

# Recognizability bias in citizen science photographs

Wouter Koch<sup>1,2\*</sup>      Laurens Hogeweg<sup>3,4</sup>  
Erlend B. Nilsen<sup>5</sup>      Robert B. O'Hara<sup>6</sup>  
Anders G. Finstad<sup>1</sup>

<sup>1</sup>Department of Natural History, Norwegian University of Science and Technology, Erling  
Skakkes gate 47b, Trondheim, Norway

<sup>2</sup>Norwegian Biodiversity Information Centre, Havnegata 9, 7010 Trondheim, Norway

<sup>3</sup>Intel Benelux, High Tech Campus 83, 5656 AE Eindhoven, The Netherlands

<sup>4</sup>Naturalis Biodiversity Center, PO Box 9517, 2300 RA, Leiden, The Netherlands

<sup>5</sup>Norwegian Institute for Nature Research, Postboks 5685 Torgarden, 7485 Trondheim,  
Norway

<sup>6</sup> Department of Mathematical Sciences, Norwegian University of Science and Technology,  
Alfred Getz' vei 1, 7034 Trondheim, Norway

\*Corresponding author. E-mail:  
wouter.koch@artsdatabanken.no

Contributing authors:  
laurens.hogeweg@naturalis.nl  
erlend.nilsen@nina.no bob.ohara@ntnu.no  
anders.finstad@ntnu.no

## Abstract

Citizen science initiatives and automated collection methods increasingly depend on image recognition in order to provide the amounts of observational data research and management needs. Training recognition models, meanwhile, also requires large amounts of data from these sources, creating a feedback loop between the methods and the tools. Species that are harder to recognize, both for humans and machine learning algorithms, are likely to be underreported, and thus be less prevalent in the training data. As a result, the feedback loop may hamper training mostly for species that already pose the greatest challenge. In this study, we trained recognition models for various taxa, and found evidence for a “*recognizability bias*”, where species that models struggle with are also generally underreported. This has implications for the kind of performance one can expect from future models that are trained with more data, including such challenging species. We consider identification methods that rely on more than photographs alone to be important in improving future identification tools.

## Introduction

There is an ever growing need for large amounts of biodiversity observation data. With an increasing awareness of the multiple crises biodiversity faces [1–3], substantial amounts of such data are essential if humanity is to monitor trends and address these issues [4–6]. Occurrence data are typically subject to spatial, temporal and taxonomic bias [7, 8], and traditional manual methods of data collection are insufficient to gather the data volume needed, or address these biases. Alternative data collection methods, ranging from citizen science (non-professional volunteers reporting observations [9]) to camera-traps automating insect monitoring [10, 11] are being deployed to gather large amounts of data. With the increased output from such initiatives, manual management and quality control become infeasible. Automated image recognition tools for species identification are increasingly used to facilitate this [12–15]. Training image recognition models, however, also requires large amounts of pictures [16]. This creates a mutual reliance between large scale image data collection and image recognition models [17].

Visual identification of species is a complex task, and taxa vary in their recognizability; while some species are unmistakable, many others are very challenging or even outright impossible to identify, regardless of picture quality [18]. As models are trained using training data reported and identified by humans, species with low recognizability among humans will be underreported and be underrepresented in the training data. This affects recognition models, as these are then being trained with data biased towards higher recognizability, consisting mostly of pictures of species that are easier to recognize. If this is the case,

training models will be hampered not only by the lower recognizability of particularly challenging species, but also by their higher absence from the training data.

To evaluate the existence of this possible bias and its consequences, we evaluated how data availability, picture quality, biological traits and data collection differs across species within 3 orders of birds, and how these differences relate to recognition model performance. All data came from a large Norwegian citizen science project, where recognition tools are not a part of the reporting or validation process. Birds are the most well-represented orders per species, allowing for the most detailed analysis. We also trained models for 9 other orders of plants, animals and fungi, to test for a general correlation between data availability and model performance, and to evaluate what this means for future recognition models.

We find evidence for a “*recognizability bias*”, where species that are more readily identified by humans and recognition models alike are more prevalent in the available image data. This pattern is present across multiple taxa, and does not appear to relate to a difference in picture quality, biological traits, or data collection metrics other than recognizability.

## Methods

We trained image recognition models using convolutional neural networks on pictures retrieved from the Norwegian citizen science platform Species Observation Service [19] for 12 orders: Agaricales, Anseriformes, Asparagales, Asterales, Charadriiformes, Coleoptera, Diptera, Lecanorales, Lepidoptera, Odonata, Passeriformes, and Polyporales [20]. A separate model was trained for each order, using 200 documented observations per species for training and validation, and a minimum of 20 for the test set. See Koch *et al.* [21] for details. From these models and various external datasets, several relevant metrics were collected (table 1).

Metric	Definition
Data availability	The total number of citizen science observations from the Norwegian citizen science platform Species Observation Service [19] for a species, containing one or more pictures. This is a more meaningful measure than simply the total number of pictures, as multiple pictures within an observation are not independent from one another and therefore do not add as much information as unique observations.
F <sub>1</sub> -score	The performance obtained for a species in a recognition model, defined as the harmonic mean of the precision and recall [16]

Species in Norway	The number of species within an order that are present in Norway, according to the Norwegian Species Nomenclature Database [22].
-------------------	--

Table 1: Metrics collected for species within all orders

More detailed analyses were done on the included bird orders; waterfowl (Anseriformes), shorebirds (Charadriiformes), and passerines (Passeriformes), as bird orders have the highest proportion of species in Norway represented in the dataset, and ample standardized available data on a range of biological traits allowing for a deeper analysis. For these analyses, a number of additional metrics were collected for the included bird species (table 2).

