1 **Non-detection during excursions by citizen scientists**
2 **modeled as a function of weather, season, list length, and individual preferences**

3 Gert W. Jacobusse[1] & Eelke Jongejans[2,3]

4 1 HZ University of Applied Sciences, Middelburg, The Netherlands

5 2 Radboud University, RIBES, Ecology, Nijmegen, The Netherlands

6 3 Netherlands Institute of Ecology, Animal Ecology, Wageningen, The Netherlands

7

8 SUMMARY

9 INTRODUCTION: Citizen science is an increasingly valuable source of information about biodiversity.
10 It is challenging to use this information for analysis of distribution and trends. The lack of a protocol
11 leads to bias in observations and therefore data are not representative. The bias is a consequence of
12 unequal detection probabilities, caused by different preferences and habits of citizen scientists.
13 METHODS: We propose to incorporate characteristics of these excursions in analyses of data collected
14 by citizen scientists to improve estimates of the probability that a species is not detected and reported,
15 even though it does occur. By limiting these models to areas that are known to be occupied, detection
16 can be modeled separately without considering variation in occupancy. We apply this idea to 150
17 common species in the Southwest Delta of The Netherlands, and illustrate the data selection, the
18 modeling process and the results using four species. RESULTS: The strongest features to predict
19 detection are the number of species during a visit (list length), earlier observations of the target
20 species by the same observer, and the day of year. We compare three approaches to predict the total
21 non-detection probability that takes all visits to an area into account. Predictions based on only the
22 number of visits were outperformed by predictions that also take the list length into account. Our
23 predictions based on all features combined consistently beat both other approaches, across all 10
24 species groups that were compared. DISCUSSION: We thus show that explicitly modelling the
25 characteristics of all visits to an occupied area results in estimation of non-detection probabilities,
26 while providing insight into the causes of detection and reporting bias. Furthermore, predictions of
27 our model provide a basis for quantifying the sampling effort in each area, which is a promising first
28 step to correct bias in citizen science data when aiming to map a species' distribution.

29

30 INTRODUCTION

31 Biodiversity plays a fundamental role in ecosystem functioning and resilience (Isbell et al.; 2015, Oliver
32 et al., 2015). Numerous threats such as habitat loss, disturbance, pollution and climate change
33 contribute to the decline of species. To address these threats effectively, it is essential to establish a
34 quantitative understanding of biodiversity across regions (Hochkirch et al., 2021). Data gathered
35 through citizen science are recognized as a valuable source of information about biodiversity (van
36 Strien et al., 2013). Online platforms and the availability of technical aids like automatic species
37 recognition from photos have made it easier to report observations during excursions. These
38 developments increased the amount and quality of available data (Luna et al., 2018). However, the
39 non-systematic nature of observations reported by citizen scientists causes the bias in the resulting
40 datasets to be extensive. Most challenging is the lack of explicit information about the absence of
41 species, as typically only species presences are reported. Biased detection and reporting need to be
42 addressed first, to unlock the full potential of citizen science as a means for quantifying biodiversity.

43  Researchers have recognized and investigated the consequences of bias in previous studies (Isaac et
44  al., 2014; Ranc et al., 2017; Jha et al., 2022), but adequate tools to account for detection and reporting
45  biases are dearly missing.

46  <u>Species distribution models</u>

47  For the quantification of biodiversity across regions, a major consideration is the occupancy
48  of areas by species. Species Distribution Models (SDM) predict where species live (occupied areas),
49  based on climatic and environmental data (Elith & Leathwick, 2009; Melo-Merino et al., 2020). Ideally,
50  SDMs are fitted to observations of both presences and absences of a species, where documented
51  absences can come from both unoccupied and occupied areas (left part of Figure 1). In practice,
52  however, many data sources, especially those gathered through citizen science, only contain presence
53  data (Johnston et al., 2023). These presence data are typically incomplete because not all areas have
54  been systematically searched or even visited (middle part of Figure 1). When presence records are not
55  representative because of unequal detection probabilities, and species are not detected in some areas
56  that they do occupy, this may lead to serious bias in parameter estimates in SDMs (Gu & Swihart,
57  2004). Unequal probabilities of a species being detected at least once in an area, are affected by
58  multiple factors, including variable sampling efforts, which is usually thought of as the number of visits
59  to an area. Here we focus on the opposite: the local probability that a species is not detected and
60  reported during any of the visits in a given year, while the area is known to be occupied (based on
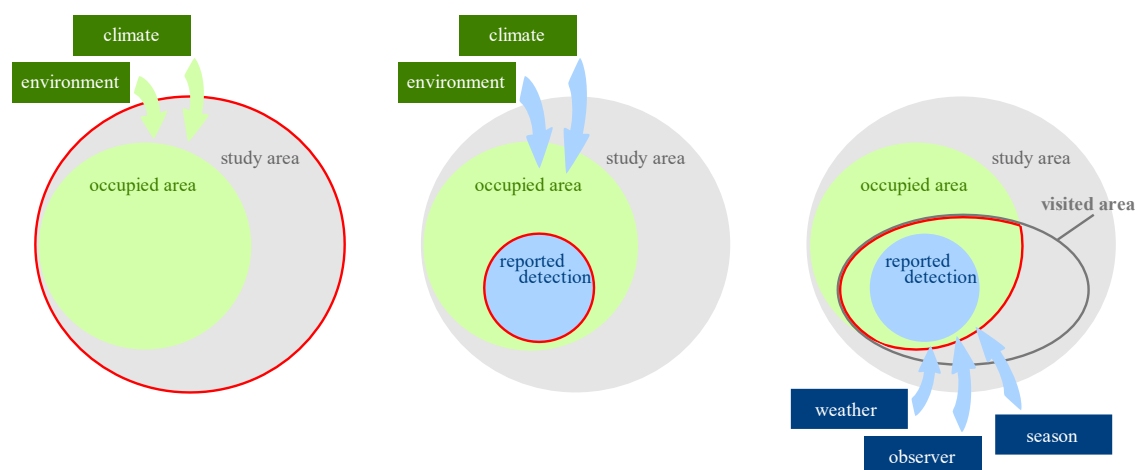61  evidence from previous years).



62

63  Figure 1. Schematic representation of a study area (grey) that is partly occupied by a species (green). Red
64  borders indicate what data are available for the three modeling approaches. Left: ideal situation for fitting a
65  Species Distribution Model with complete cover of the entire study area, resulting in both presence and absence
66  data. Here, the SDM is a function of climate and environmental variables and can generate predictions for areas
67  outside the study area. Mid: often, only presence data are available (blue), only for a subsection of the occupied
68  area. SDMs require assumptions and generation of pseudo-absences. Right: In this manuscript, we focus on the
69  modeling of non-detection probabilities, for which we focus on visits by citizen scientists to areas that are known
70  to be occupied. Data therefore consist of both detections and non-detections that can be related to the
71  frequency and seasonal timing of visits, and the characteristics of the visits like weather and observer
72  characteristics.

73

74

75 Modeling non-detection

76       Occupancy models distinguish between occupancy as the ecological quantity of interest and
77 detection as a consequence of the observation process. They exploit replicate samples in each area to
78 estimate the detection probability, assuming that the occupancy status of the area does not change
79 (Royle & Dorazio, 2008). Many occupancy models incorporate time-varying and site-specific
80 explanatory variables to account for unequal detection probabilities (MacKenzie et al., 2002; Tyre et
81 al., 2003; Bailey et al., 2014). Some authors pointed out that careful sampling design and a strict
82 definition of occupancy are necessary to achieve valid results (Kendall & White, 2009; Efford &
83 Dawson, 2012) and provide insights into how to prioritize efforts to systematically gather data
84 (Sanderlin et al., 2019). Given the lack of any strict definitions for citizen science data, this study will
85 focus on opportunities to achieve valuable insights from more messy data. This is far from trivial,
86 because even when data are gathered systematically, modeling non-detection is a challenging task.

87 The fact that observations are a consequence of both occupancy and detection introduces
88 confounding, together with a lot of free parameters and therefore a need to limit complexity of the
89 model. Research on non-detection has accomplished this by limiting the number of covariates and
90 fitting quite restricted linear models, fixing detection probabilities per species (Kéry and Royle, 2008)
91 or per visit (Outhwaite et al., 2018), applying special estimation procedures (Lele et al., 2012) and
92 filtering species that show low detection probabilities (Ferrer-Paris & Sánchez-Mercado, 2020).

