

1 **Solving sampling bias problems in presence-absence or presence-only**
2 **species data using zero-inflated models.**

3
4 **Running title: Sampling bias and zero-inflated models**

5
6 **Corresponding author: Victoria Nolan**

7 Life Sciences Building, University of Nottingham, Nottingham, NG7 2TQ, UK

8 victorianolan888@gmail.com, +447415316514

9 ORCID: 0000-0002-6069-963X

10 **Additional author: Tom Reader**

11 Life Sciences Building, University of Nottingham, Nottingham, NG7 2TQ, UK

12 ORCID: 0000-0001-7586-8814

13 **Additional author: Francis Gilbert**

14 Life Sciences Building, University of Nottingham, Nottingham, NG7 2TQ, UK

15 ORCID: 0000-0002-2727-4103

16
17
18
19
20
21
22
23
24
25

26 **Abstract**

27

28 *Aim*

29 Large databases of species records such as those generated through citizen science projects, archives or
30 museum collections, are being used with increasing frequency in species distribution modelling (SDM)
31 for conservation and land management. Despite the broad spatial and temporal coverage of the data, its
32 application is often limited by the issue of sampling bias and consequently, zero-inflation; there are
33 more zeroes (which are potentially ‘false absences’) in the data than expected. Here, we demonstrate
34 how pooling species presence data into a ‘pseudo-abundance’ count, can allow identification and
35 removal of sampling bias through the use of zero-inflated (ZI) models, and thus solves a common SDM
36 problem.

37

38 *Methods*

39 We present the results of a series of simulations based on hypothetical ecological scenarios of data
40 collection using random and non-random sampling strategies. Our simulations assume that the locations
41 of occurrence records are known at a high spatial resolution, but that the absence of occurrence records
42 may reflect under-sampling. To simulate pooling of presence-absence or presence-only data, we count
43 occurrence records at intermediate and coarse spatial resolutions, and use ZI models to predict the
44 counts (species abundance per grid cell) from environmental layers.

45

46 *Results*

47 Our results show that ZI models can successfully identify predictors of bias in species data, and produce
48 abundance prediction maps that are free from that bias. This phenomenon holds across multiple spatial
49 scales, thereby presenting an advantage over presence-only SDM methods such as binomial GLMs or
50 MaxEnt, where information about species density is lost, and model performance declines at coarser
51 scales.

52

53 *Main Conclusions*

54 Our results highlight the value of converting presence-absence or presence-only species data to ‘pseudo-
55 abundance’ and using ZI models to address the problem of sampling bias. This method has huge
56 potential for ecological researchers when using large species datasets for research and conservation.

57

58 **Key Words**

59 Conservation, Presence-absence, Presence-only, Sampling bias, Species abundance, Species
60 distribution model, Species occurrence, Zero-inflation

61

62 **Introduction**

63 Species distribution modelling (SDM) is widely used to address important ecological questions about
64 species distributions and the environment (Dormann et al., 2007; Phillips et al., 2009; Elith et al., 2011).
65 Species occurrence or abundance data from large, observational datasets such as citizen-science
66 projects, museum or herbarium collections and record lists are increasingly being used in SDM (Pearce
67 & Boyce, 2006; Schmeller et al., 2009; Tiago, Pereira & Capinha, 2017). The extensive spatial and
68 temporal coverage of the data, as well as the growing ease of online access, provides numerous benefits
69 over often costly and labour-intensive sampling methods employed in more focused scientific studies
70 of distribution (Dickinson, Zuckerglantz & Bonter, 2010; Dwyer, Carpenter, Bundhoo, Franklin &
71 Campbell, 2016; Gouraguine et al., 2019). Nevertheless, although some collections of species records
72 can be generated using hypothesis-led, systematic sampling protocols (Schmeller et al., 2009; Pocock
73 & Evans, 2014), much of these data comprise presence-only occurrence records, where there is often
74 little information about the source or survey effort accompanying the records (Boakes et al., 2010;
75 Rocchini et al., 2011). As a result sampling bias (also called sample selection or survey bias) is often
76 present: certain temporal periods, geographical areas or taxa are sampled more intensively or frequently
77 than others (Phillips et al., 2009; Dickinson et al., 2010; Bird et al., 2014).