Metric	Definition
Picture quality	Using Label Studio v1.4 [23], $\geq 50$ pictures per species were annotated by drawing rectangles approximately equal in surface area to the visible part of each individual bird. From this, we took the percentage of the picture occupied by the largest depiction of an individual of the target species, minus the percentage of the picture occupied by all individuals of other bird species. Per species, the median log value was used as a proxy for picture quality.
Urbanness	The proportion of 100 documented observations from the Species Observation Service with a location within a cell tagged as “urban” in the ESA CCI landcover dataset [24].
Hand-wing index	Wing length minus wing width, a measure positively correlated with flight efficiency and dispersal ability of a species. Retrieved from the Global-HWI dataset [25].
Body mass	The average log-transformed body mass of a species, retrieved from the Global-HWI dataset [25].
Habitat openness	A three-step scale of the openness of the habitat of a species, retrieved from the Global-HWI dataset [25].
Documentation rate	The proportion, per species, of observations in the Species Observation Service that have one or more pictures.
Picture density	The average number of pictures per observation from the Species Observation Service, from those with at least one picture.
Observation rate	The number of observations in the Species Observation Service dataset per observation in the TOV-e bird monitoring scheme [26]

---

Table 2: Metrics collected for species within the bird orders

LASSO multiple regression models were trained using Scikit-learn [27] to evaluate the effect of the biological traits, picture quality measurement, and data collection process from table 2 on the  $F_1$ -scores for birds. All LASSO models have the order as a factor. The full model for biological traits is given by

$$F_1 = \beta_0 + \beta_1 HWI + \beta_2 BM + \beta_3 H + \beta_4 U + \beta_5 DA + \epsilon + (1|Order)$$

where  $HWI$  is the hand-wing index,  $BM$  is the body mass,  $H$  is the habitat openness,  $U$  is the urbanness, and  $DA$  is the log data availability. The full model for picture quality is given by

$$F_1 = \beta_0 + \beta_1 Q + \beta_2 DA + \epsilon + (1|Order)$$

where  $Q$  is the picture quality, and  $DA$  is the log data availability. The full model for data collection parameters is given by

$$F_1 = \beta_0 + \beta_1 OR + \beta_2 DR + \beta_3 PD + \beta_4 DA + \epsilon + (1|Order)$$

where  $OR$  is the observation rate,  $DR$  is the documentation rate,  $PD$  is the picture density, and  $DA$  is the log data availability.

## Results

There is a strong positive linear correlation between log data availability and the  $F_1$ -score for bird species (figure 1). Note that data availability does not affect training, as all models were trained and evaluated using 220 documented observations per species, regardless of the total availability. A positive linear correlation was also evident in 7 of the 9 other orders (figure 2), in particular Asterales and Odonata. The beetles (Coleoptera) and lichens (Lecanorales) exhibited no apparent correlation, with an  $R^2$  of 0.06 and 0.12, and P-values of 0.27 and 0.18, respectively.

In each bird order, there is a linear relationship between species' picture density and documentation rate ( $R^2 \geq 0.52$ ,  $p \leq 1.51 \times 10^{-7}$ , see table S2). We also find a negative linear correlation between picture density and  $F_1$ -scores ( $R^2 \geq 0.23$ ,  $p \leq 2.1 \times 10^{-4}$ , see table S2), and some negative linear correlation between documentation rate and  $F_1$ -scores ( $R^2 \geq 0.11$ ,  $p \leq 4.64 \times 10^{-3}$ , see table S2). For passerines, there is a negative linear relationship between habitat openness and picture quality ( $R^2 = 0.26$ ,  $p = 3.53 \times 10^{-8}$ , see table S2). Waterfowl and shorebirds could not be evaluated as they only occur in open habitats.

LASSO models trained on biological traits, collection process parameters, and picture quality, all having and log data availability as an additional parameter and order as a factor, had  $R^2$  values of 0.60, 0.57 and 0.63, respectively.