93 Detection and occupancy

94       In the context of citizen science data, detection probabilities are usually very low, well below
95 10 percent. Typically, detection is not the result of a systematic effort but depends on the focus and
96 preference of the observer, who is free to record, or not record, any of the species observed. Detection
97 takes place only if the species is present in the area, observed during a visit, recognized, and reported.
98 Each citizen scientist may have their own definition of an area, therefore areas for modeling need to
99 be constructed in hindsight. Within these constructed areas, we define a visit list as the reported
100 species (at least 1) in an area on a certain date by a certain observer. Defining a visit in this way is
101 crucial because it means that the observer was present in the area on that date and had the chance
102 to report other species that occupy the area. As input for modeling detection, we also define an
103 excursion as the combination of all visits of an observer on one day, to one or multiple areas.

104       Occupancy also needs to be defined systematically. The incomplete overlap between home
105 range and area, together with the unpredictable timing of visits, requires us to rely on cumulative
106 observations of occupancy over time. Efford and Dawson (2012) referred to this as asymptotic
107 occupancy, contrasting it with instantaneous occupancy, where each individual can only occupy a
108 single area. Analogous to that, we will relax the assumption of closure – that there are no changes in
109 occupancy between visits. In this study, we count an area as occupied during the entire year, even if
110 the species is a migratory bird or an insect that can only be found during a specific part of the year.
111 This relaxation carries some implications and risks that will be discussed later. The intended result is
112 that presence of a species is not a necessary consequence of occupancy. Absence during the 'wrong'
113 season can be accounted for by the model, with a zero detection probability. This is in line with the
114 idea that allowing estimates of detection probability to vary between site visits has the potential to
115 'absorb' violations of the closure assumption (MacKenzie et al., 2017). On top of that, citizen scientists
116 have a tendency to overreport interesting observations like a rare species, or a migratory bird that just
117 arrived in the area. With the relaxed occupancy definition, this confounding between closure and
118 detection can be accurately modeled as a seasonal increase in detection probability.

3

119 <u>Conditional models</u>

120       Working with citizen science data and the aforementioned definitions of occupancy, visits and
121 excursions makes modeling non-detection even harder, and increases the need to limit complexity.
122 Therefore, we advocate species-specific models that are targeted at detection only, without
123 considering occupancy. This can be done by conditioning analysis on areas that are known to be
124 occupied. Like Chen et al. (2009) mentioned, inside areas that are known to be occupied, a zero
125 observation is assumed to represent the overlooking of a species rather than its absence. Then non-
126 detection can be directly modelled, without the need to account for possible non-occurrence by use
127 of an occupancy model. This corresponds to the modelling approach in the right part of Figure 1. After
128 the models are thus fitted using occupied areas only, the trained models will predict detection
129 probabilities for all areas given the characteristics of the visits to those areas. Such predictions are also
130 made for unoccupied areas: even completely unsuitable areas, like a land area for a fish, will receive
131 a prediction. Therefore, it is unnecessary to a priori make a distinction between suitable and
132 unsuitable areas.

133       In this paper we develop a model for species-specific non-detection in areas that are known
134 to be occupied by that species, taking the frequency, timing and weather conditions of visits to an
135 area into account as well as the list length and past reports of the observers. We test this non-
136 detection model on citizen science data from the south-west part of the Netherlands. Apart from
137 quantifying sampling efforts, this exercise will reveal species-specific insights into how explanatory
138 variables affect detection probabilities at the visit level, which will allow researchers to map areas with
139 insufficient sampling effort, which should eventually lead to improved species distribution modelling
140 based on citizen science data.

141

142 METHODS

143 <u>Modeling non-detection</u>

144 We developed models to predict non-detection and tested these models using citizen science-based
145 observation data from the Dutch National Database Flora and Fauna (NDFF), mostly from
146 waarneming.nl, in the province of Zeeland from 2017 to 2023. Observations in 2023 are only used as
147 test set. Observations before that are used to train the model and to determine which areas are
148 occupied.

149 The unit of analysis is the unique combination of area, observer and date, which corresponds to a visit
150 to an area during which at least one observation was done. By modeling per visit, specific effects of
151 observer and weather variables can be expected to show up as important determinants of non-
152 detection. In an approach with aggregated data per month or per area, these effects would be missed
153 by averaging over longer periods or multiple observers. Areas have been defined as squared
154 kilometers according to the Amersfoort / RD New (EPSG:28992) grid. All areas with at least 50
155 observations[1] in total in the 2020-2022 period have been included.

156 A total of 150 species were selected from the 10 most observed groups, by taking the most common
157 species within each group: 30 birds, 20 plants, 20 moths, 20 butterflies, 10 mammals, 10 dragonflies,
158 10 Diptera, 10 Hymenoptera, 10 Orthoptera and 10 mushrooms. For each species separately, a
159 modeling dataset is created including only occupied areas: these are defined as areas where the

---

[1] Based on unique observation IDs; during one excursion many observations can be recorded, even of the same species within the same area

species had been observed at least once in two out of three preceding years. For instance, observations in 2017-2019 determined which areas were known to be occupied (i.e., in at least two years) when using the 2020 data for training the model. Similarly, the 2021 and 2022 datasets were also used as training sets (see appendix 1 for an illustration). The training datasets form the basis to create a model that predicts the conditional probability of non-detection per visit in 2023.

The binary modeling target is whether or not the species has been observed during the visit. When a species is not observed, this is seen as a non-detection because only occupied areas have been selected for the modeling dataset. To predict the probability of non-detection, the following features are calculated to represent the circumstances during each visit:

- The list length, which is the number of species observed and reported during the visit
- The visit's day of the year
- Temperature deviation on the day of the visit, compared to the average temperature on the same day of the year over the last thirty years, based on weather station Vlissingen
- For the observer, percentage of all observations before 2023 that were observations of the target species

Other features about weather (precipitation) and observer (average rarity of species observed, total number of observations, total number of species, overall and within the same group as the target species) have been calculated and compared, but were dropped during feature selection. Features about the observer only exist for observers that were already active before 2023, for new observers these features were encoded as -1, which allows the model to treat them as a separate category.

A RandomForestClassifier is trained to learn the relation between these features and the target. Using the trained model, non-detection probabilities are predicted for all visits in 2023. Cross validation is performed to create holdout predictions of probabilities for areas that did not take part in the training. This means that in the test set containing visits in 2023, each predicted probability for a visit is based on a model that has not used any visits to the same area during training. For validation, the Area Under Curve (AUC; Bradley, 1997) per visit is calculated to measure how well the model succeeds in predicting the non-detection of species in areas that are occupied. The AUC metric is a number between 0.0 and 1.0 that measures how well the ordering in the predictions matches the actual binary non-detection. The closer the value is to 1.0, the better the predictions match reality. Random ordering would result in an AUC of 0.5, so the expected range of modeling results is between 0.5 for a worthless model and 1.0 for a perfect model. For each species separately, the conditional non-detection probability P(non-detection|occupancy) per area is calculated as the product of the conditional non-detection probabilities of all visits to the area in 2023, assuming that the probabilities of different visits are independent. This prediction is tested against actual non-detection during all visits to that area combined, by calculating the AUC per area. The result is compared to two benchmarks for predicting detection: one based on only the number of visits and one that also takes the list length per visit into account.

Feature selection was an iterative process, involving both feature definition and feature selection. Guiding principles were: First, to avoid multicollinearity by standardizing features. For example, temperature has been expressed as a deviation and observations of the target species as a percentage. Second, features were selected in a forward fashion, starting with one feature per category (weather, observer and season), only adding more features when the improvement of the AUC per area was at least 0.01. For each selected feature, this improvement is reported as the contribution to the final model, calculated by running the model again without the feature to measure the difference in AUC per area, averaged over all species.

5

205    To create marginal dependence plots, the detection probability is calculated on a rolling window for
206    each decile of the input features. In contrast to partial dependence plots from the Random Forest
207    Classifier, these plots do not take confounding and interaction between features into account. Feature
208    definition and selection were aimed at reducing differences between partial and marginal
209    dependence, resulting in dependence plots that are independent of the model and feature selection,
210    but still representative of the relations that the model has learnt.

211

212    RESULTS

213    To illustrate model performance for a variety of species groups and different types of observations,
214    results are shown in detail for four of the 150 species: a coastal bird species (the Eurasian
215    oystercatcher *Haematopus ostralegus*), the most observed plant in the dataset (the ribwort plantain
216    *Plantago lanceolata*), a moth that is mainly caught using special traps (the bright-line brown-eye
217    *Lacanobia oleracea*) and a common mushroom (the shaggy ink cap *Coprinus comatus*).

