* Citizen science: A new approach to advance ecology, education, and  
545 conservation. *Ecological research* **31**, 1–19 (2016).
- 546 75. Callaghan, C. T., Poore, A. G., Major, R. E., Rowley, J. J. & Cornwell, W. K. Optimizing  
547 future biodiversity sampling by citizen scientists. *Proceedings of the Royal Society B* **286**,  
548 20191487 (2019).
- 549 76. Pacifici, K. *et al.* Integrating multiple data sources in species distribution modeling: A  
550 framework for data fusion. *Ecology* **98**, 840–850 (2017).
- 551 77. Robinson, O. J. *et al.* Integrating citizen science data with expert surveys increases  
552 accuracy and spatial extent of species distribution models. *Diversity and Distributions* **26**,  
553 976–986 (2020).
- 554 78. Strien, A. J. van, Termaat, T., Groenendijk, D., Mensing, V. & Kery, M. Site-occupancy  
555 models may offer new opportunities for dragonfly monitoring based on daily species  
556 lists. *Basic and Applied Ecology* **11**, 495–503 (2010).
- 557 79. Van der Wal, R. *et al.* Mapping species distributions: A comparison of skilled naturalist  
558 and lay citizen science recording. *Ambio* **44**, 584–600 (2015).
- 559 80. Dennis, E. B., Morgan, B. J., Brereton, T. M., Roy, D. B. & Fox, R. Using citizen science  
560 butterfly counts to predict species population trends. *Conservation biology* **31**, 1350–  
561 1361 (2017).