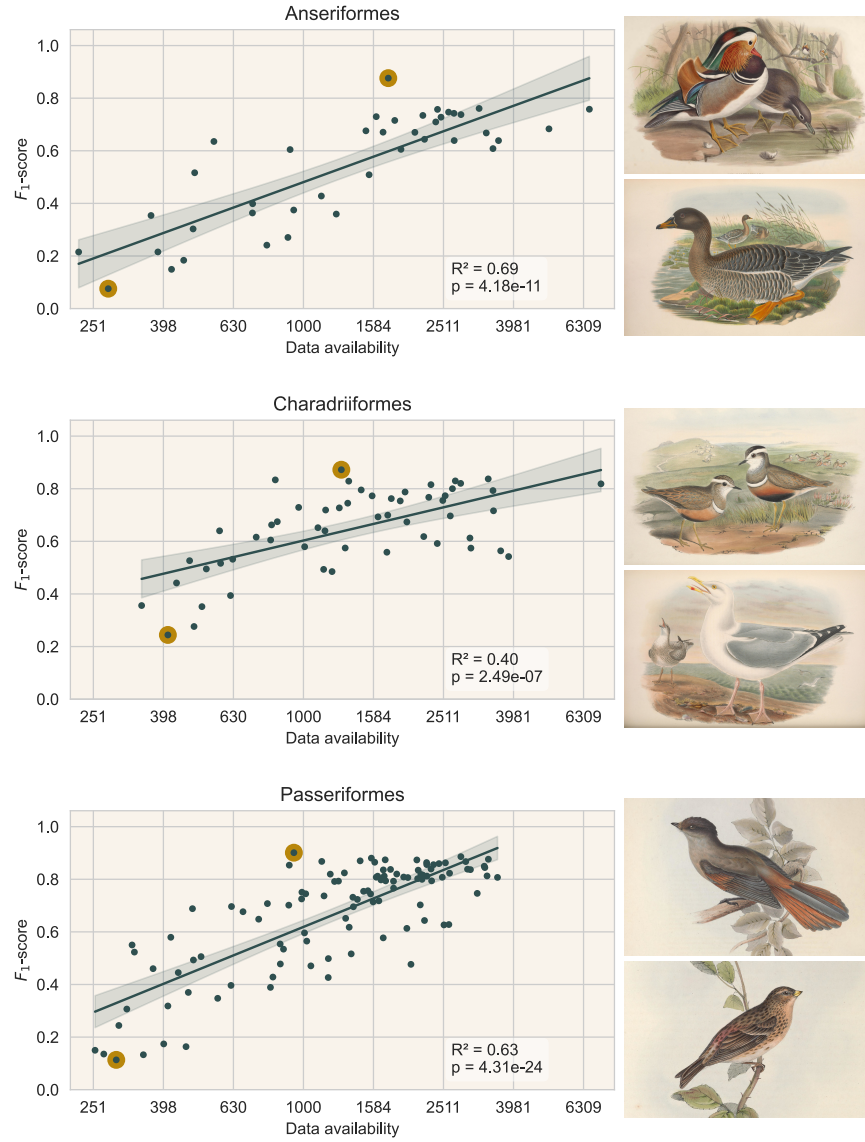


Figure 1: Effect of the total data availability per species on their  $F_1$ -scores, in models trained with 200 documented observations, for three bird orders. The top- and bottom-performing species per order (highlighted dots) are depicted, see table S1. Regressions are Ordinary Least Squares with 95% confidence intervals.

With that, none of the full model performances were substantial improvements from a LASSO model with log data availability as its only parameter ( $R^2 = 0.57$ ).

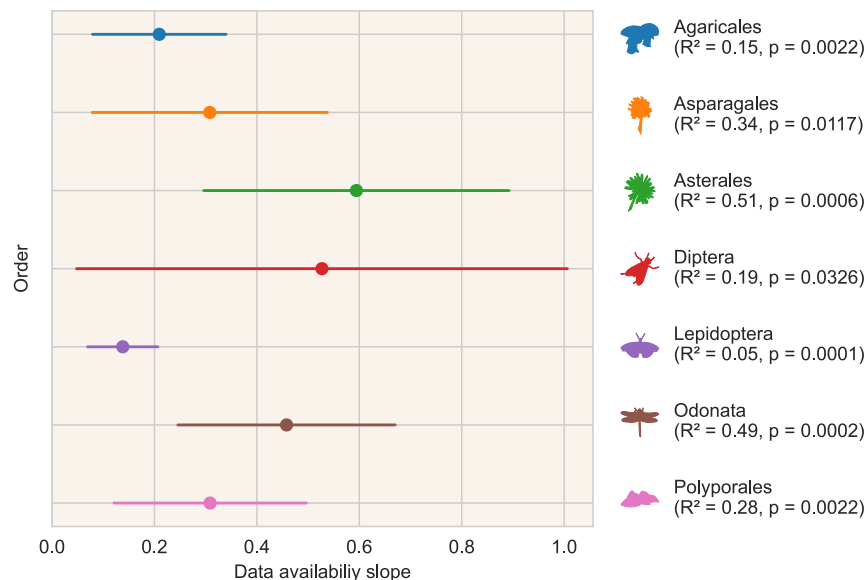


Figure 2: The slopes of the correlations between total data availability per species and their  $F_1$ -scores, in models trained with 200 documented observations, for non-bird orders with a correlation  $p < 0.05$ . Regressions are Ordinary Least Squares, lines indicate the 95% confidence intervals.

## Discussion

We find a conspicuous pattern where recognition models attain higher performances for species that are reported with pictures more frequently. It is probable that the recognizability of the species influences both their likelihood of being reported with pictures, as well as recognition model performances. The citizen science project used as a data source here does not include any recognition tools in its reporting or validation process, allowing a distinction between human and algorithm recognition biases. Unmistakable species can be recognized and reported by more citizen scientists, resulting in greater data availability for such species. A recognition model, dealing with the same information as human observers, is also proportionally more likely to reliably recognize these species.

This is supported by a qualitative comparison between species with the highest and lowest recognition model performances, where easy to recognize, characteristic species are reported more often than hard to recognize species (e.g. nondescript species or species similar to other related species) (see figure 1 and table S1). Further support comes from the fact that most of the correlation is explained by the data availability for a species, rather than the documentation rate or the picture density. Thus, there is more data available mainly when a species is recognized and reported more, rather than it being disproportionately more likely to be reported with pictures, or with many pictures when reported

with pictures.