218    A total of 2179 areas of one squared kilometer had at least 50 observations between 2020 and 2022
219    and were selected for testing the model. The number of observations in these areas varied between
220    50 and 38164 and had a median of 281.

221    For each of the 150 species, areas for testing the model were selected based on occupancy, defined
222    as an observation in at least two out of three years from 2020 to 2022 (see red squares in Fig. 2). One
223    species, the Asian Hornet, was filtered out because not one area met the occupancy definition. For
224    the remaining 149 species, the number of occupied areas varied between 16 and 1561 and had a
225    median of 147. An observation during 2023 took place if any of the visits to the area resulted in
226    detection of the species. Among occupied areas, the percentage of observations per species in 2023
227    ranged from 11% for the big sheath mushroom *Volvopluteus gloiocephalus* to 86% for the hornet
228    mimic hoverfly *Volucella zonaria* (median 59%). Outside occupied areas, this percentage ranged from
229    0.4% for the fly agaric *Amanita muscaria* to 35% for the common buzzard *Buteo buteo* (median 8%).
230    Observations thus had a higher probability inside previously occupied areas, implying a correlation
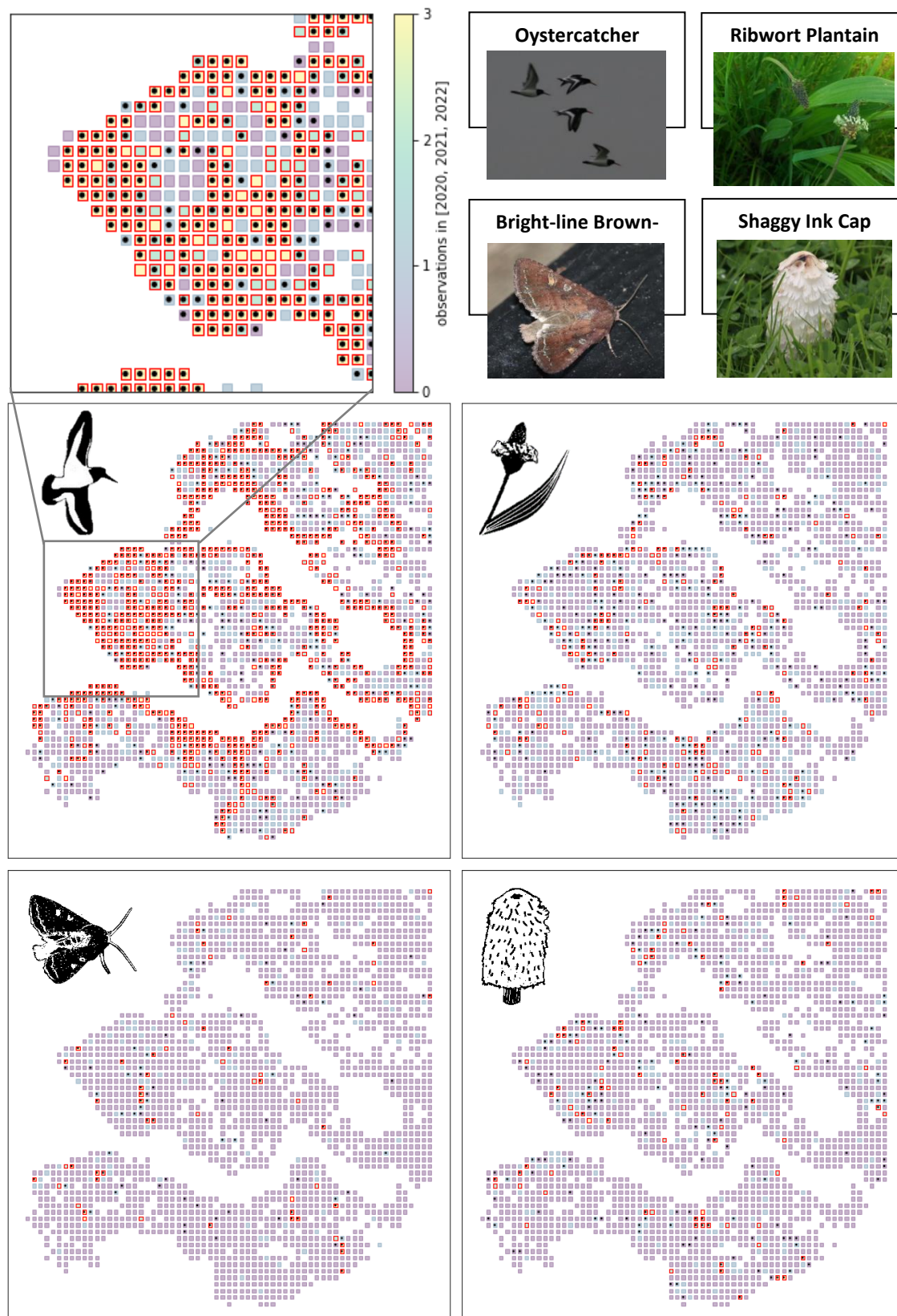231    between occupation based on 2020 to 2022 and observation in 2023.

Figure 2. Selected areas for the modeling dataset (red squares) and observations in 2023 (black dots) for the selected species. Background color of the areas represents the number of years from 2020 to 2022 that the species was observed. Only areas where the species was observed in at least two out of three years are selected for modeling.

237

238    A Random Forest Classifier was trained per species, to learn the relation between features of each
239    visit and observations in 2020, 2021 and 2022. The dependence plots offer intriguing insights into the
240    factors influencing detection and, consequently, shed light on underlying sources of bias. List length
241    is the strongest feature individually and improves the average AUC of the final model by 0.037. For all
242    species, the detection probability increases with the list length (first column in Fig. 3). This makes
243    sense as a more comprehensive effort during the visit leads to a larger probability that the focal
244    species is detected and reported. The detection probability of bird and moth species increases faster
245    with list length, which is consistent with the fact that the numbers of common birds and moths were
246    smaller than that of plants and mushrooms, resulting in larger proportions per species. The
247    oystercatcher and the bright-line brown-eye had even steeper curves than the average bird and moth,
248    because both were among the most observed species. Interestingly, the shaggy ink cap shows a flatter
249    curve that even seems to decrease when the list length increases. This may be related to the ecology
250    of the shaggy ink cap: it grows in grasslands that are typically not very species-rich. List length tended
251    to be largest for moth species, probably as a consequence of the strategy to capture moths using a
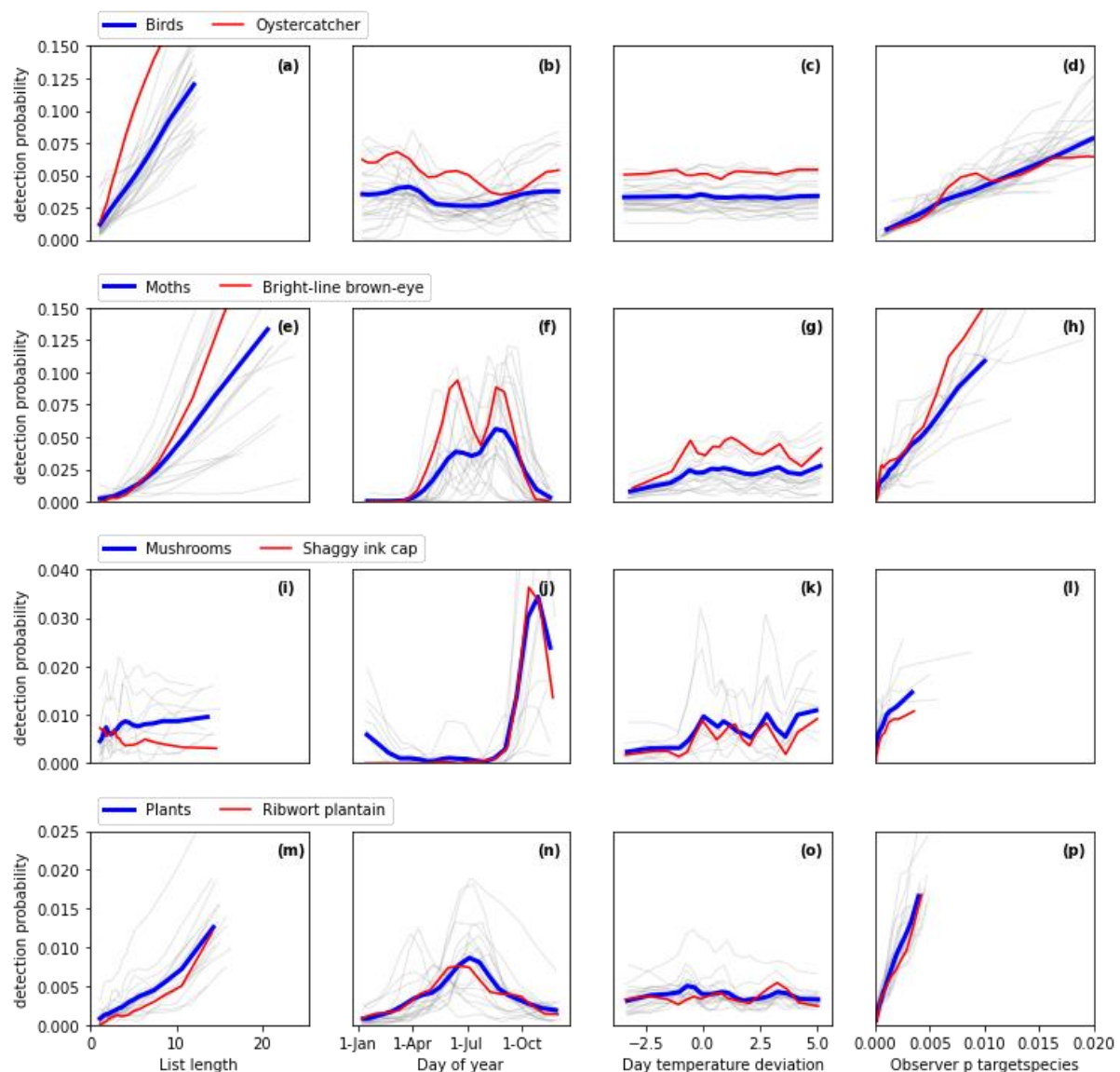252    light trap that detects all species at a single location.