78
79 Sampling bias in SDM can lead to over- or under-estimation of important species-environment
80 relationships (Syfert, Smith & Coomes, 2013), and predicted distribution maps may partly represent
81 survey effort rather than species niche requirements (Phillips et al., 2009; Mair & Ruete, 2016).
82 Proposed methods to correct for sampling bias generally rely on either spatial filtering of occurrence
83 records, or on the manipulation of background data ('pseudoabsences') (Phillips et al. 2009; Kramer-
84 Schadt et al. 2013; Fourcade et al. 2014, Boria et al. 2014). Both of these techniques have limitations:
85 the former results in a dataset of reduced sample size and statistical power (Wisz et al., 2008), whereas
86 the latter usually requires some prior knowledge of the source of the bias (Dudík, Schapire & Phillips,
87 2005; Phillips, 2008). A third option is the use of statistical models that can account for some of the
88 causes of sampling bias (Bird et al., 2014; Isaac, van Strien, August, de Zeeuw, & Roy, 2014), for
89 example Geographically Weighted Regression (GWR) (Brunsdon, Fotheringham & Charleton, 1998),
90 or Maximum Entropy (MaxEnt) with a bias layer, although again, most of these require prior knowledge
91 of the source of the bias.

92
93 One specific problem relating to sampling bias that is particularly noticeable in species abundance
94 databases is zero-inflation: the presence of more recorded zeroes or locations where data are absent than
95 expected under standard distributions (binomial, Poisson, negative binomial etc.) (Martin et al., 2005).
96 These excess zeros can arise from multiple processes. Some are considered to be 'true zeros', which
97 result from either ecological processes that render a site unsuitable for occupancy, or stochastic
98 processes, such as a sudden random extinction event in an otherwise suitable location (Cunningham &

99 Lindenmayer, 2005; Martin et al., 2005). In contrast, ‘false zeros’ are locations where a species occurs
100 but was not recorded because of errors or omissions in the sampling method (Dénes, Silveira &
101 Beissinger, 2015). These errors are either systematic and occur repeatedly throughout the survey
102 process (for example through a lack of detection or poor survey design), or are owing to sampling bias,
103 because some geographical areas have not been sampled at all (Bird et al., 2014).

104

105 Generalised Linear Models (GLMs) are a common method for analysing relationships between species
106 occurrences or abundance and environmental variables, but excess zeros are problematic for GLMs,
107 and if unaccounted for, can result in biased parameter estimates and poor predictive power (Lambert,
108 1992). As a possible solution to this problem, zero-inflated (ZI) models and their components
109 (extensions of GLMs) have been widely discussed in the literature (Lambert, 1992; Welsh,
110 Cunningham, Donnelly & Lindenmayer, 1996; Zuur, Ieno, Walker, Saveliev & Smith, 2009). ZI models
111 consist of two parts: a logistic component that models the probability of an observation being an excess
112 zero (hereafter called the “zero component”), and a “count component” that models a count (e.g. species
113 abundance) under an assumed distribution (Lambert, 1992). Both components of ZI models are capable
114 of producing zeros, and a key feature is the ability to include different predictor combinations in each
115 component. In other words, they can model the different sources of zeros independently (Wenger &
116 Freeman, 2008; Zuur et al., 2009).

117 ZI models, which require counts of occurrences (i.e. abundance), are rarely considered in SDM, because
118 most large datasets record species presences, not abundance. SDM methods that can use presence-only
119 data, such as MaxEnt, are therefore most commonly applied (Phillips & Dudík, 2008; Fitzpatrick,
120 Gotelli & Ellison, 2013; Fourcade et al., 2014). However, the ability of ZI models to separate the two
121 processes underlying the generation of zeroes in a species dataset could provide an alternative method
122 to model and account for sampling bias. ZI models can be used with any species database that records
123 abundance directly, or by aggregating presence-only or presence-absence data into counts of
124 occurrence. In this study, we therefore propose ZI models as a new, alternative method to address
125 problems of sampling bias in SDM. We present here the results of a series of simulations, based on
126 hypothetical ecological scenarios representing the large-scale collection of species occurrence data, that
127 aim to address three particular research questions.