An alternative explanation to recognizability for increased model performance might be a difference in the kind of pictures, but we find no evidence for this. Species traits, habitat use, and image quality could affect recognition model performance if pictures of more photographed birds are taken more up close, with higher zoom, or were cropped more. We found no evidence, however, for a link between model performance and either picture quality or biological traits in birds. For the passerines, where habitat openness varies among species, we do find that picture quality decreases for species associated with more open habitats. It makes intuitive sense that birds in open habitats are photographed from a greater distance than their forest dwelling counterparts, which will be hidden from view unless in close proximity. While this intercorrelation supports the validity of the picture quality metric, neither habitat nor picture quality affect recognition model performance. We conclude that differences in model performance are caused by the recognizability of the species, rather than by how, or how large species are generally depicted.

Since multiple pictures connected to a single observation are not truly independent, training data are generated based on the number of documented observations, rather than the total number of pictures. One might expect that species with a higher picture density will perform better, as observations with more pictures can provide some additional information in the training process. We find a reverse effect however, where performance for such species is substantially lower. A likely explanation is that species with high picture densities are rarities in Norway (e.g. the top 3 species being Caspian gull, Blyth’s reed warbler, and Pine bunting). Species with the lowest picture density, meanwhile, are typical common, well-known species such as corvids and titmice. Rarities are reported not because they are easy to find or identify by casual observers, but due to their popularity among avid birdwatchers, who are likely to document their observations. A strong correlation between picture density and documentation rate supports this; rarities are more often reported with pictures, and in such cases relatively often with several pictures.

While we investigated the bird orders in detail, the link between data availability and model performance is present in other orders too (figure 2). Some orders are notoriously difficult to identify to species level, e.g. flies (Diptera) and beetles (Coleoptera), but our models for these perform surprisingly well. The list of species with sufficient observations with pictures for inclusion in the experiment reveals that only relatively easy to recognize species, often with distinct colorations (e.g. ladybugs for beetles) are represented in this subset.

More generally, the requirement that species must have at least 220 citizen science observations with pictures generates a non-random subset of species, and it differs greatly per order how selective this criterion is. Bird species are most frequently reported; 48% of the species present in Norway [22] within the bird orders examined here meet the selection criterion. One of the other orders for which the pattern was found, the dragonflies and damselflies (Odonata), have only 52 species in Norway, of which 44% met the criteria for inclusion. This is in stark contrast to the beetles (1% inclusion), and lichens (2% inclusion), where



no clear correlation is found. It is reasonable to assume that for these taxa, the experiment only considers the most recognizable species. If observations were thousandfold, more challenging species could be included, giving a broader range in performances and possibly a similar positive correlation between model performance and data availability.

The consequence of the recognizability bias found here is that as more data is collected, ultimately providing the numbers of pictures needed to train models also on less reported, harder to recognize species, current performance of recognition models cannot be extrapolated to these expanded models. In other words, data that are lacking now are in part lacking because such species are harder to recognize. When such data is added in the future, the performance increase will not be as great as in the past. Besides citizen science, even methods that have no inherent reporting bias, such as automated insect camera traps and trail cameras, can still be subject to recognizability bias. There too, species that are less readily identified will result in more unidentifiable pictures, providing relatively less training data.

Image recognition tools play an important role in maintaining the quality of the large amounts of biodiversity data science and management require. There are limits to what can be identified from a picture however, and identification tools are needed that rely on more than just pixel information. Models that take into account season, location, sound, etc. can be especially beneficial for difficult species. Still, there is no substitute for the taxonomic knowledge of experts. Preserving this knowledge, and making it available in the form of identification keys is vital. These can be powerful tools to more reliably identify challenging species, in tandem with automatic identification.

## Acknowledgements

We are grateful to Rune Sørås, Ingeborg H. Bringslid, and Rienk W. Fokkema for their help in annotating pictures.

## Data accessibility

All code is available through Zenodo at <https://doi.org/10.5281/zenodo.6734696>. Bird illustrations in figure 1 are works in the Public Domain made by John Gould (1804-1881), obtained through the Biodiversity Heritage Library [28–30]

## Authors' contributions

WK: conception, experimental design, code, analysis, writing. LH: code, text revision. EBN: conception, text revision. RBOH: analysis, text revision. AGF: conception, analysis, text revision.

## References

- [1] IPBES. 2022 *Thematic assessment of the sustainable use of wild species of the Intergovernmental Science-Policy Platform on Biodiversity and Ecosystem Services*. Bonn, Germany: IPBES Secretariat. (doi:10.5281/zenodo.6448567).
- [2] IPCC. 2022 *Climate Change 2022: Impacts, Adaptation, and Vulnerability. Contribution of Working Group II to the Sixth Assessment Report of the Intergovernmental Panel on Climate Change*. Cambridge University Press. In Press.
- [3] Secretariat of the Convention on Biological Diversity. 2020. Global biodiversity outlook 5.
- [4] Xu H, Cao Y, Yu D, Cao M, He Y, Gill M, Pereira HM. 2021 Ensuring effective implementation of the post-2020 global biodiversity targets. *Nature Ecology & Evolution* **5**, 4, 411–418. (doi:10.1038/s41559-020-01375-y).
- [5] Wetzel FT, Bingham HC, Groom Q, Haase P, Köljal U, Kuhlmann M, Martin CS, Penev L, Robertson T, Saarenmaa H, *et al.* 2018 Unlocking biodiversity data: Prioritization and filling the gaps in biodiversity observation data in europe. *Biological Conservation* **221**, 78–85. (doi:10.1016/j.biocon.2017.12.024).
- [6] Scholes RJ, Mace GM, Turner W, Geller GN, Jürgens N, Larigauderie A, Muchoney D, Walther BA, Mooney HA. 2008 Toward a global biodiversity observing system. *Science* **321**, 5892, 1044–1045. (doi:10.1126/science.1162055).
- [7] Boakes EH, McGowan PJK, Fuller RA, Chang-qing D, Clark NE, O'Connor K, Mace GM. 2010 Distorted views of biodiversity: Spatial and temporal bias in species occurrence data. *PLoS Biology* **8**, 6, e1000385. (doi:10.1371/journal.pbio.1000385).
- [8] Troudet J, Grandcolas P, Blin A, Vignes-Lebbe R, Legendre F. 2017 Taxonomic bias in biodiversity data and societal preferences. *Scientific Reports* **7**, 1. (doi:10.1038/s41598-017-09084-6).
- [9] Silvertown J. 2009 A new dawn for citizen science. *Trends in Ecology & Evolution* **24**, 9, 467–471. (doi:10.1016/j.tree.2009.03.017).
- [10] Hansen OLP, Svenning JC, Olsen K, Dupont S, Garner BH, Iosifidis A, Price BW, Høye TT. 2019 Species-level image classification with convolutional neural network enables insect identification from habitus images. *Ecology and Evolution* **10**, 2, 737–747. (doi:10.1002/ece3.5921).
- [11] Kirkeby C, Rydhmer K, Cook SM, Strand A, Torrance MT, Swain JL, Prangmsma J, Johnen A, Jensen M, Brydegaard M, *et al.* 2021 Advances in

