



OPEN

## 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 J. Roberts<sup>2</sup> & Henrique M. Pereira<sup>1,3</sup>

Citizen science platforms are quickly accumulating hundreds of millions of biodiversity observations around the world annually. Quantifying and correcting for the biases in citizen science datasets remains an important first step before these data are used to address ecological questions and monitor biodiversity. One source of potential bias among datasets is the difference between those citizen science programs that have unstructured protocols and those that have semi-structured or structured protocols for submitting observations. To quantify biases in an unstructured citizen science platform, we contrasted bird observations from the unstructured iNaturalist platform with that from a semi-structured citizen science platform—eBird—for the continental United States. We tested whether four traits of species (body size, commonness, flock size, and color) predicted if a species was under- or over-represented in the unstructured dataset compared with the semi-structured dataset. We found strong evidence that large-bodied birds were over-represented in the unstructured citizen science dataset; moderate evidence that common species were over-represented in the unstructured dataset; strong evidence that species in large groups were over-represented; and no evidence that colorful species were over-represented in unstructured citizen science data. Our results suggest that biases exist in unstructured citizen science data when compared with semi-structured data, likely as a result of the detectability of a species and the inherent recording process. Importantly, in programs like iNaturalist the detectability process is two-fold—first, an individual organism needs to be detected, and second, it needs to be photographed, which is likely easier for many large-bodied species. Our results indicate that caution is warranted when using unstructured citizen science data in ecological modelling, and highlight body size as a fundamental trait that can be used as a covariate for modelling opportunistic species occurrence records, representing the detectability or identifiability in unstructured citizen science datasets. Future research in this space should continue to focus on quantifying and documenting biases in citizen science data, and expand our research by including structured citizen science data to understand how biases differ among unstructured, semi-structured, and structured citizen science platforms.

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

Data quality from citizen science projects has been questioned<sup>11,12</sup>, and such concerns can act as a barrier to the widespread use of citizen science data in ecology and conservation<sup>13</sup>. These concerns arise because citizen science data can be biased temporally, spatially, and/or taxonomically. Temporally, many citizen science datasets

<sup>1</sup>German Centre for Integrative Biodiversity Research (iDiv) Halle-Jena-Leipzig, Puschstr. 4, 04103 Leipzig, Germany. <sup>2</sup>Ecology and Evolution Research Centre, School of Biological, Earth and Environmental Sciences, UNSW Sydney, Sydney, NSW, Australia. <sup>3</sup>Institute of Biology, Martin Luther University Halle-Wittenberg, Am Kirchtor 1, 06108 Halle (Saale), Germany. ✉email: corey.callaghan@idiv.de

are biased because participants frequently sample on weekends<sup>14</sup> or disproportionately during specific times of the year such as spring migration for birds<sup>15</sup>, or during specific times of day, such as the morning period when birds are most active. Spatially, there is often a disproportionate number of sightings from areas with large human populations<sup>16</sup>, areas with more accessibility<sup>17</sup>, regions with high biodiversity that attract observers<sup>18</sup>, and regions of the world with higher socioeconomic status<sup>19</sup>. Taxonomic biases also exist as some taxa receive orders of magnitude more citizen science observations than other taxa, evidenced by the fact that birds represent a disproportionate amount of data in the Global Biodiversity Information Facility<sup>2</sup>. Even within citizen science projects focused on specific taxa, there can be strong taxonomic biases towards particularly charismatic species or those that are readily identified<sup>20–23</sup>.

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

In addition to being an excellent resource for scientists to better understand ecological questions, citizen science projects can encourage increased engagement of the general public with science<sup>37,38</sup>. Many citizen science projects provide learning opportunities for their volunteers. For example, participants in citizen science projects have increased their knowledge about invasive weeds<sup>39–41</sup>, increased their knowledge of bird biology and behavior<sup>42</sup>, and even enhanced their conservation awareness and sense of place<sup>42,43</sup>. The ecological advances derived from citizen science data, combined with the important role it plays in community engagement with science, suggests that citizen science data will continue to play an important role in ecological and conservation research in the future<sup>2,4,38,44</sup>. However, what motivates volunteers to participate in science, and contribute observations, has important implications for the quality of the data obtained<sup>45</sup>, particularly if there are biases towards certain species, places, or times of sampling.

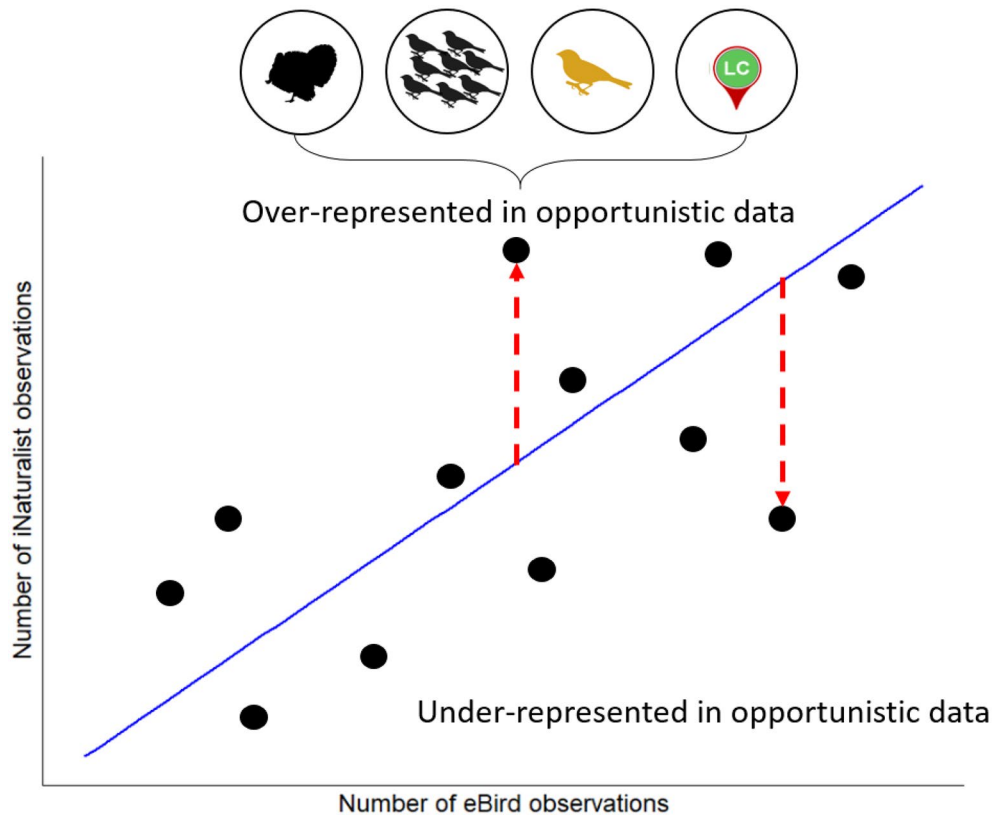
To ensure the continued and expanded use of citizen science data in ecology and conservation, it is important to document and understand the different biases present in citizen science datasets. Importantly, the degree of bias in a particular dataset will be influenced by the level of structure of that citizen science project. For example, unstructured projects (e.g., iNaturalist, [www.inaturalist.org](http://www.inaturalist.org)) or semi-structured projects (e.g., eBird, [www.ebird.org](http://www.ebird.org)) will generally be more spatially biased than structured projects that have pre-defined spatial sampling locations (e.g., Breeding Bird Surveys). Or, a citizen science project that collects incidental presence-only data, such as iNaturalist, is likely more susceptible to individual observer preferences compared with a semi-structured or structured project that requires all species encountered to be recorded by the observers. Charismatic species<sup>21</sup> can be over-represented in citizen science data because observers are more likely to record species that they, or society, consider more interesting or relevant for monitoring<sup>46</sup>. Similarly, rare species are more likely to be the subject of conservation monitoring or more likely to be actively searched for by amateur naturalists<sup>47,48</sup> and thus can be over-represented in biodiversity datasets. In contrast, in some citizen science projects, abundant species can form a disproportionate number of records (e.g.,<sup>49</sup>) because species' abundance can lead to an increase in the number of records by casual observers<sup>50</sup>. Differences in species detectability<sup>50</sup>, and the ease of making the observations, also lead to taxonomic biases in citizen science datasets. Therefore, species traits (e.g., body size, color, flock size) may have an additive effect, influencing both the detectability of a species<sup>51–53</sup>, and in turn, the likelihood of a species being submitted to an unstructured citizen science database.