253

Figure 3. Marginal dependence plots, showing the conditional probability of detection given occupancy per species for each visit, as a function of the selected features, per species group. In grey, the individual curves of all species. In blue, the average curve over all species. In red, the four species that are covered in more detail throughout the results.

Each species group had its own typical distribution of detection probability over the year (second column in Fig. 3). The "day of year" feature strongly contributes to the final model by improving the average AUC by 0.020. Mushrooms had the highest probability of detection during visits in autumn, while moths and plants had their peak during summer. On average, birds had more or less constant detection probabilities throughout the year. Individual bird species do show curves as expected, dependent on the period that the bird usually stays in the province. Migratory birds like the barn swallow *Hironda rustica*, the western marsh harrier *Circus aeruginosus* and the common chiffchaff *Phylloscopus collybita* have their peaks during the warm season, while most other birds have higher detection probabilities during winter. Like most moths, the bright-line brown-eye is most often detected during the summer, and even two separate generations (Vajgand, 2009; Vlinderstichting, 2024) are reflected in the detection probabilities. Some plants had their peak in spring instead of in summer (see grey lines in Fig. 3n), most notably cow parsley *Anthriscus sylvestris*, ground-ivy

271 *Glechoma hederacea* and the red dead-nettle *Lamium purpureum*, this seems to be related to the
272 month when these plants (first) bloom.

273 The effect of temperature deviations in the third column of Fig. 3 was generally weak. This feature
274 improved the average AUC of the final model by only 0.004 and was only selected because it is the
275 strongest one compared to other weather features. However, for moths there is a trend with
276 increased detection probabilities when temperatures were higher than the 30-year mean for that time
277 of the year. To a lesser extent, this also seemed to be the case for mushrooms. Note that the
278 temperature deviation could also determine how many visits took place. Our analysis does not look
279 into that because detection probability is calculated only on the occasions when a visit did take place.

280 All species tend to have a higher detection probability when the observer had detected and reported
281 the focal species often during visits in earlier years (fourth column in Fig. 3). This feature assists the
282 model to distinguish between observers with a different focus, as some observers only report birds
283 while others are mainly focused on butterflies. Among all observer features, the selected feature
284 about detecting the same species is individually the strongest predictor. It is also the second most
285 important feature in the final model, improving the average AUC by 0.034. Compared to other
286 observer features, it shows the most consistent and interpretable dependence plots, in the sense that
287 it shows about the same pattern for all species. Other features about the observer history were
288 created and tested in addition to the final model (see Appendix 2), but none of the added features
289 improved the average AUC by more than 0.01.

290 Validation

291 The AUC has been calculated for non-detection per visit, by testing predicted non-detection
292 probabilities against actual non-detection during the visit. Over all species, this AUC per visit varies
293 between 0.47 and 0.97 and has a median of 0.77. These outcomes can be compared to an AUC of 0.5
294 that would reflect a random probability per visit, so they prove that the model is quite successful in
295 predicting visits that result in a detection.

296 The non-detection probability of an area (conditional on being occupied in previous years) is the
297 product of the conditional non-detection probabilities of all visits to the area in 2023. Over all species,
298 the resulting AUC per area varied between 0.48 and 0.96 and had a median of 0.77. These results are
299 compared to two benchmarks that are simpler because they do not use all available information. The
300 first benchmark takes only the number of visits into account. For this benchmark, areas are ordered
301 only by the number of visits, expecting a larger detection probability in areas that were visited more
302 often. This benchmark results in a median AUC of 0.69. The second benchmark applies the same
303 modeling approach, but only uses the list length feature, disregarding information about weather,
304 observer and season. This more competitive benchmark results in a median AUC of 0.73. Instead of
305 looking at the median AUC over all species, results were also compared within each species group
306 separately. Figure 4 shows that the predictions of the full model consistently outperform both
307 benchmarks for all the species groups, although the difference is less clear for Diptera.
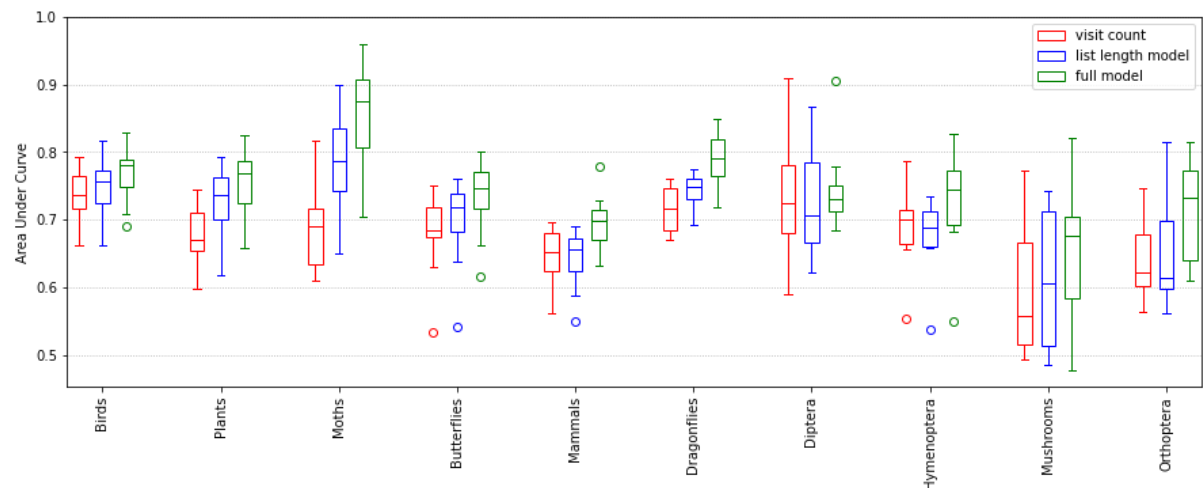
308



309

310 Figure 4. Boxplots summarizing the Area Under Curve (AUC) per species group, given three different
311 approaches to rank the areas. Each box contains the middle 50% between the first and third quartile
312 of the AUC outcomes for an approach carried out on a species group (for example, the red box for the
313 visit count approach applied to birds on the left). The horizontal line in the middle of the box is the
314 median, whiskers indicate the full range between minimum and maximum, with the exception of
315 outliers, indicated with circle markers outside this range. The AUC is expected to vary between 0.5 for
316 a worthless model and 1.0 for a perfect model. The figure shows that the full model beats both
317 benchmark approaches within each species group.