128

129 Our first research question is to test our main theory of whether sampling bias (resulting in excess ‘false’
130 zeroes) can be modelled and accounted for using ZI models, in order to improve species distribution
131 predictions. ZI models have been used effectively to model true and false zeros in ecological count data,
132 such as when modelling the abundance of rare species (Welsh et al., 1996; Cunningham &
133 Lindenmayer, 2005; Martin et al., 2005). They are also particularly prevalent in the field of occupancy-
134 abundance modelling (Sileshi et al., 2009; Smith, Anderson & Millar, 2012), especially when there are

135 false zeros in the data owing to systematic sampling errors from imperfect detection (Wenger &
136 Freeman, 2008; Sólymos, Lele & Bayne, 2012; Williams, Yates, Stock, Barrett & Finn, 2016).
137 However, research into zero-inflation caused by spatio-temporal sampling bias in species occurrence
138 data is scarce. A few studies have used ZI models to identify and quantify sources of bias in species
139 data (Dwyer et al., 2016; Williams et al., 2016; Tiago, Ceia-Hasse et al., 2017), yet none have tested
140 the ability of the models to produce accurate predictions of species distributions from biased data. We
141 outline through our simulations how accurate distribution maps can be produced using ZI models in
142 this way, and we describe the required criteria during model fitting and prediction for this to occur. In
143 particular, our simulations also address our second research question: under what levels of zero-inflation
144 is our ZI model method most appropriate?

145 Our final research question considers the issue of scale, and the benefits of pooling fine-scale occurrence
146 data to model occurrence density across coarser spatial scales. Species presence is normally modelled
147 at the smallest spatial scale (grid cell size) possible, given the resolution of the records and
148 environmental layers used to build the model. Counting or aggregating presences across grid cells at a
149 larger spatial scale to generate “abundance” data intuitively seems to be a bad idea, because it throws
150 away information about the precise location of the records. However, this may be inevitable if predictor
151 layers have lower spatial resolution than occurrence location data, and we propose here that it may
152 actually present considerable advantages. Aggregated counts of occurrences are commonly not a direct
153 measure of true abundance (the total number of individuals of the target species), since each raw
154 occurrence often represents a locality which is home to several or many individuals. Regardless,
155 modelling ‘abundance’, and any zero-inflation therein, may give important clues to sources of bias in
156 the data which are not obvious in the raw occurrences, and the benefits of being able to identify and
157 eliminate bias could outweigh the costs of any loss of spatial resolution caused by aggregation.
158 Therefore, counting occurrence records at larger spatial scales in order to model “occurrence density”
159 may be a better alternative to traditional presence-only SDM methods. Indeed, abundance models have
160 been shown to perform better than presence-absence models fitted using the same data across multiple
161 spatial scales (Howard, Stephens, Pearce-Higgins, Gregory & Willis, 2014; Johnston et al., 2015).

162
163 Other methods do exist that propose aggregating occurrences into counts of ‘abundance’ that may also
164 provide advantages when using spatially biased species data, including Poisson point models (Renner
165 et al., 2015; Komori et al., 2020). These models can incorporate bias predictors when modelling
166 intensity rather than occurrences across the study area. Nevertheless, they still require a-priori
167 knowledge about potential bias predictors, whereas we show here that ZI models are able to provide an
168 indication of potential sources of sampling bias in the data when the exact sources are unknown.

169

170 We do not attempt to provide a detailed statistical summary of ZI models and theory (there is much
171 associated literature already available), but aim to draw attention to the main modelling methods and
172 usefulness of ZI models for ecological researchers and species distribution modellers dealing with large,
173 biased databases. We argue that ZI models can provide insight into, and correction methods for, the
174 bias in large species databases, and that they can be powerful and effective SDM tools.