Quantifying biases in citizen science datasets can help (1) researchers using these data to better account for biases when drawing ecological conclusions, (2) the design and implementation of future citizen science projects, and (3) understand what species or regions may need data collection from professional scientists by understanding the 'limits' of citizen science projects<sup>19</sup>. Here, we quantify biases in bird observation data from an unstructured, citizen science project—iNaturalist—with that from a semi-structured one—eBird. We restricted our comparison to birds because (1) birds are among the most popular taxa with the non-scientific public, ensuring large sample sizes in both citizen science projects, and (2) data on the species traits that may influence the likelihood of unstructured observations are readily available for birds. We assessed the over-representation or under-representation of bird species' observations in the unstructured citizen science project compared to the semi-structured project (see Fig. 1). We then tested the following predictions: that (1) more colorful species; (2) larger species; (3) species with the 'least concern' IUCN status; and (4) more gregarious species (i.e., with larger flock sizes) are over-represented in the unstructured citizen science dataset (iNaturalist) in contrast to the semi-structured citizen science dataset (eBird). Our analysis highlights the importance of considering species' traits when using citizen science data in ecological research.

## Methods

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

*iNaturalist citizen science data.* iNaturalist is a multi-taxon citizen science project hosted by the California Academy of Sciences. It is an unstructured citizen science project where volunteers contribute opportunistic photos or sound recordings through a smart-phone or web-portal. Photos are then identified to the lowest possible taxonomic resolution using a community identification process, and once two users, or more than two-thirds, confirm the species-level identification of an organism it is considered "research grade". Observations that are



**Figure 1.** A conceptual figure depicting the methods used in our analysis. We used the residual from the relationship between the number of log<sub>10</sub>-transformed eBird observations (i.e., semi-structured citizen science observations) and log<sub>10</sub>-transformed iNaturalist observations (i.e., unstructured citizen science observations) to quantify the over- or under-representation of a species in unstructured citizen science data. We predicted that species which were over-represented in unstructured iNaturalist data would be larger in size, occur more frequently in large flocks, be brighter in color, and be categorized as Least Concern IUCN status (a proxy for commonness).

research grade are then uploaded to the Global Biodiversity Information Facility. We downloaded iNaturalist observations from the Global Biodiversity Information Facility for the contiguous United States<sup>55</sup> for the period from January 2010 to May 2019, on December 3rd, 2019. For more details on the iNaturalist methodology, see here: <https://www.inaturalist.org/pages/getting+started>.

**eBird citizen science data.** eBird is one of the most successful citizen science projects in the world, with > 1 billion bird observations globally. It was launched in 2002 by the Cornell Lab of Ornithology and focuses on collecting reliable data on the distributions and relative abundance of birds throughout the world<sup>54</sup>. It is a semi-structured project where volunteers submit ‘checklists’ of birds seen and/or heard on birding outings, following different protocols (e.g., stationary, incidental, or travelling). An important component of eBird that differentiates it from unstructured data collection is that users are required to indicate whether the checklist is ‘complete’—meaning they included all species they were able to identify during that birding outing. When using complete checklists only in an analysis, a user can infer non-detections in the dataset for any species not recorded. Observers can submit checklists at any time and place of their choosing, and for any duration and distance travelled. All complete checklists additionally include the duration and distance travelled while birding. Filters are set—based on spatiotemporal coordinates—which restrict the species and their associated counts that can be added to the downloadable eBird dataset without approval from a regional expert reviewer<sup>56</sup>. We used the eBird basic dataset (version ebd\_May-2019) and restricted our analysis to data from the contiguous United States for the period from January 2010 to May 2019. We also restricted the data used to those of the best ‘quality’ by excluding incomplete checklists, checklists that were incidental or historical, which travelled > 5 km, lasted < 5 min, and lasted > 240 min, minimizing the leverage of outliers on analyses<sup>57,58</sup>.