318 The non-detection probabilities were averaged per species to compare them to the proportion of
319 areas where the species was observed at least once in the test year 2023. The average was calculated
320 for occupied (at least 2 years in 2020 to 2022) and 'unoccupied' areas separately. As expected, inside
321 occupied areas there is a strong negative relation between non-detection probability and observations
322 (blue dots in Figure 5). The dashed grey line shows the ideal relation for a perfect model. The results
323 were scattered around it, but the observation percentage is a bit lower than expected when the non-
324 detection probability gets lower and the variance seems to increase towards non-detection
325 probabilities around 50%. In areas that were not occupied, the observation percentage was generally
326 lower and not as strongly related to the non-detection probability: the reason for not observing a
327 species is (at least partly) that the area was not occupied in 2023.
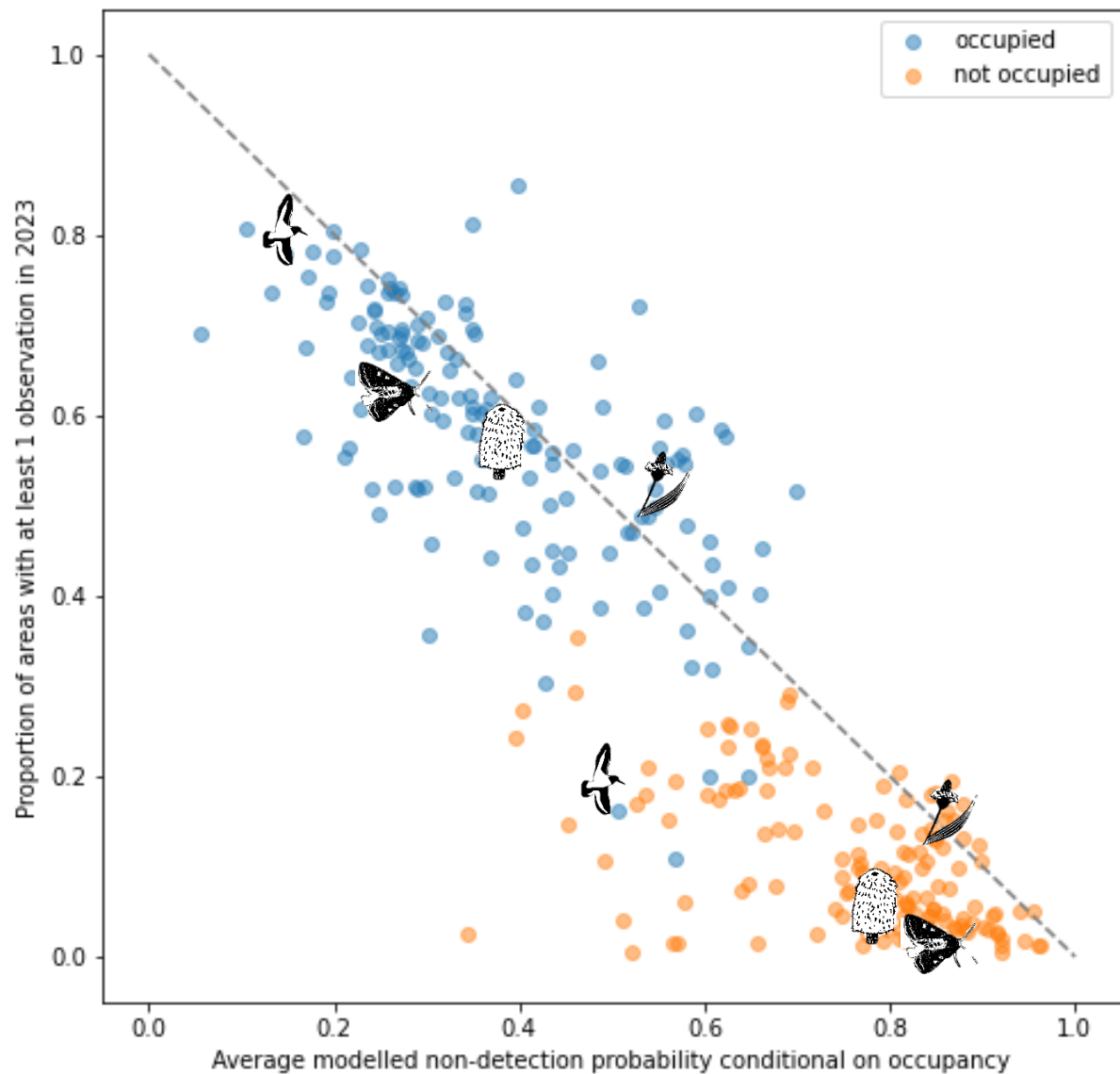
328

Figure 5. The relation between the average non-detection probability per species conditional on occupancy in previous years (horizontal axis) and the proportion of areas with an observation in the test year 2023 (vertical axis). The blue dots represent species-specific averages for areas occupied in at least 2 years in the 2020-2022 period. Model predictions for conditional non-detection probabilities were also calculated for the other (previously 'unoccupied') areas, which had fewer observations in 2023 and showed a weaker relationship with the proportion of observations (orange dots). In 'unoccupied areas' (orange dots) the ribwort plantain has relatively high proportion of observations, suggesting that a lack of detection may be the sole reason for the unoccupied status. Both the bright-line brown-eye and the shaggy ink cap show similar (high) non-detection probabilities, but a lower proportion of observations, meaning that the small effort led to even fewer observations, suggesting that the occupation status has a correlation with true occupancy. The same applies to the Eurasian oystercatcher, in a more conclusive manner because even in unoccupied areas there has been a considerable effort, given the non-detection probability around 50%, resulting in only 20% observations. In occupied areas, all four species have a proportion of observations that is in line with the complement of the non-detection probability, as expected.

345 The probability of non-detection (conditional on occupancy) that the model reveals, quantifies the
346 sampling effort. In unoccupied areas, this probability is not the actual probability to not detect a
347 species, but only the probability that a species would not be detected if it were present, given the
348 sampling effort. Therefore, the probability of detection is a suitable quantification of (the lack of)
349 sampling effort, both for occupied and unoccupied areas. Darker areas in figure 6 represent lower
350 non-detection probabilities, and had a more complete sampling effort. Given all visits by citizen
351 scientists in 2023, the sampling effort for oystercatchers turned out to be the most comprehensive
352 one. The spatial distribution of the sampling effort looks about the same for these four species. As
353 expected, the effort is clustered in space and highest in densely populated areas like the cities of Goes
354 and Middelburg, and areas with a high biodiversity like areas along the coast and nature reserves that
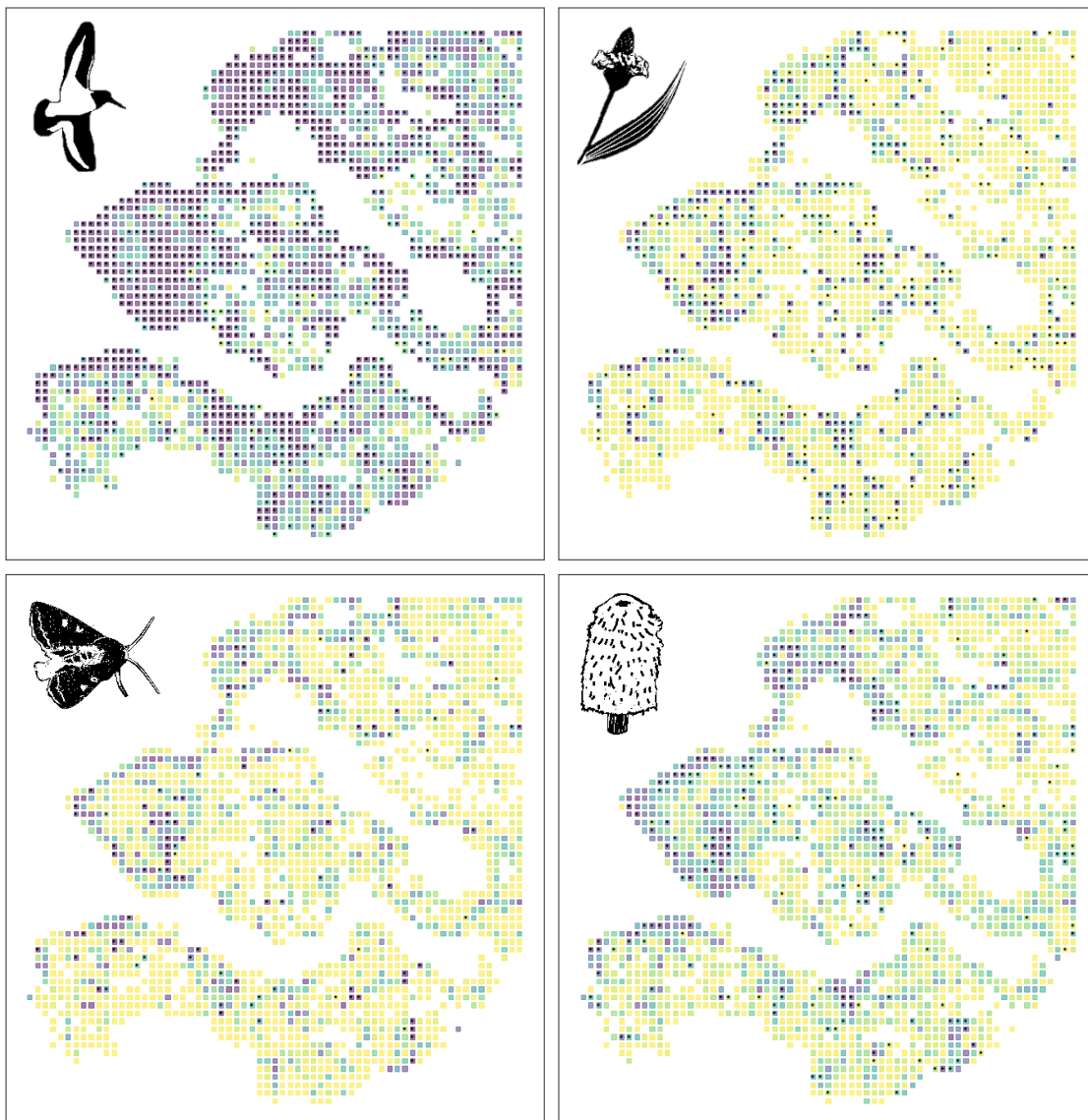355 attract many observers.