175

176 **Materials and Methods**

177 Our general approach was to use ZI models to predict the observed number of species occurrences per
178 grid cell for a series of simulated species using predictors of either the biology of the species and/or
179 sampling bias in the data. We envisaged a large species for which it is theoretically possible to survey
180 all individuals in a landscape (e.g. trees, large animals). The true distribution of all individuals was
181 simulated for each species, and this distribution was then sampled incompletely, with or without spatial
182 sampling bias. Before sampling, the true abundance of the species could be calculated by summing
183 occurrences per grid square. But with incomplete sampling, the observed or “sampling abundance” per
184 grid cell is an underestimate. An alternative way to view our simulations, which is more realistic for
185 species which are small or hard to enumerate (e.g. smaller plants, most insects), is to consider each
186 occurrence in the raw data to represent a recorded encounter with the species at a local site which may
187 contain many individuals. In such cases, the models do not strictly predict abundance, but instead they
188 predict what we might call “occurrence density”.

189

190 As a result of the two-part nature of ZI models, two types of abundance predictions can be produced.
191 Assuming that all excess zeros arise from incomplete sampling, the first type of prediction is of true,
192 biological abundance (or occurrence density) across the study area, created only from the count
193 component of the model, which we call here the ‘count abundance prediction’. This is likely to be the
194 desired modelling outcome, especially for conservation and land-management planning. The second
195 type of prediction, which we here call the ‘sampling abundance prediction’, comes from the whole
196 model (combining both the count and zero component) and therefore represents the predicted
197 abundance (or occurrence density) that would be recorded if sampling were carried out in the same way
198 as when collecting the data that were used to fit the model. Bias in sampling will be reflected in this
199 second prediction. However, if some excess zeros arise also from biological zero-inflation, for example
200 if a species is clustered, the zero component will reflect some of the underlying biological processes as
201 well as the sampling bias. In this case, the count abundance prediction will only partially reflect the true
202 species abundance. The best type of prediction to use will therefore depend on the estimated strength
203 of biological zero-inflation versus the bias in the data.

204

205

206

207 *Simulation study area and predictor variables*

208 We simulated the occurrence of a hypothetical species in a study area that consisted of a 100 x 100 cell
209 grid at 1-km² resolution placed randomly within the boundary of England (Fig. 1a). The total area
210 covered by the grid is therefore 10,000 km² and there are 10,000 individual grid cells. Two predictor
211 variables were selected across this area. The first was a ‘biological predictor’ that we chose to be
212 ‘altitude’, which we used to define the relationship between the simulated species occurrences and
213 environment (Meynard et al., 2019). Real values for altitude (m) across the study area were obtained
214 from WorldClim DEM (accessed 10/05/18) at a 1-km² resolution and ranged from 0 to 284 m above
215 sea level (Fig. 1b). The choice of biological predictor for a simulation study of this sort is necessarily
216 somewhat arbitrary, but we chose altitude because it is both a plausible predictor of occurrence for a
217 range of organisms, and it is quite strongly spatially auto-correlated, an important possible source of
218 biological zero inflation in the abundance data formed when occurrences are counted across grid cells
219 at intermediate spatial scales. The actual biological mechanism underlying the relationship between
220 altitude and species occurrences is not important for this study, but altitude is a good proxy for a suite
221 of environmental variables such as temperature or precipitation commonly used in SDM which have
222 direct effects on species distributions.

223

224 Because altitude is spatially autocorrelated, and so is the sampling bias we wanted to investigate (see
225 below), there was a risk that biological and sampling bias predictors in our simulations could correlate:
226 depending on the positions of the simulated towns on our map, there could be a strong correlation
227 between real altitude and sampling effort. Thus, in order to allow us to investigate the impact of
228 sampling bias completely independently of the biological predictor, we also generated an alternative
229 ‘biological predictor’ with no autocorrelation: a spatially random control variable. This control variable
230 (henceforth labelled ‘altitude_randomised’) was created by randomising the real altitude values across
231 the study area at a 1-km² resolution (Fig. 1b), and hence removed any correlation between altitude and
232 distance from town.