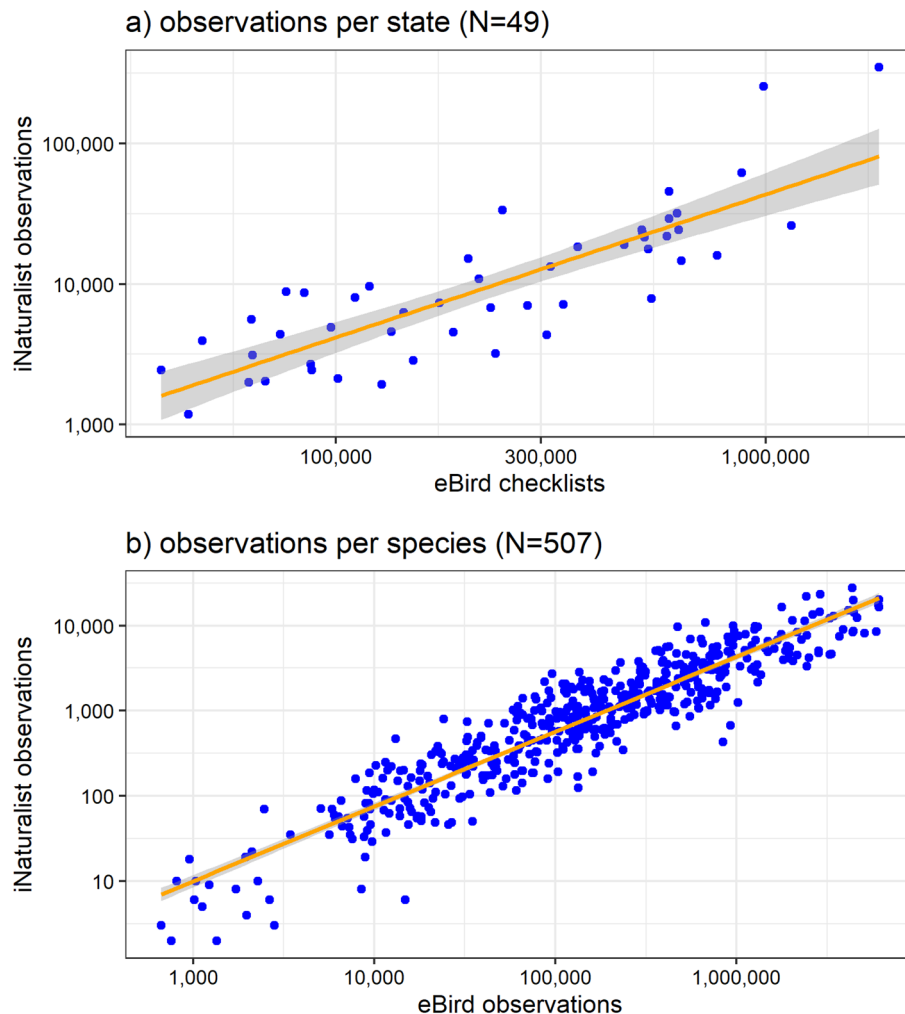
**Filtering and aggregating the citizen science datasets.** Although both datasets are global in scope, we restricted our analysis to the contiguous United States as both of these citizen science projects initiated in the United States, and thus the data are most numerous from there. For comparisons, we aggregated data at the state level. This was done to account for differences that may exist throughout the entirety of the United States including differences in user behavior and the species pools that differ geographically. We used the eBird Clements taxonomy (version 2018) and all species from iNaturalist were matched with this taxonomy. A total of 1030 species was initially collated from the eBird checklists, but many of these only occurred once or a few times—possibly

representing misidentifications that had not yet been fixed by local reviewers or escaped and exotic birds which are incorporated in the eBird dataset but not considered part of the local avifauna or of interest to our analysis here. Although, these could represent scarce and uncommon species in a state as well, albeit these are rarely sampled by iNaturalist. To account for these biases, we removed species that were on < 1% of eBird checklists for a given state; trimming the eBird observations to the ‘core’ suite of species that occur in a state (*sensu*<sup>57</sup>). After trimming the species and harmonizing the taxonomy with iNaturalist, there were 507 species remaining which were considered in our main analyses presented throughout the results. Although our results here are presented using the 1% cutoff level, we tested the sensitivity of this cutoff level and found comparable results across 0, 0.5, 1, and 1.5% cutoffs. For each state, the eBird and iNaturalist data were summarized by calculating the total number of observations in that state for every species where an observation represents a single unique sighting of a species for both iNaturalist and eBird. Using these aggregated data, we conducted preliminary comparisons of the unstructured and semi-structured datasets by quantifying the relationship between the number of eBird checklists and iNaturalist observations at the state level, and the number of unique observations at the species level. We also explored the relationship between the proportion of checklists a species was found on and the proportion of all iNaturalist observations a species represented at the state level.

**Species-specific over- or under-representation in iNaturalist.** Our first analytical step was to model the log–log linear relationship between the total number of observations in iNaturalist and total number of observations in eBird for a species (Fig. 1). This linear model was repeated separately for each state, where the response variable was log-transformed number of iNaturalist observations and the predictor variable was log-transformed number of eBird observations. We repeated the model by state to account for inherent differences among states that were not of interest in our particular analysis, such as (1) the number of observers in a state, (2) the different relative abundance of a species throughout the United States, and (3) any other intrinsic differences that might exist among states that was not of interest in our analysis. A species with a high (i.e., positive) residual would be over-represented in iNaturalist relative to eBird, whereas a species with a low (i.e., negative) residual would be under-represented in iNaturalist relative to eBird (Fig. 1). First, we visualized these state-specific residuals along the trait variables to empirically summarize and visualize the results of this methodological approach. Second, we took the residuals from these models and used these as the response variables in our subsequent analyses of species trait characteristics (see below).

**Species-specific trait data.** We tested whether four predictor variables (see Fig. 1) would explain the over- or under-representation of bird species in the unstructured citizen science data. For each species, we used a proxy for their commonness/abundance, categorized according to IUCN status, taken from HBW BirdLife international checklist version 3 (<http://datazone.birdlife.org/species/taxonomy>). This variable was treated as an ordinal variable in our models (see below) and encompassed Least Concern, Vulnerable, and Near Threatened species. The three species recorded as endangered were removed from this analysis due to a lack of statistical power at this level with so few observations. For each species we used the continuous predictor variables of (1) body size; (2) color; and (3) average flock size. Body sizes (adult body mass in grams) were taken from the amniote life history database compiled by Myhrvold et al.<sup>59</sup> and were log-transformed to meet normality assumptions. Color was taken from Dale et al.<sup>60</sup> and was extracted as RGB values for six patches per species (upper breast, lower breast, crown, forehead, nape, throat). To define a continuum of color where the brightest/most colorful (and likely most detectable species based on color) had the highest value we combined both the ‘distance from brown’ and the ‘brightness’ of a species for the data from Dale et al.<sup>60</sup>. Distance from brown was defined as the maximum Euclidian distance in the cubic RGB color space from brown (R = 102, B = 68, G = 0) from any of the six patches on a species, regardless of sex (i.e., the highest value across both sexes). Brightness was defined as the maximum relative luminance (i.e.,  $0.2126R + 0.7152G + 0.0722B$ ) from any of the six patches on a species, regardless of sex. These two variables were combined and scaled from 0 to 1 for all species in Dale et al.<sup>60</sup> and this value was used as our measure of color. Calculations were done in “Lab” space, an approximately perceptually uniform color space standardized by the Commission Internationale d’Eclairage. Exploratory analyses showed similar results with HSV color space. Flock size—an approximation of the gregariousness of a species—was taken from eBird as the average number of reported individuals among all checklists where a species was reported, across all data. We acknowledge that the number of a species reported on an eBird checklist likely encompasses both the gregariousness of a species as well as the density of a species in an area, as birders can travel through multiple territories. However, qualitative exploration of the flock size variable aligned with a priori expectations of average flock size (Table S1).