- automatic identification of flying insects using optical sensors and machine learning. *Scientific Reports* **11**, 1. (doi:10.1038/s41598-021-81005-0).
- [12] Christin S, Hervet É, Lecomte N. 2019 Applications for deep learning in ecology. *Methods in Ecology and Evolution* **10**, 10, 1632–1644. (doi:10.1111/2041-210x.13256).
  - [13] Weinstein BG. 2017 A computer vision for animal ecology. *Journal of Animal Ecology* **87**, 3, 533–545. (doi:10.1111/1365-2656.12780).
  - [14] Wäldchen J, Rzanny M, Seeland M, Mäder P. 2018 Automated plant species identification—trends and future directions. *PLOS Computational Biology* **14**, 4, e1005993. (doi:10.1371/journal.pcbi.1005993).
  - [15] Ceccaroni L, Bibby J, Roger E, Flemons P, Michael K, Fagan L, Oliver JL. 2019 Opportunities and risks for citizen science in the age of artificial intelligence. *Citizen Science: Theory and Practice* **4**, 1. (doi:10.5334/cstp.241).
  - [16] Goodfellow I, Bengio Y, Courville A. 2016 *Deep Learning*. MIT Press.
  - [17] Lotfian M, Ingensand J, Brovelli MA. 2021 The partnership of citizen science and machine learning: Benefits, risks, and future challenges for engagement, data collection, and data quality. *Sustainability* **13**, 14, 8087. (doi:10.3390/su13148087).
  - [18] Lukhtanov VA. 2019 Species delimitation and analysis of cryptic species diversity in the XXI century. *Entomological Review* **99**, 4, 463–472. (doi:10.1134/s0013873819040055).
  - [19] The Norwegian Biodiversity Information Centre. 2022. Norwegian species observation service. (doi:10.15468/zjbbzel).
  - [20] GBIForg. 2021. Gbif occurrence download. (doi:10.15468/DL.TC4W55).
  - [21] Koch W, Hogeweg L, Nilsen EB, Finstad AG. 2022 Maximizing citizen scientists’ contribution to automated species recognition. *Scientific Reports* **12**, 1. (doi:10.1038/s41598-022-11257-x).
  - [22] Norwegian Biodiversity Information Centre. 2021. Species nomenclature database.
  - [23] Tkachenko M, Malyuk M, Holmanyuk A, Liubimov N. 2020-2022. Label Studio: Data labeling software. Open source software available from <https://github.com/heartexlabs/label-studio>.
  - [24] ESA. 2017 Land cover cci product user guide version 2. *Tech. Rep.* .
  - [25] Sheard C, Neate-Clegg MHC, Alioravainen N, Jones SEI, Vincent C, MacGregor HEA, Bregman TP, Claramunt S, Tobias JA. 2020 Ecological drivers of global gradients in avian dispersal inferred from wing morphology. *Nature Communications* **11**, 1. (doi:10.1038/s41467-020-16313-6).

- [26] Kålås JA, Øien IJ, Stokke B, Vang R. 2022. Tov-e bird monitoring sampling data. version 1.6. (doi:10.15468/6jmw2e).
- [27] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, Blondel M, Prettenhofer P, Weiss R, Dubourg V, *et al.* 2011 Scikit-learn: Machine learning in python. *J. Mach. Learn. Res.* **12**, 2825–2830.
- [28] Gould E, Gould J, Lear E. 1837 *The Birds of Europe*. London, Printed by R. and J.E. Taylor, published by the author. (doi:10.5962/bhl.title.65989).
- [29] Gould J, Sharpe RB, Richter HC, Hart WM, Wolf J, Francis T, Walton H. 1850 *The Birds of Asia*. London, Printed by Taylor and Francis, published by the author. (doi:10.5962/bhl.title.54727).
- [30] Gould J, Wolf J, Richter HC, Hart WM, Walter of England, Archbishop of Palermo. 1873 *The birds of Great Britain*. London, Printed by Taylor and Francis, published by the author. (doi:10.5962/bhl.title.127814).