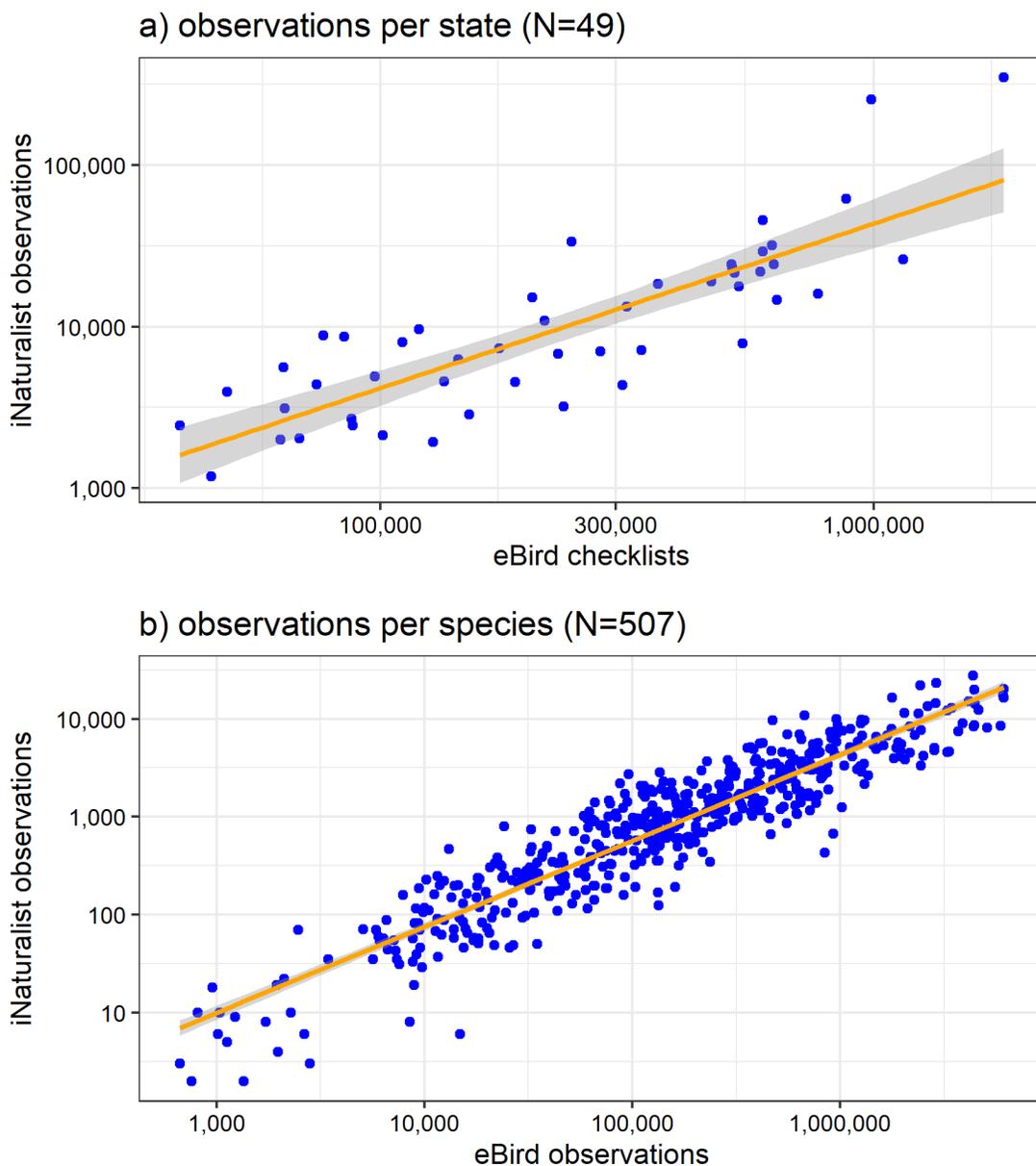
562

563

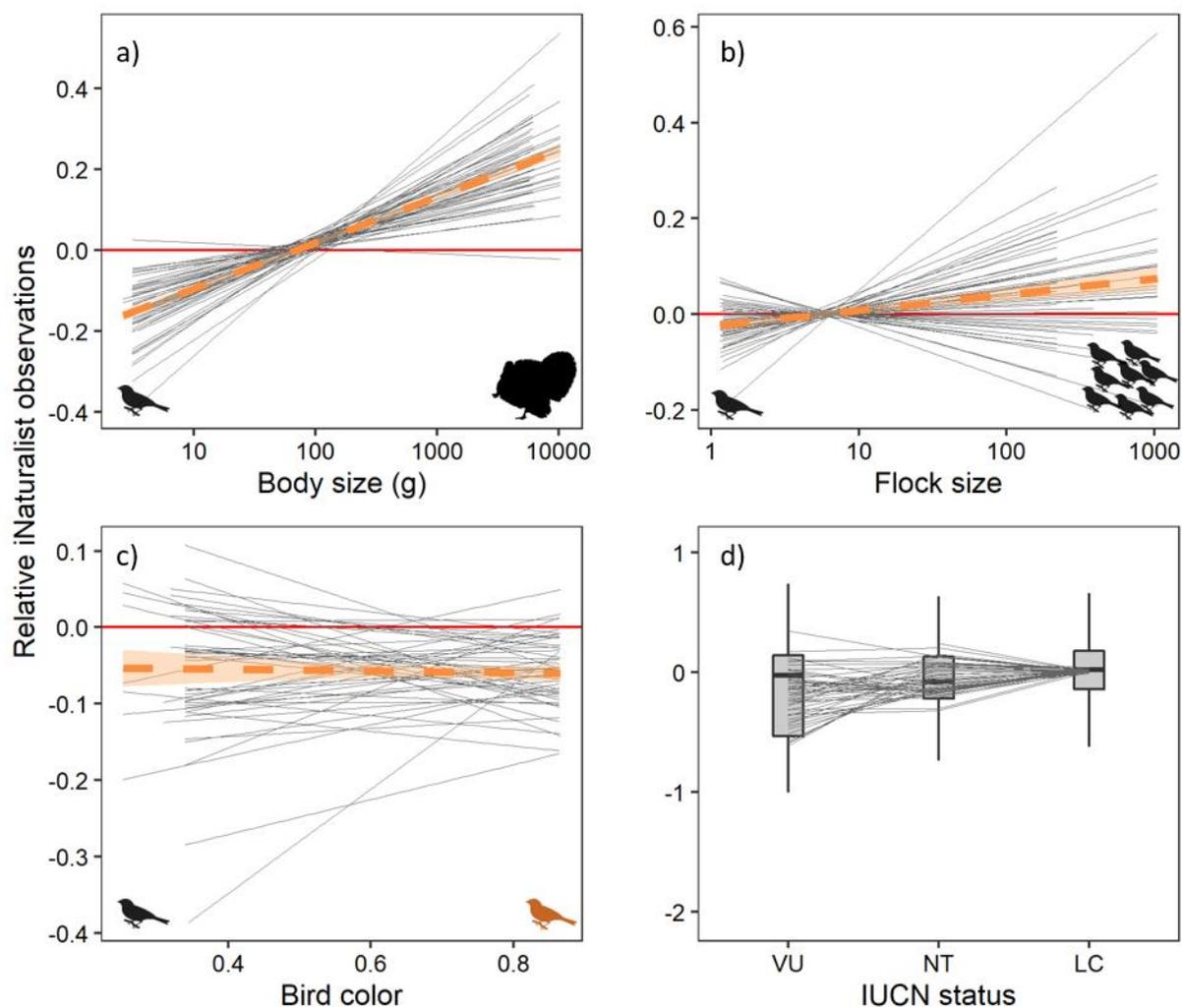
564 FIGURES  
565



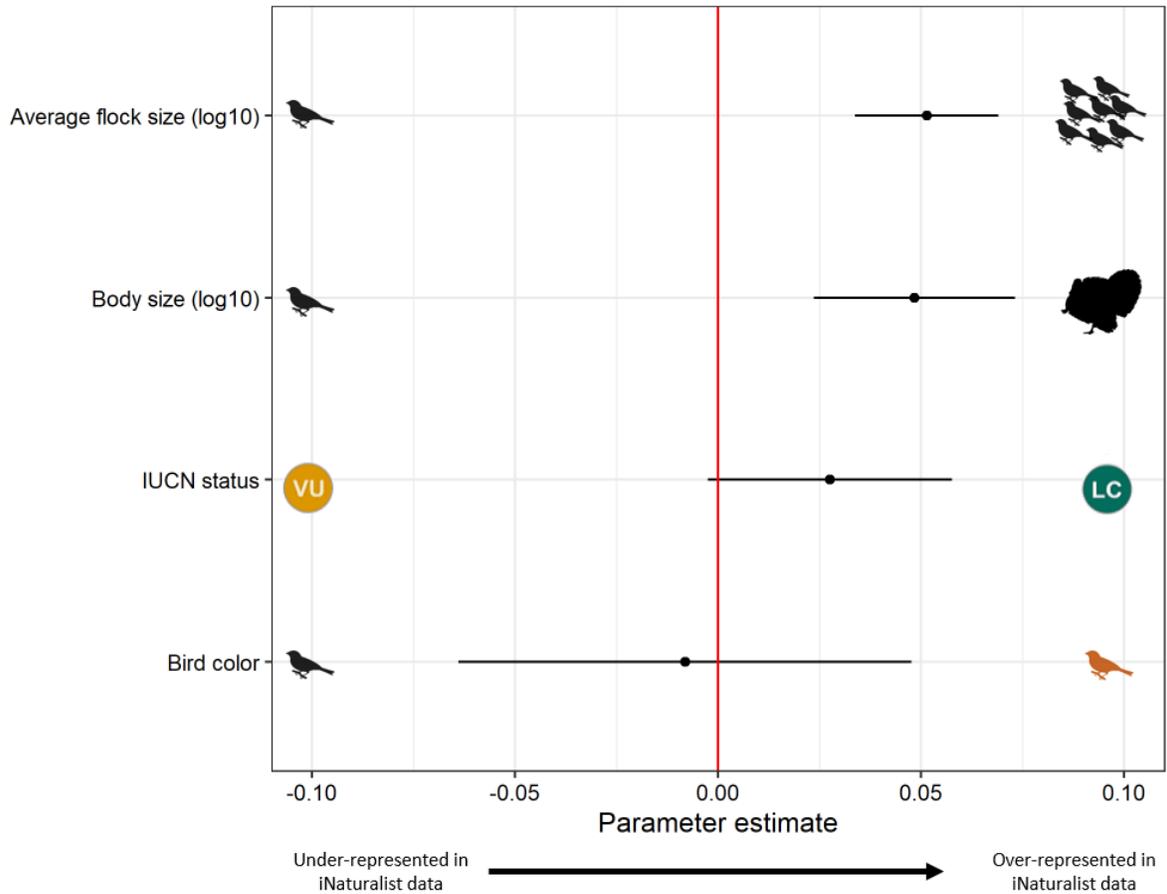
566  
567 **Figure 1.** A conceptual figure depicting the methods used in our analysis. We used the residual  
568 from the relationship between the number of eBird observations (i.e., semi-structured citizen  
569 science observations) and iNaturalist observations (i.e., opportunistic citizen science  
570 observations) to quantify the over- or under-representation of a species in opportunistic citizen  
571 science data. We predicted that species which were over-represented in opportunistic iNaturalist  
572 data would be larger in size, occur more frequently in large flocks, be brighter in color, and be  
573 categorized as Least Concern IUCN status (a proxy for commonness).