**Statistical analysis.** We used mixed effects models to examine the effects of species traits on the relative bias between our unstructured and semi-structured citizen science datasets. The response variable was the residuals from a log–log linear model fit between the eBird observations and the iNaturalist observations for a given species (described above), the predictor variables were the respective traits, and the random effect (i.e., random intercept) was state. By using state as a random effect, we accounted for any replication of some species across states and the varying degree of over, or under-representation of that species across states. First, we ran a global model where all traits were included as predictor variables: log<sub>10</sub>-transformed body size, log<sub>10</sub>-transformed flock size, color, and IUCN status treated as an ordinal variable where Least Concern was coded as the highest variable and Near Threatened as the lowest. Second, to confirm the results of this global model, we ran four separate models—one for each trait as listed above—because there was much missing data for species’ traits. This approach allowed us to test the relationship of a predictor given the other predictor variables (i.e.,



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

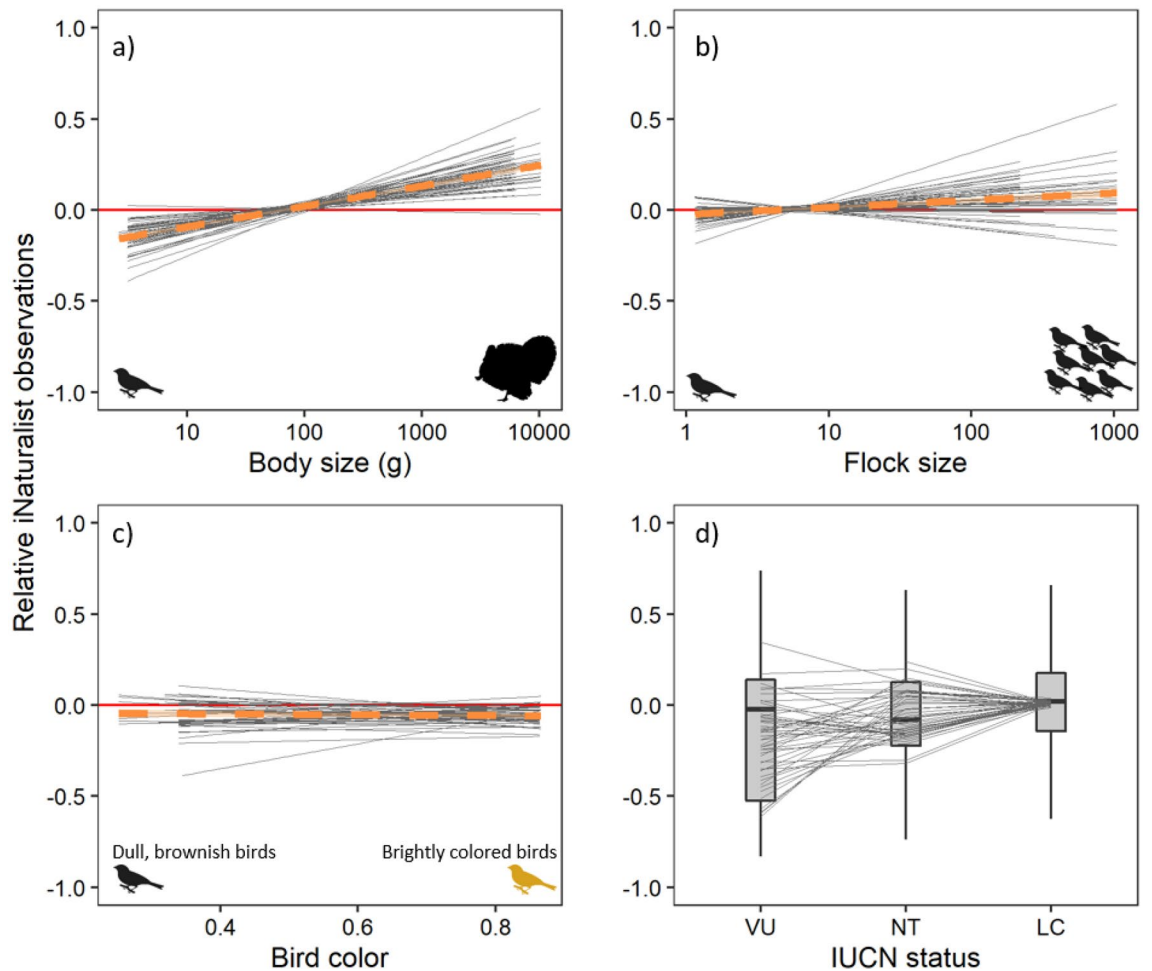
all predictors against the response variable simultaneously) as well as the independent relationships (i.e., each predictor separately against the response variable).

**Data analysis and availability.** All analyses were carried out in R software<sup>61</sup> and relied heavily on the tidyverse workflow<sup>62</sup>. Mixed-effects models were fitted using the lme4 package<sup>63</sup> and  $p$ -values were computed using the lmerTest package<sup>64</sup>. Data and code to reproduce these analyses are available here: <https://doi.org/10.5281/zenodo.5509770>.

## Results

A total of 255,727,592 eBird and 1,107,224 iNaturalist observations were used in our analysis. At the state level, the number of eBird checklists and the number of iNaturalist observations were strongly correlated (Fig. 2a;  $R^2 = 0.58$ ,  $p$ -value  $< 0.001$ ). Similarly, at the species level, the total number of iNaturalist observations and eBird observations for a given species was strongly correlated (Fig. 2b;  $R^2 = 0.9$ ), and for both datasets the number of observations per species was positively-skewed (Figure S1). We also found that the percent of eBird checklists a species was found on and the percent of total iNaturalist observations a species comprised was strongly correlated among states (Figure S2), suggesting that species are sampled to a proportionally similar extent in unstructured and semi-structured citizen science projects.