356

357 Figure 6. Per squared kilometer area, the probability of non-detection given occupancy, with
358 observations in 2023 as black dots. All areas are shown, because even in (previously) unoccupied areas
359 this probability quantifies the sampling effort. The brightest (yellow) color represents 100% non-
360 detection probability, all four selected species have non-detection probabilities ranging from 0.00%
361 to 100% with an average of 35.8% for the Eurasian oystercatcher, 82.9% for the ribwort plantain,
362 80.9% for the bright-line brown-eye and 74.3% for the shaggy ink cap.

13

363

364 DISCUSSION

365 Key results

366 Citizen science data about biodiversity are a growing source of information that is often left unused
367 because it is hard to derive representative insights from observations that do not follow a known
368 protocol. The fact that only presence data are recorded is challenging because it is unknown what lead
369 the observer to select some and not other species for reporting. The method proposed in this paper
370 gives insights into causes of bias, the non-linear effects of these causes and their relative importance.
371 Based on input features about the visit, the season, the weather conditions and the observer, a model
372 predicts the probability of non-detection when the area is occupied by the species of interest.
373 Validation results confirm that the model succeeds in predicting the non-detection probability based
374 on the available features.

375 The model and its predictions focus on detection only, because the model has been fit on a subset of
376 areas that are known to be occupied by the species. Therefore, occupancy does not play a role in the
377 relations that the model finds. This implies that the predictions are suitable to quantify sampling
378 effort. This effort is not just based on the number of visits, but also takes circumstances during the
379 visit that have an impact on detection probability into account.

380 By fitting models for detection inside occupied areas only, complexity has been considerably reduced.
381 In contrast to other studies, we did not need to limit the number of covariates, assume linear
382 relationships, fix detection probabilities per species (Kéry and Royle, 2008) or per visit (Outhwaite et
383 al., 2018) or filter species that show low detection probabilities (Ferrer-Paris & Sánchez-Mercado,
384 2020). Moreover, by using external data to define occupancy, we avoided the reliance on a strict
385 definition of occupancy to estimate the detection probability from replicate samples.

386 Selection of occupied areas

387 Occupied areas need to be selected in order to fit the model. It is required that occupancy is static
388 during the time period of the modeling dataset. This selection can be challenging for species that are
389 fluctuating, expanding or declining in range. Because of this challenge, occupied areas for the rapidly
390 expanding Asian hornet could not be established based on historical records.

391 The modeling approach allows a very loose definition of occupancy. Occupancy does not need to mean
392 that a species is present at the time of the visit, like regularly assumed (Efford & Dawson, 2012). When
393 a bird of prey went hunting somewhere else, a butterfly is unfindable because it is a pupa or a
394 migratory bird is spending the winter in Africa, their area would still count as occupied. Given the
395 features of the visit, the model will soon recognize that the detection probability is lower or zero
396 during a certain season or weather condition. Within occupied areas, all complexity in patterns of
397 species presence is taken care of by the model. The model could even find interactions between
398 species behavior and observer habits, for example an observer who is more likely to detect a migratory
399 bird when it first returns to its breeding area. A different type of occupation among areas within the
400 study area is challenging, however. For instance, many oystercatchers occupy the tidal areas in winter,
401 while in the breeding season they are spread out from the coast to far inland (Allen et al., 2019). A
402 pattern like that can confuse the model, because all areas share the same parameters and seasonal
403 patterns are assumed to be the same among areas.

404 Selection of areas that are actually not occupied may introduce bias in the modeling results. The
405 dataset would then contain visits with non-detection that are not caused by detection, but by (the

406 lack of) occupancy. That would lead to overestimation of the non-detection probabilities. Treating
407 unoccupied areas as occupied areas could be a consequence of misidentifications (Miller et al., 2011),
408 but otherwise it is not likely given our approach: the relaxed definition of occupancy includes
409 occasional visitors, and could be made more stringent with the definition of 'previously occupied', e.g.
410 reported in all three previous years, or multiple times per year. The opposite, missing occupied areas
411 for the selection, would not need to be a problem. The only assumption required to have realistic
412 results is that the areas in the dataset are representative of all the occupied areas.

413 In this study, the selection of occupied areas is a consequence of detection in earlier years. When the
414 distribution of sampling effort over areas is very skewed and persistent over years, the
415 representativity assumption may be violated. Moths, for example, are often caught using light traps
416 that are only applied at certain sampling locations for years in a row. Within the selected areas that
417 are occupied by moths, the detection probability is quite high. In other areas that are not part of the
418 modeling dataset, the detection probability is much lower and will be overestimated by the model
419 that is mostly trained using areas with light traps. The challenge is to add features about observer or
420 area that will help the model to recognize the difference between areas with and without light traps.
421 The current feature about the percent of target species in the observer history does help to partly
422 relieve this problem, because it distinguishes visits of observers who previously reported the same
423 moth species – these are the observers involved with light traps and likely to find and report moths
424 even when they become active in new areas. Attempts to define features at the area level, like the
425 percent of observations in the species group of the focal species, did not improve the results and
426 introduced the challenge to avoid that the model learns to distinguish between areas instead of area
427 characteristics.

428 Another option to select occupied areas would be to use an external dataset that functions as the gold
429 standard. For example, in parts of this study area, shorebirds are systematically counted multiple
430 times a year by professionals following a standardized protocol (Lilipaly & Sluijter, 2023). Using these
431 data to define occupancy by shorebird species would lower the risk of systematically selecting more
432 areas as occupied when they have a higher sampling effort that is persistent over years, because the
433 protocol makes sure that the sampling effort is kept constant. An additional advantage is that
434 occupancy could be established during the same time period, instead of the history within the same
435 dataset. But there are also additional challenges when using an external dataset: the resolution in
436 time and space may not match well and there may be possible overlap between datasets because
437 observers submit their observations to multiple platforms.

438 When observations in the test year 2023 were compared between occupied and previously
439 unoccupied areas, it turned out that species were more often found in the same areas where they
440 were previously found. This correlation can be caused by occupancy itself: a previously occupied area
441 has a larger probability of being occupied again. But an alternative explanation might be the sampling
442 effort: areas that are consistently visited more often have a higher probability of observation, both in
443 2020 to 2022 and in 2023. Figure 5 shows that both things may play a role and that the model helps
444 to distinguish between the two. Regarding the sampling effort, the average non-detection probability
445 tends to be higher in unoccupied areas: this confirms that some of the unoccupied areas are not well
446 investigated (with regard to the focal species) in 2023, which makes it more likely that their
447 unoccupied status based on 2020 to 2022 has been a matter of sampling effort. On the other hand,
448 there is a clear distinction between occupied and previously unoccupied areas that have the same
449 average non-detection probability: given the same sampling effort, observations in 2023 are more
450 likely in areas that were occupied in at least two years in the 2020-2022 period.

451

15

452 <u>Extensions</u>

453 All calculations and models applied in this paper are species-specific. Therefore, attributes of species
454 are constant values and not suitable as features for modeling. Some attributes like size, visibility,
455 recognizability, species group and rarity are likely to have an impact on detection. In the current
456 approach, this will affect only the average level of predictions per species. An interesting extension
457 would be a meta-model to describe and predict how attributes of species influence the detection
458 probability. This would be particularly useful to create predictions for species that have very few
459 records of observations.