574  
 575 **Figure 2.** a) The relationship between the total number of eBird checklists and total number of  
 576 iNaturalist observations for 49 states, including the District of Columbia. There was strong  
 577 evidence that these variables were correlated ( $R^2=0.58$ ,  $p$ -value  $<0.001$ ) suggesting that sampling  
 578 between datasets is correlated among states. b) The relationship between the number of  
 579 observations for a species from eBird (x-axis) and the number of observations for a species from  
 580 iNaturalist (y-axis) for only eBird species which were found on  $>1\%$  of eBird checklists.



581  
 582 **Figure 3.** The relationship between a) body size of a species, b) flock size, c) color and d)  
 583 commonness and the residuals of a linear model fit between iNaturalist and eBird observations  
 584 (see Figure 1). These results demonstrate that there is a strong bias of body size in iNaturalist  
 585 compared with eBird. Positive values on the y-axis mean over-represented in iNaturalist and  
 586 negative values on the y-axis mean under-represented in iNaturalist. Body size and flock size are  
 587 represented on a log<sub>10</sub> scale. Each line represents a state (N=49). For a-c), the overall  
 588 relationship pooling states is represented by the orange fitted line and 95% confidence interval.  
 589



590  
591  
592  
593  
594  
595  
596

**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.

597 TABLES

598 **Table 1.** A summary of the average number of observations in a model among states and the  
 599 standard deviation of the number of observations in a model. The N for the mixed effects models  
 600 represents the total number of observations in each model.

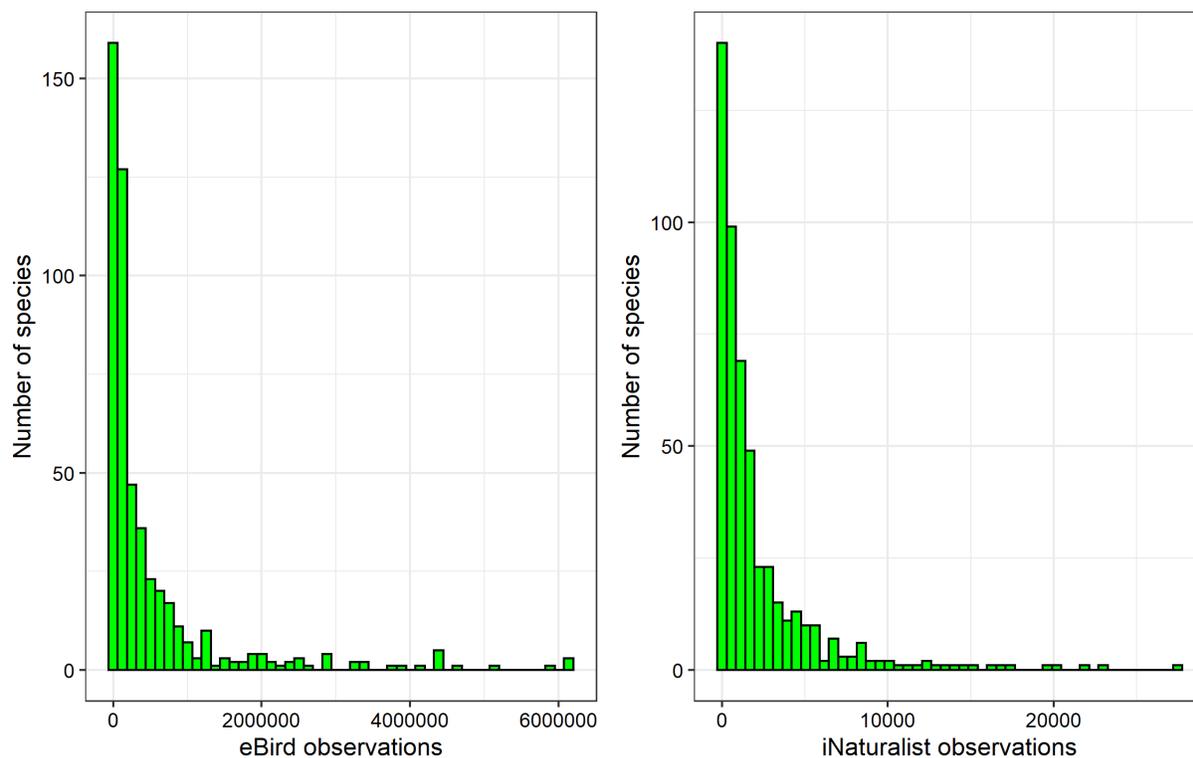
	State-specific models (N=49)		Mixed effects model
	Mean number of obs	SD of obs	Number of obs
Body size	158.02	18.11	7743
Color	92.69	10.62	4542
Flock size	177.59	21.44	8702
IUCN status	155.76	17.78	7629
All variables			3986