Across the 507 species included in our analyses (Table S1), we showed that larger species were more likely to be over-represented in the unstructured citizen science dataset, and this was true in most states, as illustrated by the empirical comparison (Fig. 3a). The empirical comparison also showed over-representation of flock size in the unstructured dataset, although some states showed a negative relationship indicating the possibility that



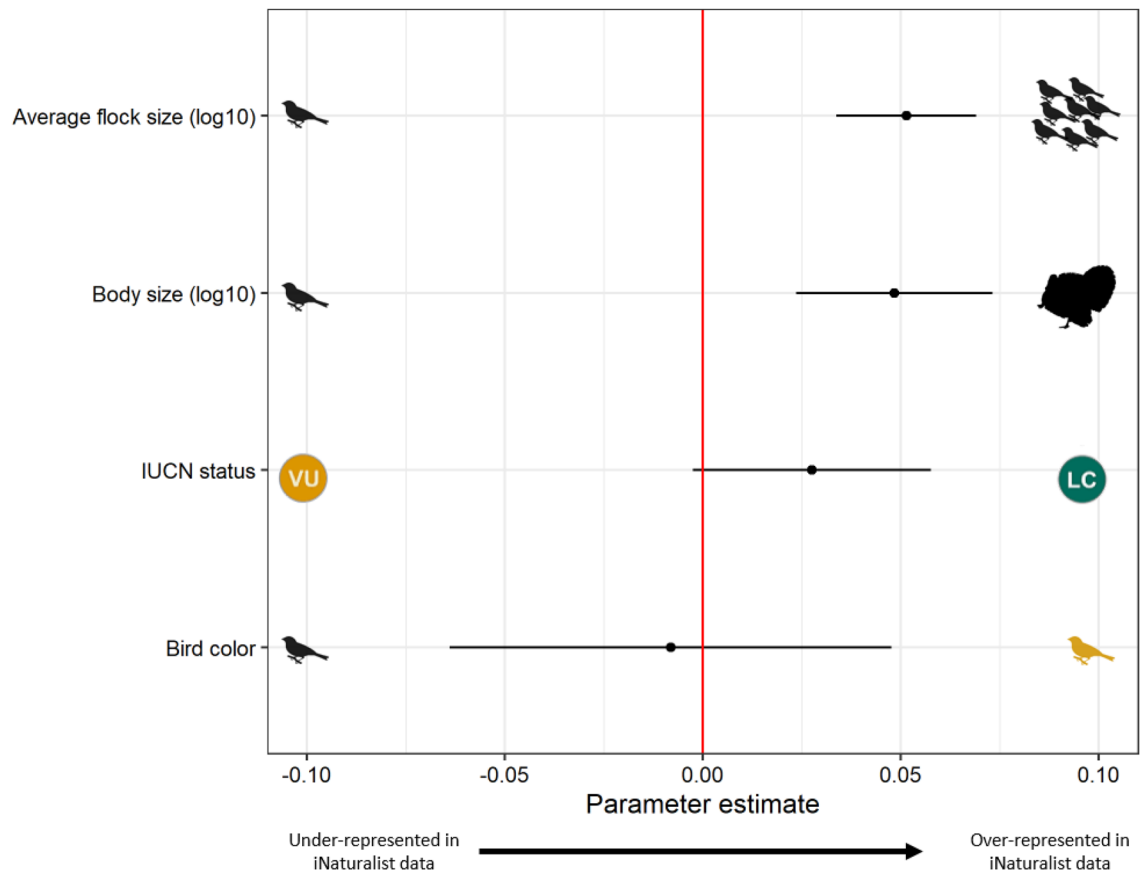
**Figure 3.** The relationship between (a) body size of a species, (b) flock size, (c) color and (d) commonness and the residuals of a linear model fit between iNaturalist and eBird observations (see Fig. 1). These empirical results demonstrate that there is a strong bias of body size in iNaturalist compared with eBird. Positive values on the y-axis mean over-represented in iNaturalist and negative values on the y-axis mean under-represented in iNaturalist. Body size and flock size are represented on a log<sub>10</sub> scale. Each line represents a state (N = 49). For (a)–(c), the overall relationship pooling states is represented by the orange fitted line and 95% confidence interval. The data represented here were then used in our mixed effects model (see Fig. 4 for results).

this trait varies in space (Fig. 3b). There was no discernible pattern in the relationship between color and over- or under-representation in iNaturalist data (Fig. 3c), and there was some evidence that Least Concern species were over-represented in the iNaturalist data (Fig. 3d).

Using the data visualized in our empirical comparisons (Fig. 3), our mixed effects multiple regression linear model (N = 3986 observations and 222 species) with state as a random effect (Fig. 4) found strong evidence that body size (parameter estimate = 0.049; 95% CI = 0.023, 0.073) and flock size (parameter estimate = 0.051; 95% CI = 0.034, 0.069) were over-represented in iNaturalist compared with eBird; moderate evidence that common species were over-represented (parameter estimate = 0.027; 95% CI = -0.003, 0.058); and no evidence that color influenced the over- or under-representation of a species in iNaturalist (parameter estimate = -0.008; 95% CI = -0.064, 0.048). The patterns found in the multiple regression model qualitatively matched that of the individual trait models, where more observations were included in some instances (see Table S2).

## Discussion

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



**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.

over-represented. In contrast to our prediction, however, we found no evidence that brightly-colored species were over-represented in unstructured citizen science data.

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

Related to detectability, an important distinction between iNaturalist and eBird is how identifications are made. For an observer to make a record in iNaturalist, usually a photograph is uploaded (although sound recordings are also accepted). Because a photograph is needed, the submission process is two-fold—first, it needs to be detected, and second, it needs to be photographed, which is likely easier for many large-bodied birds. Longer lenses, often restricted to serious photographers, may be needed to photograph smaller-bodied birds whereas smartphones can usually capture a sufficient image of a larger-bodied bird. In contrast to iNaturalist, in eBird, a lot of identifications are made acoustically, and identification can sometimes also use contextual clues such as behavior or habitat of the bird—often difficult to capture in a photograph. Most traits analyzed here are related to visual encounter/identification, thus likely explaining the differences found between the unstructured iNaturalist and the semi-structured eBird data. To illustrate this difference, in New York state, the most under-represented species in iNaturalist (i.e., with the lowest residuals) are Marsh Wren, American Crow, Warbling Vireo, Least Flycatcher, Willow Flycatcher—all species that are identified largely acoustically. In contrast, the most over-represented species in iNaturalist (i.e., with the highest residuals) are House Sparrow, American Robin, Palm Warbler, Northern Mockingbird—all species that are easy to visually see and thus detect and photograph (Table S1). Therefore, the bias towards large-bodied birds in the unstructured data is probably a result of detectability and the ability to capture a photograph<sup>53</sup>. Photographs can also be uploaded to eBird, and a further test of this hypothesis could interrogate the species in eBird which have photographs uploaded. This process is similar in insects, for example, which are generally small, but larger insects (e.g., butterflies) are both easier to observe, photograph, and identify—making it likely that the biases we found in birds generalize to insects as well. Indeed,

a study of bugs and beetles found that smaller species are typically less represented in citizen science data<sup>68</sup>. Importantly, because this body size bias is systematic, it may be easier to model, as we know that these data are not missing at random (e.g.,<sup>70</sup>) and thus body size can be included in various modelling processes when using unstructured citizen science data (e.g.,<sup>67</sup>).