460 The modeling of non-detection facilitates improved information about the distribution of species, by
461 estimating and correcting how various features cause bias. The steep increase in the number of
462 observers that participate in citizen science makes it even more challenging to estimate trends over
463 time (Fink et al., 2023). In principle, the model could recognize changes in detection probability as a
464 consequence of more observers and observers with different characteristics. It would require
465 additional feature extraction and validation to extend the model to apply it for estimation of trends.

466 In practice, abundance of species will have an impact on the detection probability. In the dataset for
467 this study, abundance estimates are not standardized and often missing. Therefore, species presence
468 has been simplified to occupancy only, disregarding abundance. It would be interesting to extend the
469 model to take abundance into account as a feature that influences the detection probability. This may
470 be a complicated puzzle because occupancy is also used to select areas for modeling. Incomplete
471 information in existing records is not uncommon in citizen science data. Other potentially interesting
472 features that have not been used because of incomplete data include time of day and life stage of the
473 observed specimen.

474 This study does not look into the causes of relations between visit features and detection probability.
475 The model learns the relations between input features and detection to produce a prediction, without
476 necessarily understanding what the relations mean. Some of the relations may be based on artefacts
477 in the data that do not actually help to correct bias. Better understanding of the relations will help to
478 create more accurate models. This requires both expert insights and investigation of how the model
479 creates its predictions.

480 Finally, our detection model only focuses on non-detection bias, without considering misidentification
481 that could lead to false positives. Miller et al (2011) show that uncertain detections can be used to
482 improve occupancy estimates. For the approach in this study, uncertain detections based on indirect
483 observation (scat or tracks) would not be problematic because the definition of occupancy does not
484 require instantaneous occupancy. Uncertain detections because of unskilled citizen scientists would
485 usually be detected by validators who check feasibility and evidence in the form of photos that are
486 uploaded.

487

488 FURTHER RESEARCH

489 The bias in citizen science data is a consequence of the specific interests and habits of observers. For
490 direct inference about occupancy, the entangled causes of occupancy and detection need to be
491 unraveled. To achieve that in the context of citizen science data, both limitations of model complexity
492 and additional assumptions would be required. By focusing on detection within occupied areas, we
493 deliberately took an observation-driven view (Royle & Dorazio, 2008). The consequence is that our
494 model does not allow for direct inference about occupancy. The quantification of sampling effort

495  needs to be combined with other information to arrive at conclusions about occupancy, the quantity
496  of interest. For many species, this kind of information is available in datasets that do contain
497  systematically gathered data, but only in a relatively small subset of areas. Koshkina et al. (2017) prove
498  that a considerable improvement can be attained by combining citizen science and systematically
499  gathered data. They combine both sources in an integrated model, that is quite complex and has issues
500  to correctly identify different sets of covariates. Following the approach presented in this paper, it
501  would be quite natural to have separate models to exploit the information from these sources before
502  it is combined.

### Combining conditional models

504  The conditional model for detection purposely neglected occupancy by focusing on occupied areas
505  only. To model suitability of the areas for occupancy, we suggest fitting separate SDMs (Elith &
506  Leathwick, 2009) using explanatory variables that represent the physical environment and the local
507  climate. Ideally, these models are fitted using some set of areas where complete presence and
508  absence data are available, in a way that the assumption of perfect detection makes sense. This
509  corresponds to the model displayed in the left part of Figure 1. In contrast, the detection models
510  discussed in this paper only use visit specific explanatory variables like weather, observer and season.
511  Although some authors argue and prove that the physical environment does have an impact on
512  detection (Gu & Swihart, 2004; Chen et al., 2009) there is a risk that these variables would mainly be
513  related to abundance as a cause of detection. For example, think of a model that uses the number of
514  trees in an area as an input to predict detection of a bird species that lives in trees. It is likely that
515  more trees would provide more hiding places for the birds, making detection harder – which would
516  be a reason to include the number of trees in the model. But it is also quite likely that the number of
517  trees would change the abundance of the tree-living bird species in that area. The model simplifies
518  presence of the bird to occupancy, so this increase in abundance would not change the representation
519  of the bird in the data. But abundance of the bird might very well change the likelihood of detection,
520  introducing a relation between environmental characteristics and detection that is mediated by
521  abundance and opposite in direction to the direct effect of trees on detection probability. Including
522  environmental characteristics to predict detection might improve predictions of detection probability,
523  but it contradicts both the choice to estimate occupancy instead of abundance and the choice to build
524  separate conditional models for detection and occupancy.

### Correcting bias

526  A promising application of the conditional model presented in this paper is the potential to correct
527  bias in citizen science data. To do so, the SDMs in the previous paragraph need to be combined to the
528  detection models we presented. First, suitability-based probabilities of occupancy can be estimated
529  for any area, based on SDMs that take climate and environment into account. Next, citizen science
530  data can be used to update these expectations. In areas without presence records, non-detection
531  probabilities according to the detection model provide information about the sampling effort, based
532  on visit specific explanatory variables. Non-detection inside an occupied area is quite likely if the
533  sampling effort has been small, but becomes less likely with an increasing sampling effort. This gives
534  an excellent opportunity to combine the information that is available in systematically gathered data
535  on one hand and citizen science on the other hand. Probabilities of occupancy can be estimated for
536  all areas by applying the SDMs that were learnt in a subset of areas where complete presence absence
537  data are available from systematic research. These estimates are based on environmental
538  characteristics and therefore represent the suitability of the area for occupancy. To estimate the
539  distribution of a species, this suitability needs to be translated into an occupancy score. This could be
540  done using a fixed threshold (Liu et al., 2005) but it would be more sophisticated to take into account

541 how well the area has been investigated. If the sampling effort has been very substantial in an area
542 without detecting the species, it is less likely that the area is occupied, even when it is very suitable
543 given the environmental characteristics. That is where citizen science data can be combined to
544 improve the estimates for areas outside the systematically monitored areas. The sampling effort,
545 quantified by taking the circumstances during all visits into account, tells how well the area has been
546 investigated – and how likely it is that the area is occupied when a species has not been reported after
547 this effort took place. To summarize this, the suitability-based probability of occupancy needs to be
548 adjusted downward when citizen science data do not contain presence records while there has been
549 a substantial sampling effort.

550 This approach to correct bias in occupancy estimates based on citizen science can be formalized using
551 a mathematical theorem known as Bayes' rule (Bayes, 1763). To estimate the posterior probability of
552 occupancy in areas where presence has not been ascertained, this rule updates the suitability-based
553 prior probability of occupancy using the conditional probability of non-detection given occupancy. The
554 resulting posterior probability of occupancy given non-detection is suitable to draw conclusions about
555 occupancy in areas where the species has not been detected. The posterior probabilities can be
556 validated by comparing them to external data about occupancy.

557

558 REFERENCES

559 Allen, A. M., Ens, B. J., van de Pol, M., van der Jeugd, H., Frauendorf, M., Oosterbeek, K., & Jongejans,
560 E. (2019). Seasonal survival and migratory connectivity of the Eurasian Oystercatcher revealed by
561 citizen science. *The Auk: Ornithological Advances*, 136(1), uky001.

562 Bailey, L. L., Simons, T. R., & Pollock, K. H. (2004). Comparing population size estimators for
563 plethodontid salamanders. *Journal of Herpetology*, *38*(3), 370-380.

564 Bailey, L. L., MacKenzie, D. I., & Nichols, J. D. (2014). Advances and applications of occupancy
565 models. *Methods in Ecology and Evolution*, *5*(12), 1269-1279.

566 Bayes, T. (1763). LII. An essay towards solving a problem in the doctrine of chances. By the late Rev.
567 Mr. Bayes, FRS communicated by Mr. Price, in a letter to John Canton, AMFR S. *Philosophical*
568 *transactions of the Royal Society of London*, (53), 370-418.

569 Bradley, A. P. (1997). The use of the area under the ROC curve in the evaluation of machine learning
570 algorithms. *Pattern Recognition*, *30*(7), 1145-1159.

571 Chen, G., Kery, M., Zhang, J., & Ma, K. (2009). Factors affecting detection probability in plant
572 distribution studies. *Journal of Ecology*, *97*(6), 1383-1389.

573 Efford, M. G., & Dawson, D. K. (2012). Occupancy in continuous habitat. *Ecosphere*, *3*(4), 1-15.

574 Elith, J., & Leathwick, J. R. (2009). Species distribution models: ecological explanation and prediction
575 across space and time. *Annual Review of Ecology, Evolution, and Systematics*, *40*, 677-697.