## Supplementary information

Order	Species	F <sub>1</sub> -score
Passeriformes	<i>Perisoreus infaustus</i>	0.901
Passeriformes	<i>Cinclus cinclus</i>	0.886
Passeriformes	<i>Periparus ater</i>	0.88
Passeriformes	<i>Bombycilla garrulus</i>	0.876
Anseriformes	<i>Aix galericulata</i>	0.876
Passeriformes	<i>Certhia familiaris</i>	0.874
Passeriformes	<i>Aegithalos caudatus</i>	0.873
Charadriiformes	<i>Charadrius morinellus</i>	0.872
Passeriformes	<i>Regulus regulus</i>	0.87
Passeriformes	<i>Lophophanes cristatus</i>	0.868
Passeriformes	<i>Emberiza citrinella</i>	0.867
Passeriformes	<i>Garrulus glandarius</i>	0.865
Passeriformes	<i>Pyrrhula pyrrhula</i>	0.863
Passeriformes	<i>Pinicola enucleator</i>	0.863
Passeriformes	<i>Cyanistes caeruleus</i>	0.86
Passeriformes	<i>Sitta europaea</i>	0.856
Passeriformes	<i>Turdus merula</i>	0.855
Passeriformes	<i>Phylloscopus sibilatrix</i>	0.854
Passeriformes	<i>Coccothraustes coccothraustes</i>	0.85
Passeriformes	<i>Carduelis carduelis</i>	0.845
Passeriformes	<i>Motacilla cinerea</i>	0.839
Passeriformes	<i>Erithacus rubecula</i>	0.838
Passeriformes	<i>Parus major</i>	0.837
Charadriiformes	<i>Haematopus ostralegus</i>	0.837
Passeriformes	<i>Motacilla alba</i>	0.837
Passeriformes	<i>Prunella modularis</i>	0.836
Passeriformes	<i>Lanius collurio</i>	0.835
Charadriiformes	<i>Phalaropus lobatus</i>	0.834
Charadriiformes	<i>Calidris maritima</i>	0.83
Charadriiformes	<i>Arenaria interpres</i>	0.829
Passeriformes	<i>Phylloscopus inornatus</i>	0.824
Passeriformes	<i>Saxicola rubicola</i>	0.823
Charadriiformes	<i>Charadrius hiaticula</i>	0.82
Passeriformes	<i>Sylvia atricapilla</i>	0.82
Passeriformes	<i>Turdus philomelos</i>	0.819
Charadriiformes	<i>Gallinago gallinago</i>	0.819
Passeriformes	<i>Luscinia svecica</i>	0.818
Passeriformes	<i>Plectrophenax nivalis</i>	0.816
Charadriiformes	<i>Tringa totanus</i>	0.815
Passeriformes	<i>Lanius excubitor</i>	0.813

Passeriformes	<i>Turdus viscivorus</i>	0.812
Passeriformes	<i>Fringilla coelebs</i>	0.812
Passeriformes	<i>Nucifraga caryocatactes</i>	0.812
Passeriformes	<i>Passer montanus</i>	0.808
Passeriformes	<i>Turdus iliacus</i>	0.808
Passeriformes	<i>Emberiza schoeniclus</i>	0.808
Passeriformes	<i>Oenanthe oenanthe</i>	0.807
Passeriformes	<i>Saxicola rubetra</i>	0.807
Passeriformes	<i>Chloris chloris</i>	0.802
Charadriiformes	<i>Calidris alpina</i>	0.8
Passeriformes	<i>Anthus petrosus</i>	0.8
Passeriformes	<i>Motacilla flava</i>	0.798
Charadriiformes	<i>Cephus grylle</i>	0.795
Passeriformes	<i>Fringilla montifringilla</i>	0.794
Passeriformes	<i>Troglodytes troglodytes</i>	0.794
Charadriiformes	<i>Vanellus vanellus</i>	0.793
Passeriformes	<i>Carpodacus erythrinus</i>	0.793
Passeriformes	<i>Turdus torquatus</i>	0.793
Passeriformes	<i>Eremophila alpestris</i>	0.792
Charadriiformes	<i>Actitis hypoleucos</i>	0.788
Charadriiformes	<i>Pluvialis apricaria</i>	0.773
Charadriiformes	<i>Tringa glareola</i>	0.773
Charadriiformes	<i>Limosa lapponica</i>	0.767
Passeriformes	<i>Ficedula hypoleuca</i>	0.766
Charadriiformes	<i>Charadrius dubius</i>	0.762
Anseriformes	<i>Cygnus olor</i>	0.761
Anseriformes	<i>Mergellus albellus</i>	0.758
Anseriformes	<i>Clangula hyemalis</i>	0.757
Passeriformes	<i>Muscicapa striata</i>	0.756
Charadriiformes	<i>Calidris pugnax</i>	0.755
Passeriformes	<i>Phoenicurus ochruros</i>	0.754
Charadriiformes	<i>Tringa nebularia</i>	0.754
Passeriformes	<i>Acrocephalus schoenobaenus</i>	0.751
Anseriformes	<i>Branta leucopsis</i>	0.747
Passeriformes	<i>Sturnus vulgaris</i>	0.746
Charadriiformes	<i>Limosa limosa</i>	0.745
Passeriformes	<i>Calcarius lapponicus</i>	0.745
Passeriformes	<i>Phoenicurus phoenicurus</i>	0.744
Anseriformes	<i>Mergus merganser</i>	0.742
Anseriformes	<i>Bucephala clangula</i>	0.738
Passeriformes	<i>Curruca communis</i>	0.737
Anseriformes	<i>Tadorna tadorna</i>	0.734
Passeriformes	<i>Poecile montanus</i>	0.732