Similar to body size, our mixed effects model found that birds which occur in larger groups (i.e., flocks) are over-represented in the unstructured dataset. This, again, may be inherently linked to the recording process, rather than a specific bias or preference of the observers themselves. This is because common birds, that occur in large flocks, are more likely to be seen and thus submitted to the unstructured citizen science project<sup>65</sup>. A larger flock will likely also provide more opportunities to capture a photograph than when observing a single individual, as has been illustrated in the detectability of animals from aerial surveys by professionals<sup>71</sup>. However, we note that our measure of flock size (i.e., average checklist count) may not truly represent the gregariousness of a species when observed by an eBirder or iNaturalist contributor. The variability among states in the empirical comparison of flock size (e.g., Fig. 3) illustrates that this overall bias is not uniform across all regions in the United States. Similarly, for body size there was some differentiation in the empirical comparison across states, although to a lesser extent than flock size. Such differences in bias between eBird and iNaturalist can be driven by many external factors such as the avifauna in the region, the number of eBird or iNaturalist users and the characteristics of those users. Future work should look to investigate what factors influence the biases we highlight here in space and/or time.

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

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

While citizen science data are undoubtedly valuable for ecology and conservation<sup>4,76,77</sup>, there remain limits to the use of citizen science datasets<sup>13,78</sup>. The ability to sample remote regions, for example, will likely remain a limitation in citizen science data, and this has been well-recognized<sup>17</sup>. Quantifying the limits of citizen science datasets for use in ecology and conservation remains an important step for the future widespread use of citizen science data in ecology and conservation. Data-integration—where noisy citizen science data are integrated with professionally-curated datasets—will likely be increasingly important in the future use of citizen science data<sup>79,80</sup>. By knowing the biases present in citizen science data, experts can preferentially generate data that maximizes the integration process, for example by collecting data from remote regions. Further, professional scientists could use limited funding to target species that are likely to be under-represented in some citizen science datasets—i.e., rare, small-bodied, species.

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



biases change from unstructured to semi-structured to structured citizen science platforms. Nevertheless, our results demonstrate the importance of using species-specific traits when modelling citizen science data<sup>27,29,52,81–84</sup>.