233

234 The second predictor of observed species occurrence was a ‘bias predictor’ (‘distance from nearest
235 town’) which affected the virtual sampling of the simulated species. We assumed that the greater the
236 distance from a town, the lower the feasibility and likelihood of sampling occurring, as has previously
237 been seen in ecological studies (Reddy & Dávalos, 2003; Parnell et al., 2003; Kadmon, Farber & Danin,
238 2004). Unlike with altitude, we chose to simulate a hypothetical bias layer rather than use values based
239 on the locations of real towns, in order to ensure the lowest possible correlation between the two
240 predictors, although some correlation between them was likely because of spatial autocorrelation in
241 both. Within the study area, 10 points representing ‘town centres’ were randomly placed, and the
242 distance from the nearest town (m) was calculated for each grid cell, creating a continuous predictor
243 layer at 1-km² resolution across the study area. To reduce the influence of collinearity between

244 predictors, the process of generating the ‘town centres’ was repeated 10 times, creating 10 sets of
245 randomly placed ‘town centres’ (Fig. S1.1). As a result, mean Pearson’s correlation coefficients across
246 the 10 repetitions show weak correlations between the bias predictor ‘distance from nearest town’ and
247 the biological predictor ‘altitude’ ($r = -0.0499$, $SD \pm 0.228$), and even weaker correlations with the
248 biological predictor ‘altitude_randomised’ (-0.0044 , $SD \pm 0.012$).

249

250 To summarise, we had three variables in total across the simulation study area: two biological predictors
251 (‘altitude’ and ‘altitude_randomised’), and one bias predictor (‘distance from nearest town’). All
252 predictors were centred (the mean of each predictor was subtracted from each value of the predictor)
253 and scaled (the centred values were divided by the standard deviation of the predictor values) so that
254 the differences in units of the predictors was removed.

255

256 *Simulating the virtual species*

257 To obtain counts of ‘abundance’ to use in ZI models, we first simulated species occurrences across the
258 study area and then aggregated them into counts of “abundance” (alternatively interpreted as
259 occurrence density – see above). Because we assumed that the simulated distribution of occurrences
260 was the complete true distribution, all other locations are assumed to be ‘true absences’. Therefore,
261 when aggregating the raw occurrence points into ‘abundance’ counts, a value of 0 represented a true
262 absence and any value greater than 0 a true presence.

263

264 The recommended first step in a simulation study is to define the relationship between the environment
265 and occurrence points (Meynard et al., 2019). We modelled the distributions of three simulated species
266 each with 5,000 occurrence points (Fig. 1c). The occurrence points of the first species (‘random
267 species’) were simulated randomly across the study area, and show no preference for any environmental
268 condition. The second and third species were simulated based on the two biological predictors
269 (‘altitude’ and ‘altitude_randomised’) and were assumed to favour high altitudes; these species were
270 named ‘altitude species’ and ‘altitude_randomised species’ respectively. We chose these three scenarios
271 in order to create datasets in which different kinds of zero-inflation occur. For the random species, zero-
272 inflation can only occur as a result of sampling (where sites which are not sampled might be incorrectly
273 recorded as zeros), whilst for the altitude species and altitude_randomised species, zero-inflation can
274 result both from sampling and from the fact that grid cells are potentially not suitable for the species
275 because of environmental conditions.

276

277 We then simulated the effect of the relationship between our biological predictors and species
278 occurrences by creating layers of the probability of occurrence which varied according to altitude or
279 altitude_randomised (see Meynard and Kaplan 2013; Meynard et al., 2019). Initially we tried using a
280 linear relationship between the altitude predictor layers and probability of occurrence, but this

281 introduced relatively little zero-inflation in the data. For the purposes of investigating sampling bias
282 and zero-inflation we therefore chose to use a logarithmic relationship, whereby probability of
283 occurrence rapidly increases initially with small increases in altitude, but gradually tapers off at higher
284 altitudes. This heavily disfavours low altitude values, and the majority of these will be assigned low
285 probability values close to zero. Hence, biological aggregation of the occurrence points was effectively
286 increased, yielding greater zero-inflation. Each biological predictor was resampled to a 100 m x 100 m
287 resolution across the study area, and were then rescaled using the 'rescale by function' tool in ARCGIS
288 version 10.3.1 (ESRI, 2013), such that the new probability of occurrence layers (ranging between 0 and
289 1) were logarithmically related to the biological predictors.