Anseriformes	<i>Melanitta fusca</i>	0.73
Charadriiformes	<i>Tringa erythropus</i>	0.729
Anseriformes	<i>Anas acuta</i>	0.728
Charadriiformes	<i>Calidris minuta</i>	0.727
Passeriformes	<i>Poecile palustris</i>	0.725
Passeriformes	<i>Passer domesticus</i>	0.723
Charadriiformes	<i>Calidris temminckii</i>	0.719
Passeriformes	<i>Hirundo rustica</i>	0.718
Charadriiformes	<i>Numenius arquata</i>	0.716
Anseriformes	<i>Mareca penelope</i>	0.715
Passeriformes	<i>Corvus frugilegus</i>	0.714
Anseriformes	<i>Mergus serrator</i>	0.71
Passeriformes	<i>Sylvia curruca</i>	0.708
Passeriformes	<i>Turdus pilaris</i>	0.703
Passeriformes	<i>Pica pica</i>	0.702
Charadriiformes	<i>Uria aalge</i>	0.699
Charadriiformes	<i>Chroicocephalus ridibundus</i>	0.697
Passeriformes	<i>Hippolais icterina</i>	0.696
Passeriformes	<i>Loxia leucoptera</i>	0.696
Charadriiformes	<i>Calidris canutus</i>	0.693
Passeriformes	<i>Emberiza pusilla</i>	0.688
Anseriformes	<i>Cygnus cygnus</i>	0.683
Passeriformes	<i>Panurus biarmicus</i>	0.677
Anseriformes	<i>Somateria spectabilis</i>	0.676
Charadriiformes	<i>Alca torda</i>	0.675
Charadriiformes	<i>Rissa tridactyla</i>	0.674
Anseriformes	<i>Branta canadensis</i>	0.671
Anseriformes	<i>Aythya fuligula</i>	0.67
Anseriformes	<i>Anas platyrhynchos</i>	0.668
Charadriiformes	<i>Alle alle</i>	0.663
Charadriiformes	<i>Calidris alba</i>	0.652
Passeriformes	<i>Alauda arvensis</i>	0.651
Passeriformes	<i>Corvus monedula</i>	0.648
Anseriformes	<i>Anas crecca</i>	0.644
Passeriformes	<i>Phylloscopus collybita</i>	0.643
Charadriiformes	<i>Pluvialis squatarola</i>	0.64
Charadriiformes	<i>Fratercula arctica</i>	0.64
Anseriformes	<i>Somateria mollissima</i>	0.639
Anseriformes	<i>Anser anser</i>	0.639
Anseriformes	<i>Anser indicus</i>	0.635
Passeriformes	<i>Acanthis flammea</i>	0.628
Passeriformes	<i>Anthus pratensis</i>	0.627
Charadriiformes	<i>Sterna hirundo</i>	0.618

Passeriformes	<i>Corvus cornix</i>	0.618
Charadriiformes	<i>Stercorarius parasiticus</i>	0.616
Passeriformes	<i>Phylloscopus trochilus</i>	0.613
Charadriiformes	<i>Larus hyperboreus</i>	0.613
Anseriformes	<i>Anser brachyrhynchus</i>	0.608
Anseriformes	<i>Aythya marila</i>	0.605
Charadriiformes	<i>Calidris ferruginea</i>	0.605
Anseriformes	<i>Aythya ferina</i>	0.605
Passeriformes	<i>Corvus corax</i>	0.596
Charadriiformes	<i>Larus fuscus</i>	0.592
Charadriiformes	<i>Tringa ochropus</i>	0.58
Passeriformes	<i>Locustella naevia</i>	0.579
Passeriformes	<i>Carduelis spinus</i>	0.577
Charadriiformes	<i>Numenius phaeopus</i>	0.575
Charadriiformes	<i>Larus canus</i>	0.574
Passeriformes	<i>Carduelis flavirostris</i>	0.565
Charadriiformes	<i>Larus glaucoides</i>	0.564
Charadriiformes	<i>Larus marinus</i>	0.559
Passeriformes	<i>Anthus trivialis</i>	0.554
Passeriformes	<i>Regulus ignicapilla</i>	0.55
Charadriiformes	<i>Larus argentatus</i>	0.542
Passeriformes	<i>Riparia riparia</i>	0.534
Charadriiformes	<i>Scolopax rusticola</i>	0.532
Charadriiformes	<i>Phalaropus fulicarius</i>	0.526
Passeriformes	<i>Sylvia nisoria</i>	0.523
Anseriformes	<i>Polysticta stelleri</i>	0.517
Passeriformes	<i>Acanthis hornemanni</i>	0.516
Charadriiformes	<i>Stercorarius skua</i>	0.516
Anseriformes	<i>Melanitta nigra</i>	0.509
Passeriformes	<i>Sylvia borin</i>	0.506
Passeriformes	<i>Loxia pytyopsittacus</i>	0.498
Charadriiformes	<i>Calidris falcinellus</i>	0.495
Charadriiformes	<i>Sterna paradisaea</i>	0.493
Passeriformes	<i>Lullula arborea</i>	0.493
Charadriiformes	<i>Hydrocoloeus minutus</i>	0.485
Passeriformes	<i>Pastor roseus</i>	0.478
Passeriformes	<i>Loxia curvirostra</i>	0.477
Passeriformes	<i>Acanthis cabaret</i>	0.471
Passeriformes	<i>Turdus atrogularis</i>	0.46
Passeriformes	<i>Ficedula parva</i>	0.445
Charadriiformes	<i>Stercorarius longicaudus</i>	0.442
Passeriformes	<i>Corvus corone</i>	0.428
Anseriformes	<i>Anas clypeata</i>	0.428