Received: 5 March 2021; Accepted: 6 September 2021

Published online: 24 September 2021

## References

- Pocock, M. J., Tweddle, J. C., Savage, J., Robinson, L. D. & Roy, H. E. The diversity and evolution of ecological and environmental citizen science. *PLoS ONE* **12**, e0172579 (2017).
- Chandler, M. *et al.* Contribution of citizen science towards international biodiversity monitoring. *Biol. Cons.* **213**, 280–294 (2017).
- Chandler, M. *et al.* Involving citizen scientists in biodiversity observation. In *The GEO Handbook on Biodiversity Observation Networks* 211–237 (Springer, 2017).
- McKinley, D. C. *et al.* Citizen science can improve conservation science, natural resource management, and environmental protection. *Biol. Cons.* **208**, 15–28 (2017).
- Pereira, H. M. *et al.* Monitoring essential biodiversity variables at the species level. In *The GEO Handbook on Biodiversity Observation Networks* 79–105 (Springer, 2017).
- Wiggins, A. & Crowston, K. From conservation to crowdsourcing: A typology of citizen science. in *2011 44th Hawaii International Conference on System Sciences* 1–10 (IEEE, 2011).
- Haklay, M. Citizen science and volunteered geographic information: Overview and typology of participation. In *Crowdsourcing Geographic Knowledge* 105–122 (Springer, 2013).
- Kelling, S. *et al.* Using semistructured surveys to improve citizen science data for monitoring biodiversity. *Bioscience* **69**, 170–179 (2019).
- Welvaert, M. & Caley, P. Citizen surveillance for environmental monitoring: Combining the efforts of citizen science and crowdsourcing in a quantitative data framework. *Springerplus* **5**, 1890 (2016).
- Callaghan, C. T., Rowley, J. J., Cornwell, W. K., Poore, A. G. & Major, R. E. Improving big citizen science data: Moving beyond haphazard sampling. *PLoS Biol.* **17**, e3000357 (2019).
- Bonter, D. N. & Cooper, C. B. Data validation in citizen science: A case study from project FeederWatch. *Front. Ecol. Environ.* **10**, 305–307 (2012).
- Kosmala, M., Wiggins, A., Swanson, A. & Simmons, B. Assessing data quality in citizen science. *Front. Ecol. Environ.* **14**, 551–560 (2016).
- Burgess, H. K. *et al.* The science of citizen science: Exploring barriers to use as a primary research tool. *Biol. Cons.* **208**, 113–120 (2017).
- Courter, J. R., Johnson, R. J., Stuyck, C. M., Lang, B. A. & Kaiser, E. W. Weekend bias in citizen science data reporting: Implications for phenology studies. *Int. J. Biometeorol.* **57**, 715–720 (2013).
- Sullivan, B. L. *et al.* The eBird enterprise: An integrated approach to development and application of citizen science. *Biol. Cons.* **169**, 31–40 (2014).
- Kelling, S. *et al.* Can observation skills of citizen scientists be estimated using species accumulation curves?. *PLoS ONE* **10**, e0139600 (2015).
- Tiago, P., Ceia-Hasse, A., Marques, T. A., Capinha, C. & Pereira, H. M. Spatial distribution of citizen science casuistic observations for different taxonomic groups. *Sci. Rep.* **7**, 1–9 (2017).
- Geldmann, J. *et al.* What determines spatial bias in citizen science? Exploring four recording schemes with different proficiency requirements. *Divers. Distrib.* **22**, 1139–1149 (2016).
- Callaghan, C. T. *et al.* Three frontiers for the future of biodiversity research using citizen science data. *Bioscience* **71**, 55–63 (2021).
- Ward, D. F. Understanding sampling and taxonomic biases recorded by citizen scientists. *J. Insect Conserv.* **18**, 753–756 (2014).
- Troutet, J., Grandcolas, P., Blin, A., Vignes-Lebbe, R. & Legendre, F. Taxonomic bias in biodiversity data and societal preferences. *Sci. Rep.* **7**, 1–14 (2017).
- Martín-López, B., Montes, C., Ramírez, L. & Benayas, J. What drives policy decision-making related to species conservation? *Biol. Conserv.* **142**, 1370–1380 (2009).
- Boakes, E. H. *et al.* Distorted views of biodiversity: Spatial and temporal bias in species occurrence data. *PLoS Biol.* **8**, e1000385 (2010).
- Aceves-Bueno, E. *et al.* The accuracy of citizen science data: A quantitative review. *Bull. Ecol. Soc. Am.* **98**, 278–290 (2017).
- Davies, T. K., Stevens, G., Meekan, M. G., Struve, J. & Rowcliffe, J. M. Can citizen science monitor whale-shark aggregations? Investigating bias in mark-recapture modelling using identification photographs sourced from the public. *Wildl. Res.* **39**, 696–704 (2013).
- Crall, A. W. *et al.* Assessing citizen science data quality: An invasive species case study. *Conserv. Lett.* **4**, 433–442 (2011).
- van Strien, A. J., van Swaay, C. A. & Termaat, T. Opportunistic citizen science data of animal species produce reliable estimates of distribution trends if analysed with occupancy models. *J. Appl. Ecol.* **50**, 1450–1458 (2013).
- Johnston, A., Moran, N., Musgrove, A., Fink, D. & Baillie, S. R. Estimating species distributions from spatially biased citizen science data. *Ecol. Model.* **422**, 108927 (2020).
- Tiago, P., Pereira, H. M. & Capinha, C. Using citizen science data to estimate climatic niches and species distributions. *Basic Appl. Ecol.* **20**, 75–85 (2017).
- Sullivan, B. L. *et al.* Using open access observational data for conservation action: A case study for birds. *Biol. Cons.* **208**, 5–14 (2017).
- Callaghan, C. T. *et al.* Citizen science data accurately predicts expert-derived species richness at a continental scale when sampling thresholds are met. *Biodivers. Conserv.* **29**, 1323–1337 (2020).
- Birkin, L. & Goulson, D. Using citizen science to monitor pollination services. *Ecol. Entomol.* **40**, 3–11 (2015).
- Delaney, D. G., Sperling, C. D., Adams, C. S. & Leung, B. Marine invasive species: Validation of citizen science and implications for national monitoring networks. *Biol. Invasions* **10**, 117–128 (2008).
- Schultz, C. B., Brown, L. M., Pelton, E. & Crone, E. E. Citizen science monitoring demonstrates dramatic declines of monarch butterflies in western north america. *Biol. Cons.* **214**, 343–346 (2017).
- Bird, T. J. *et al.* Statistical solutions for error and bias in global citizen science datasets. *Biol. Cons.* **173**, 144–154 (2014).
- Isaac, N. J., van Strien, A. J., August, T. A., de Zeeuw, M. P. & Roy, D. B. Statistics for citizen science: Extracting signals of change from noisy ecological data. *Methods Ecol. Evol.* **5**, 1052–1060 (2014).
- Dickinson, J. L. *et al.* The current state of citizen science as a tool for ecological research and public engagement. *Front. Ecol. Environ.* **10**, 291–297 (2012).
- Bonney, R. *et al.* Next steps for citizen science. *Science* **343**, 1436–1437 (2014).
- Jordan, R. C., Gray, S. A., Howe, D. V., Brooks, W. R. & Ehrenfeld, J. G. Knowledge gain and behavioral change in citizen-science programs. *Conserv. Biol.* **25**, 1148–1154 (2011).