601

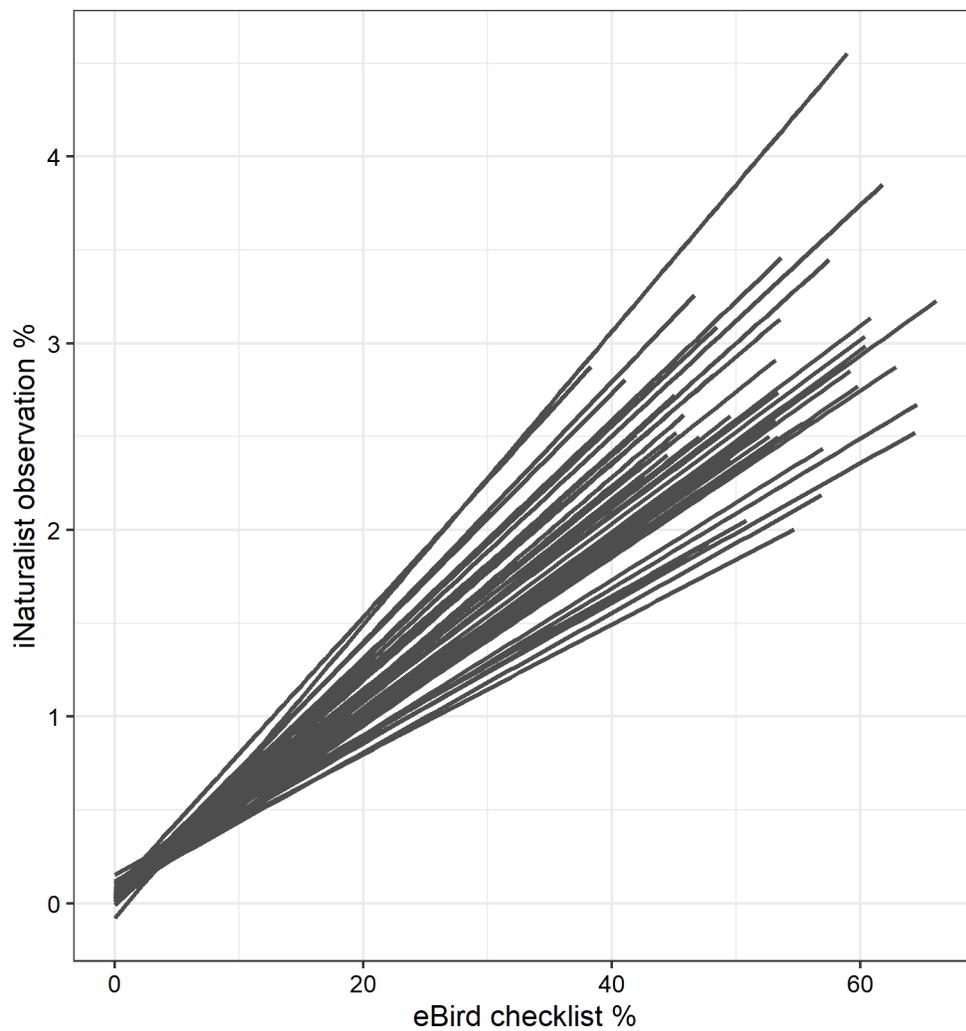
602

603

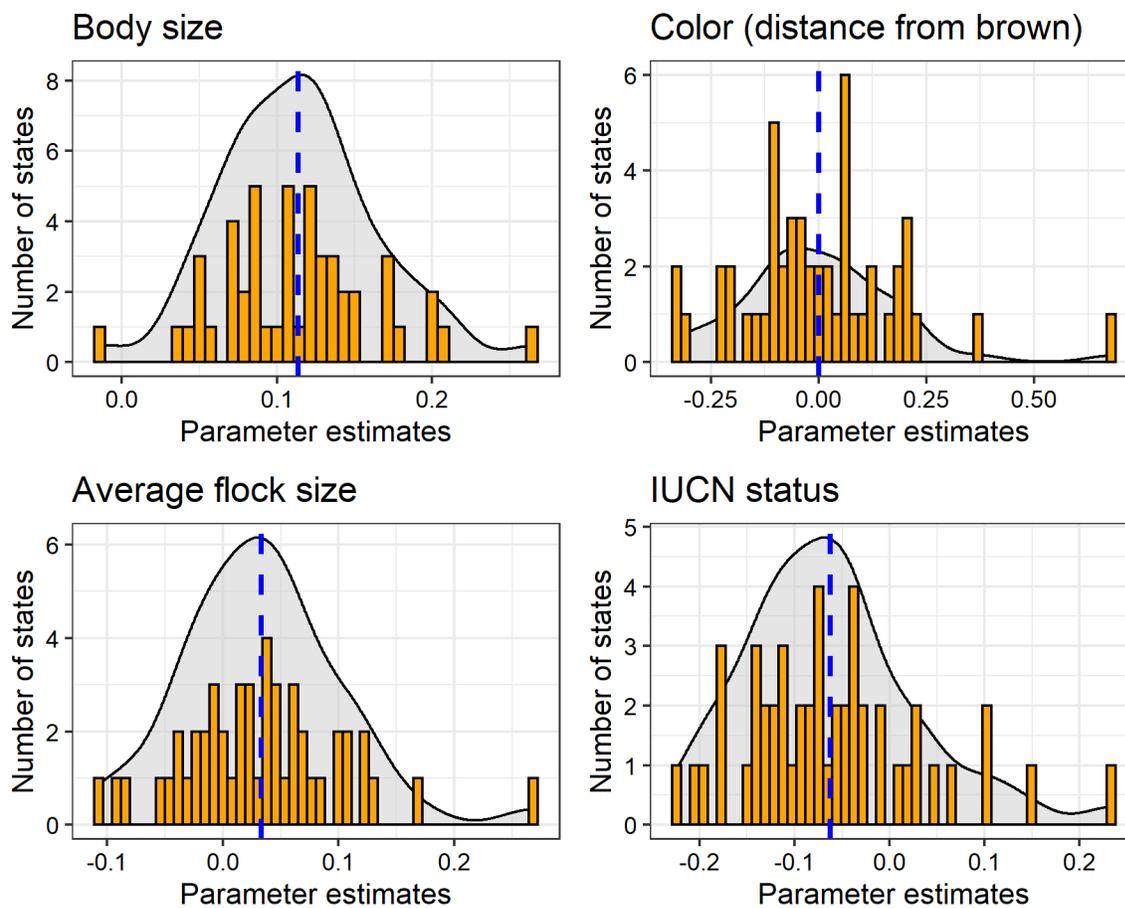
## SUPPLEMENTARY FIGURES



**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.



**Figure S3.** Distributions of parameter estimates of our individual models for our four predictor variables of interest. The x-axis represents the parameter estimate for a linear model between the residuals and the associated predictor variables, and the y-axis represents the number of states (i.e., models) which are associated with that histogram bin.