290

291 Five thousand occurrence points were placed across the study area (using the ArcGIS tool: 'Create
292 Spatially Balanced Points') based on these altitude and altitude_randomised occurrence probability
293 layers. Due to computation limitations of the 'Create Spatially Balanced Points' tool, only one
294 occurrence point can be placed within a single raster cell. Therefore a resolution of 100 m x 100 m was
295 chosen for the probability layers so that up to 100 species occurrences could be placed in each 1-km²
296 grid cell. Although visually the altitude_randomised species appears to be randomly distributed across
297 the study area, it is actually the underlying altitude grid square values that are randomised: occurrences
298 of the altitude_randomised species still occur at higher densities in grid squares with higher altitude
299 values. As we used a logarithmic species response to the altitude_randomised layer, significant
300 (biological) zero-inflation still occurs in the raw data: occurrences are unlikely in low altitude grid cells,
301 generating lots of true zeros when occurrences were counted per grid cell (Tab. 1). Only the random
302 species distribution is completely random across the study area.

303

304 Finally, true (raw) species abundance (total number of occurrence points) was calculated for each 1-
305 km² grid cell. We felt the chosen grid scale was appropriate because, although the maximum abundance
306 per grid cell is strictly 100, no grid cells reached this value (the maximum was six occurrences per 1-
307 km grid cell), and we therefore assumed that it was unlikely that the shape of the distribution of
308 abundances would be significantly affected by the upper bound (i.e. unbounded distributions such as
309 Poisson or negative binomial were likely to be appropriate). In addition, using this grid scale sets up a
310 situation where location data are available at a higher resolution than the environmental predictors.
311 Hence, we are simulating a situation in which modellers must make a decision about how to aggregate
312 high resolution data across grid cells to create models which predict species distributions based on lower
313 resolution environmental predictors.

314

315 *Simulating the sampling strategies*

316 We considered two sampling strategies across the study area to represent alternative scenarios of
317 ecological data collection. The first is random sampling, where every 1-km grid cell has an equal chance

318 of being visited and sampled. If visited, we assume all species occurrences in the cell are recorded (i.e.
319 there is no detection error) and the result is the true (raw) abundance (count of all occurrences) for each
320 visited grid cell. The second sampling strategy is affected by spatial sampling bias and relates to the
321 ‘bias predictor’, where the probability of a grid cell being sampled decreases as distance from the nearest
322 ‘town centre’ increases. The grid cells selected for this strategy were chosen based on a probability
323 layer created using a logarithmic scaler of the ‘distance from nearest town’ predictor, again using the
324 ‘rescale by function’ ArcGIS tool. This time high probability values close to one were assigned to cells
325 with small numerical values i.e. cells closer to towns and more likely to be sampled, whereas low
326 probability values close to zero were assigned to cells with large ‘distance from nearest town’ values.
327 For each strategy, 2,000 grid cells (20% of the total) were sampled and species abundance was noted
328 for each one. All other (unsampled) squares were assigned an observed abundance of zero, creating a
329 zero-inflated dataset. All sources of zero-inflation in the simulated species abundance data before and
330 after sampling are shown in Tab. 1.

331

332 *Simulation 1: Investigating the accuracy of species distribution maps from ZI models*

333 To address our first question regarding the accuracy of ZI model predictions of abundance, we focused
334 initially on the performance of ZI Poisson models, and how this compared with equivalent conventional
335 Poisson GLMs. We include comparisons between a) ZI and GLM models, b) count and sampling
336 abundance predictions from ZI models, and c) alternative ZI models fitted using different combinations
337 of biological and bias predictors.