40. Crall, A. W. *et al.* The impacts of an invasive species citizen science training program on participant attitudes, behavior, and science literacy. *Public Underst. Sci.* **22**, 745–764 (2013).
41. Jordan, R. C., Ballard, H. L. & Phillips, T. B. Key issues and new approaches for evaluating citizen-science learning outcomes. *Front. Ecol. Environ.* **10**, 307–309 (2012).
42. Evans, C. *et al.* The neighborhood nestwatch program: Participant outcomes of a citizen-science ecological research project. *Conserv. Biol.* **19**, 589–594 (2005).
43. Haywood, B. K., Parrish, J. K. & Dolliver, J. Place-based and data-rich citizen science as a precursor for conservation action. *Conserv. Biol.* **30**, 476–486 (2016).
44. Pocock, M. J. *et al.* A vision for global biodiversity monitoring with citizen science. In *Advances in Ecological Research* vol. 59, 169–223 (Elsevier, 2018).
45. Tiago, P., Gouveia, M. J., Capinha, C., Santos-Reis, M. & Pereira, H. M. The influence of motivational factors on the frequency of participation in citizen science activities. *Nat. Conserv.* **18**, 61 (2017).
46. Isaac, N. J. & Pocock, M. J. Bias and information in biological records. *Biol. J. Lin. Soc.* **115**, 522–531 (2015).
47. Angulo, E. & Courchamp, F. Rare species are valued big time. *PLoS ONE* **4**, e5215 (2009).
48. Booth, J. E., Gaston, K. J., Evans, K. L. & Armsworth, P. R. The value of species rarity in biodiversity recreation: A birdwatching example. *Biol. Cons.* **144**, 2728–2732 (2011).
49. Rowley, J. J. *et al.* FrogID: Citizen scientists provide validated biodiversity data on frogs of Australia. *Herpetol. Conserv. Biol.* **14**, 155–170 (2019).
50. Boakes, E. H. *et al.* Patterns of contribution to citizen science biodiversity projects increase understanding of volunteers recording behaviour. *Sci. Rep.* **6**, 33051 (2016).
51. Garrard, G. E., McCarthy, M. A., Williams, N. S., Bekessy, S. A. & Wintle, B. A. A general model of detectability using species traits. *Methods Ecol. Evol.* **4**, 45–52 (2013).
52. Denis, T. *et al.* Biological traits, rather than environment, shape detection curves of large vertebrates in neotropical rainforests. *Ecol. Appl.* **27**, 1564–1577 (2017).
53. Sólymos, P., Matsuoka, S. M., Stralberg, D., Barker, N. K. & Bayne, E. M. Phylogeny and species traits predict bird detectability. *Ecography* **41**, 1595–1603 (2018).
54. Wood, C., Sullivan, B., Iliff, M., Fink, D. & Kelling, S. eBird: Engaging birders in science and conservation. *PLoS Biol* **9**, 1001220 (2011).
55. GBIF.org (3rd December 2019). *GBIF occurrence download*. <https://doi.org/10.15468/dl.lpwczr>
56. Gilfedder, M. *et al.* Brokering trust in citizen science. *Soc. Nat. Resour.* **32**, 292–302 (2019).
57. Callaghan, C., Lyons, M., Martin, J., Major, R. & Kingsford, R. Assessing the reliability of avian biodiversity measures of urban greenspaces using eBird citizen science data. *Avian Conserv. Ecol.* **12**, 66 (2017).
58. Johnston, A. *et al.* Best practices for making reliable inferences from citizen science data: Case study using eBird to estimate species distributions. *BioRxiv* 574392 (2019).
59. Myhrvold, N. P. *et al.* An amniote life-history database to perform comparative analyses with birds, mammals, and reptiles: Ecological archives E096–269. *Ecology* **96**, 3109–3109 (2015).
60. Dale, J., Dey, C. J., Delhey, K., Kempnaers, B. & Valcu, M. The effects of life history and sexual selection on male and female plumage colouration. *Nature* **527**, 367–370 (2015).
61. R Core Team. *R: A Language and Environment for Statistical Computing* (R Foundation for Statistical Computing, 2020).
62. Wickham, H. *et al.* Welcome to the tidyverse. *J. Open Source Softw.* **4**, 1686 (2019).
63. Bates, D., Mächler, M., Bolker, B. & Walker, S. Fitting linear mixed-effects models using lme4. *J. Stat. Softw.* **67**, 1–48 (2015).
64. Kuznetsova, A., Brockhoff, P. B. & Christensen, R. H. B. lmerTest package: Tests in linear mixed effects models. *J. Stat. Softw.* **82**, 1–26 (2017).
65. Johnston, A. *et al.* Species traits explain variation in detectability of UK birds. *Bird Study* **61**, 340–350 (2014).
66. Steen, V. A., Elphick, C. S. & Tingley, M. W. An evaluation of stringent filtering to improve species distribution models from citizen science data. *Divers. Distrib.* **25**, 1857–1869 (2019).
67. Henckel, L., Bradter, U., Jönsson, M., Isaac, N. J. & Snäll, T. Assessing the usefulness of citizen science data for habitat suitability modelling: Opportunistic reporting versus sampling based on a systematic protocol. *Divers. Distrib.* **26**, 1276–1290 (2020).
68. Caley, P., Welvaert, M. & Barry, S. C. Crowd surveillance: Estimating citizen science reporting probabilities for insects of biosecurity concern. *J. Pest. Sci.* **93**, 543–550 (2020).
69. Périquet, S., Roxburgh, L., le Roux, A. & Collinson, W. J. Testing the value of citizen science for roadkill studies: A case study from South Africa. *Front. Ecol. Evol.* **6**, 15 (2018).
70. Nakagawa, S. & Freckleton, R. P. Model averaging, missing data and multiple imputation: A case study for behavioural ecology. *Behav. Ecol. Sociobiol.* **65**, 103–116 (2011).
71. Schlossberg, S., Chase, M. & Griffin, C. Using species traits to predict detectability of animals on aerial surveys. *Ecol. Appl.* **28**, 106–118 (2018).
72. Aristeidou, M., Scanlon, E. & Sharples, M. Profiles of engagement in online communities of citizen science participation. *Comput. Hum. Behav.* **74**, 246–256 (2017).
73. Troscianko, J., Skelhorn, J. & Stevens, M. Quantifying camouflage: How to predict detectability from appearance. *BMC Evol. Biol.* **17**, 1–13 (2017).
74. Schuetz, J. G. & Johnston, A. Characterizing the cultural niches of North American birds. *Proc. Natl. Acad. Sci.* **22**, 10868–10873 (2019).
75. Lišková, S. & Frynta, D. What determines bird beauty in human eyes?. *Anthrozoös* **26**, 27–41 (2013).
76. Tulloch, A. I., Possingham, H. P., Joseph, L. N., Szabo, J. & Martin, T. G. Realising the full potential of citizen science monitoring programs. *Biol. Cons.* **165**, 128–138 (2013).
77. Kabori, H. *et al.* Citizen science: A new approach to advance ecology, education, and conservation. *Ecol. Res.* **31**, 1–19 (2016).
78. Callaghan, C. T., Poore, A. G., Major, R. E., Rowley, J. J. & Cornwell, W. K. Optimizing future biodiversity sampling by citizen scientists. *Proc. R. Soc. B* **286**, 20191487 (2019).
79. Pacifici, K. *et al.* Integrating multiple data sources in species distribution modeling: A framework for data fusion. *Ecology* **98**, 840–850 (2017).
80. Robinson, O. J. *et al.* Integrating citizen science data with expert surveys increases accuracy and spatial extent of species distribution models. *Divers. Distrib.* **26**, 976–986 (2020).
81. van Strien, A. J., Termaat, T., Groenendijk, D., Mensing, V. & Kery, M. Site-occupancy models may offer new opportunities for dragonfly monitoring based on daily species lists. *Basic Appl. Ecol.* **11**, 495–503 (2010).
82. Van der Wal, R. *et al.* Mapping species distributions: A comparison of skilled naturalist and lay citizen science recording. *Ambio* **44**, 584–600 (2015).
83. Dennis, E. B., Morgan, B. J., Brereton, T. M., Roy, D. B. & Fox, R. Using citizen science butterfly counts to predict species population trends. *Conserv. Biol.* **31**, 1350–1361 (2017).
84. Stoudt, S., Goldstein, B. R. & De Valpine, P. Identifying charismatic bird species and traits with community science data. <https://doi.org/10.1101/2021.06.05.446577>

## Acknowledgements

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

## Author contributions

C.T.C. conceived and led the study with input from all authors. C.T.C. performed the analyses with input from all authors. C.T.C. wrote the first version of the manuscript and all authors contributed to editing the manuscript.

## Funding

Open Access funding enabled and organized by Projekt DEAL.

## Competing interests

The authors declare no competing interests.

## Additional information

**Supplementary Information** The online version contains supplementary material available at <https://doi.org/10.1038/s41598-021-98584-7>.

**Correspondence** and requests for materials should be addressed to C.T.C.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).

**Publisher's note** Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2021