576 Fink, D., Johnston, A., Strimas-Mackey, M., Auer, T., Hochachka, W. M., Ligocki, S., ... & Rodewald, A.
577 D. (2023). A double machine learning trend model for citizen science data. *Methods in Ecology and*
578 *Evolution*, *14*(9), 2435-2448.

579 Ferrer-Paris, J. R., & Sánchez-Mercado, A. D. A. (2020). Making inferences about non-detection
580 observations to improve occurrence predictions in Venezuelan Psittacidae. *Bird Conservation*
581 *International*, *30*(3), 406-422.

18

Gu, W., & Swihart, R. K. (2004). Absent or undetected? Effects of non-detection of species occurrence on wildlife–habitat models. *Biological Conservation*, *116*(2), 195-203.

Hochkirch, A., Samways, M. J., Gerlach, J., Böhm, M., Williams, P., Cardoso, P., ... & Dijkstra, K. D. B. (2021). A strategy for the next decade to address data deficiency in neglected biodiversity. *Conservation Biology*, *35*(2), 502-509.

Isaac, N. J., van Strien, A. J., August, T. A., de Zeeuw, M. P., & Roy, D. B. (2014). Statistics for citizen science: extracting signals of change from noisy ecological data. *Methods in Ecology and Evolution*, *5*(10), 1052-1060.

Isbell, F., Craven, D., Connolly, J., Loreau, M., Schmid, B., Beierkuhnlein, C., ... & Eisenhauer, N. (2015). Biodiversity increases the resistance of ecosystem productivity to climate extremes. *Nature*, 526(7574), 574-577.

Jha, A., Praveen, J., & Nameer, P. O. (2022). Contrasting occupancy models with presence-only models: Does accounting for detection lead to better predictions? *Ecological Modelling*, *472*, 110105.

Johnston, A., Matechou, E., & Dennis, E. B. (2023). Outstanding challenges and future directions for biodiversity monitoring using citizen science data. *Methods in Ecology and Evolution*, 14(1), 103-116.

Kendall, W. L., & White, G. C. (2009). A cautionary note on substituting spatial subunits for repeated temporal sampling in studies of site occupancy. *Journal of Applied Ecology*, *46*(6), 1182-1188.

Kéry, M., & Royle, J. A. (2008). Hierarchical Bayes estimation of species richness and occupancy in spatially replicated surveys. *Journal of Applied Ecology*, *45*(2), 589-598.

Koshkina, V., Wang, Y., Gordon, A., Dorazio, R. M., White, M., & Stone, L. (2017). Integrated species distribution models: combining presence-background data and site-occupancy data with imperfect detection. *Methods in Ecology and Evolution*, *8*(4), 420-430.

Lele, S. R., Moreno, M., & Bayne, E. (2012). Dealing with detection error in site occupancy surveys: what can we do with a single survey? *Journal of Plant Ecology*, *5*(1), 22-31.

Lilipaly S.J. & M. Sluijter (2023). Kustbroedvogels in het Deltagebied in 2022. Rijkswaterstaat, Centrale informatievoorziening Rapport BM 23.04. Deltamilieu Projecten Rapportnr. 2023-05, Vlissingen.

Liu, C., Berry, P. M., Dawson, T. P., & Pearson, R. G. (2005). Selecting thresholds of occurrence in the prediction of species distributions. *Ecography*, 28(3), 385-393.

Luna, S., Gold, M., Albert, A., Ceccaroni, L., Claramunt, B., Danylo, O., ... & Sturm, U. (2018). Developing mobile applications for environmental and biodiversity citizen science: considerations and recommendations. *Multimedia Tools and Applications for Environmental & Biodiversity Informatics*, 9-30.

MacKenzie, D. I., Nichols, J. D., Lachman, G. B., Droege, S., Andrew Royle, J., & Langtimm, C. A. (2002). Estimating site occupancy rates when detection probabilities are less than one. *Ecology*, *83*(8), 2248-2255.

MacKenzie, D. I., Nichols, J. D., Royle, J. A., Pollock, K. H., Bailey, L., & Hines, J. E. (2017). Occupancy estimation and modeling: inferring patterns and dynamics of species occurrence. Elsevier.

619    Melo-Merino, S. M., Reyes-Bonilla, H., & Lira-Noriega, A. (2020). Ecological niche models and species
620    distribution models in marine environments: A literature review and spatial analysis of evidence.
621    *Ecological Modelling*, 415, 108837.

622    Miller, D. A., Nichols, J. D., McClintock, B. T., Grant, E. H. C., Bailey, L. L., & Weir, L. A. (2011). Improving
623    occupancy estimation when two types of observational error occur: Non-detection and species
624    misidentification. *Ecology*, *92*(7), 1422-1428.

625    Oliver, T. H., Heard, M. S., Isaac, N. J., Roy, D. B., Procter, D., Eigenbrod, F., ... & Bullock, J. M. (2015).
626    Biodiversity and resilience of ecosystem functions. *Trends in Ecology & Evolution*, 30(11), 673-684.

627    Outhwaite, C. L., Chandler, R. E., Powney, G. D., Collen, B., Gregory, R. D., & Isaac, N. J. (2018). Prior
628    specification in Bayesian occupancy modelling improves analysis of species occurrence
629    data. *Ecological Indicators*, *93*, 333-343.

630    Ranc, N., Santini, L., Rondinini, C., Boitani, L., Poitevin, F., Angerbjörn, A., & Maiorano, L. (2017).
631    Performance tradeoffs in target-group bias correction for species distribution
632    models. *Ecography*, *40*(9), 1076-1087.

633    Royle, J. A., & Dorazio, R. M. (2008). Hierarchical modeling and inference in ecology: the analysis of
634    data from populations, metapopulations and communities. Elsevier.

635    Sanderlin, J. S., Block, W. M., Strohmeyer, B. E., Saab, V. A., & Ganey, J. L. (2019). Precision gain versus
636    effort with joint models using detection/non-detection and banding data. *Ecology and Evolution*, *9*(2),
637    804-817.

638    Tyre, A. J., Tenhumberg, B., Field, S. A., Niejalke, D., Parris, K., & Possingham, H. P. (2003). Improving
639    precision and reducing bias in biological surveys: estimating false-negative error rates. *Ecological
640    Applications*, *13*(6), 1790-1801.

641    Vajgand, D. R. A. G. A. N. (2009). Flight dynamic of economically important Lepidoptera in Sombor
642    (Serbia) in 2009 and forecast for 2010. *Acta Entomologica Serbica*, 14(2), 175-184.

643    Van Strien, A. J., Van Swaay, C. A., & Termaat, T. (2013). Opportunistic citizen science data of animal
644    species produce reliable estimates of distribution trends if analysed with occupancy models. *Journal
645    of Applied Ecology*, *50*(6), 1450-1458.

646    Vlinderstichting    (2024,    July    11).    Groente-uil    *Lacanobia    oleracea*.
647    https://www.vlinderstichting.nl/vlinders/overzicht-vlinders/details-vlinder/groente-uil

648

649 Appendix 1. Illustration of how the modeling dataset was created, using the preceding years to
650 determine occupancy per area for each year separately. For example, to qualify an area as occupied
651 by a species in the year 2020, observations in at least two years of 2017, 2018 and 2019 were required.
652 Training set visits to occupied areas in 2020, 2021 and 2022 were used for training the model. Test set
653 visits to occupied areas in the year 2023 were used for testing the model, by calculating the AUC
654 metric.

| 2017 | 2018 | 2019 → | 2020 | | | |
| | 2018 | 2019 | 2020 → | 2021 | | |
| | | 2019 | 2020 | 2021 → | 2022 | |
| | | | 2020 | 2021 | 2022 → | 2023 |

655

| | preceding years to determine occupancy |
| | training set visits |
| | test set visits |

656

657

658

659

660 Appendix 2. Additional features that were created but not selected, listing the improvement of the
661 average AUC when only that feature was added to the final model.

| Feature | Average AUC improvement, when added to the final model that includes the 4 selected features |
|---|---|
| Total number of species, reported by the observer before 2023 | 0.005 |
| Average number of observations per species, reported by the observer before 2023 | 0.004 |
| Mean species rarity of observations, reported by the observer before 2023 | 0.001 |
| Mm precipitation on the day of the visit (transformed using the natural logarithm of mm+1) | 0.001 |

662

663

Appendix 3. Marginal dependence plots, showing the conditional probability of detection given occupancy per species for each visit, as a function of the additional features that were created but not selected, per species group. In grey, the individual curves of all species. In blue, the average curve over all species. In red, the four species that are covered in more detail throughout the results.