Passeriformes	<i>Carduelis cannabina</i>	0.426
Anseriformes	<i>Anser albifrons</i>	0.399
Passeriformes	<i>Acrocephalus dumetorum</i>	0.397
Charadriiformes	<i>Calidris melanotos</i>	0.394
Passeriformes	<i>Acrocephalus palustris</i>	0.389
Anseriformes	<i>Anas strepera</i>	0.375
Passeriformes	<i>Delichon urbicum</i>	0.37
Anseriformes	<i>Mareca strepera</i>	0.364
Anseriformes	<i>Anser fabalis</i>	0.359
Charadriiformes	<i>Thalasseus sandvicensis</i>	0.356
Anseriformes	<i>Branta bernicla</i>	0.354
Charadriiformes	<i>Larus melanocephalus</i>	0.352
Passeriformes	<i>Acrocephalus scirpaceus</i>	0.347
Passeriformes	<i>Motacilla citreola</i>	0.318
Passeriformes	<i>Luscinia luscinia</i>	0.307
Anseriformes	<i>Aythya collaris</i>	0.303
Charadriiformes	<i>Lymnocyptes minimus</i>	0.276
Anseriformes	<i>Spatula clypeata</i>	0.271
Passeriformes	<i>Emberiza leucocephalos</i>	0.244
Charadriiformes	<i>Larus cachinnans</i>	0.244
Anseriformes	<i>Anas querquedula</i>	0.241
Anseriformes	<i>Anas carolinensis</i>	0.215
Anseriformes	<i>Tadorna ferruginea</i>	0.215
Anseriformes	<i>Cygnus columbianus</i>	0.184
Passeriformes	<i>Anthus richardi</i>	0.174
Passeriformes	<i>Spinus spinus</i>	0.164
Passeriformes	<i>Anthus hodgsoni</i>	0.15
Anseriformes	<i>Spatula querquedula</i>	0.149
Passeriformes	<i>Anthus cervinus</i>	0.136
Passeriformes	<i>Linaria cannabina</i>	0.133
Passeriformes	<i>Linaria flavirostris</i>	0.113
Anseriformes	<i>Anser serrirostris</i>	0.076

Table S1: Metrics collected for species within the bird orders

Dependent variable	Parameters	Slope	Intercept	R <sup>2</sup>	P-value
Agaricales F <sub>1</sub> -score	Data availability (log)	0.21	0.27	0.15	$2.15 \times 10^{-3}$
Anseriformes documentation rate	Picture density	0.38	-0.47	0.52	$1.51 \times 10^{-7}$

Anseriformes F <sub>1</sub> -score	Data availability (log)	0.48	-0.97	0.69	$4.18 \times 10^{-11}$
Anseriformes F <sub>1</sub> -score	Documentation rate	-1.52	0.63	0.19	$4.64 \times 10^{-3}$
Anseriformes F <sub>1</sub> -score	Picture density	-1.02	1.95	0.31	$2.10 \times 10^{-4}$
Asparagales F <sub>1</sub> -score	Data availability (log)	0.31	0.01	0.34	0.0117
Asterales F <sub>1</sub> -score	Data availability (log)	0.59	-0.71	0.51	$5.92 \times 10^{-4}$
Charadriiformes documentation rate	Picture density	0.34	-0.43	0.76	$5.51 \times 10^{-18}$
Charadriiformes F <sub>1</sub> -score	Data availability (log)	0.32	-0.34	0.4	$2.49 \times 10^{-7}$
Charadriiformes F <sub>1</sub> -score	Documentation rate	-0.9	0.7	0.28	$3.45 \times 10^{-5}$
Charadriiformes F <sub>1</sub> -score	Picture density	-0.42	1.25	0.39	$3.52 \times 10^{-7}$
Coleoptera F <sub>1</sub> -score	Data availability (log)	0.13	0.57	0.06	0.273
Diptera F <sub>1</sub> -score	Data availability (log)	0.53	-0.63	0.19	0.0326
Lecanorales F <sub>1</sub> -score	Data availability (log)	0.2	0.31	0.12	0.118
Lepidoptera F <sub>1</sub> -score	Data availability (log)	0.14	0.49	0.05	$9.50 \times 10^{-5}$
Odonata F <sub>1</sub> -score	Data availability (log)	0.46	-0.53	0.49	$2.02 \times 10^{-4}$
Passeriformes documentation rate	Picture density	0.3	-0.36	0.55	$1.59 \times 10^{-19}$
Passeriformes F <sub>1</sub> -score	Data availability (log)	0.54	-1	0.63	$4.31 \times 10^{-24}$
Passeriformes F <sub>1</sub> -score	Documentation rate	-0.85	0.71	0.11	$5.19 \times 10^{-4}$
Passeriformes F <sub>1</sub> -score	Picture density	-0.5	1.37	0.23	$1.68 \times 10^{-7}$
Passeriformes picture quality	Habitat openness	-0.12	5.93	0.26	$5.53 \times 10^{-8}$
Polyporales F <sub>1</sub> -score	Data availability (log)	0.31	-0.11	0.28	$2.18 \times 10^{-3}$

Table S2: Metrics collected for species within the bird orders