338

339 We chose to fit four GLMs and six ZI models for each of the three sets of species abundances per 1-
340 km² (random, altitude and altitude_randomised), all fitted with a Poisson distribution but with different
341 combinations of the biological or bias predictors (Tab. 2). These included combinations where different
342 predictors were tested in the count and zero components of the ZI models. Where the biological
343 predictor was included, models for the “altitude species” were fitted using altitude as a predictor, and
344 models for the altitude_randomised species were fitted using altitude_randomised. Model fitting was
345 repeated 10 times, each time using a different set of simulated ‘town centres’ (Fig. S1.1). Thus, there
346 are three species (random, altitude, altitude_randomised), two sampling strategies (random and biased)
347 and 10 repeats, resulting in 60 total simulation runs. All ZI and GLM models were fitted in R version
348 3.6.3 (R Core Team, 2019) using packages ‘stats’ (R Core Team, 2019) and ‘pscl’ (Zeileis et al., 2008).

349

350 Abundance predictions from each model were created using 10-fold cross-validation, where the data
351 were split into 10 subsets and each subset was used iteratively as the test data for which predictions
352 were created and the other nine subsets as training data. For the ZI models both count abundance and
353 sampling abundance predictions were evaluated. Model predictions were evaluated using a novel metric
354 based on the probability of obtaining the model predictions, that we named ‘deviation from the best

355 model' (D) (See Appendix S3 for more information). We used this metric, rather than conventional
356 measures of performance (e.g. root mean square) typically employed in presence-only or presence-
357 absence modelling, because it produces a measure of fit for count or abundance predictions which is
358 independent of the mean. D ranges from a minimum of one for a perfect model where model predictions
359 are equal to the true raw abundance data, and increases without limit as model predictive performance
360 decreases. Spearman's rank correlation coefficients (r_s) were also used to compare model abundance
361 predictions to the original model covariates.

362

363 To check that our results were not overly sensitive to the choice of predictor, simulations using average
364 temperature ($^{\circ}\text{C}$) (WorldClim, accessed 10/05/18) at a 1- km^2 resolution, as an alternative biological
365 predictor, were also carried following the same methodology (see Appendix S2): the results parallel
366 those of altitude, and so were omitted from the main results and discussion.

367

368 ***Simulation 2: Examining the impact of the extent of zero-inflation in the data***

369 To address our second question, about the effect of varying the extent of zero-inflation in the data (both
370 as a result of biological processes and sampling bias) on the effectiveness of the ZI models, we carried
371 out a second simulation. In our first simulation, we assumed 20% of grid cells were sampled, but in
372 Simulation 2 zero-inflation resulting from sampling bias was adjusted by varying the number of cells
373 sampled from the grid, ranging from 1000 (10%) to 10,000 (100%) at 10% increments. Therefore, the
374 highest level of zero-inflation occurred when 1000 cells were sampled, and thus 9000 cells were
375 assigned an abundance of zero simply because they were not sampled, and the lowest level of zero-
376 inflation occurred when 10,000 cells were sampled and none were assigned an abundance of zero for
377 this reason. At the same time zero-inflation resulting from biological processes was adjusted by adding
378 a threshold below which the altitude species can no longer survive, but keeping constant the number of
379 true occurrence points generated each time. With higher altitude thresholds, the species occurrences
380 were increasingly aggregated, and more cells were classified as true zeros. Altitude across the study
381 area ranged from 0 to 284 m, so we tested threshold values of 0 m, 50 m, 100 m, 125 m, 150 m, 175 m
382 and 200m (see Tab. S1.1 for number of cells above each threshold). Above these thresholds, species
383 occurrences were placed in a similar way based on weighted probability calculated from a logarithmic
384 scaler of the original altitude predictor as described previously. Both the random species and altitude
385 species were examined in scenarios with varying sample sizes, but obviously only the latter was tested
386 using the altitude threshold method.

387

388 Based on the results of Simulation 1, we selected three predictor combinations to fit the models and
389 create predictions. These included the GLM with both the bias and biological predictor (GLM4) and
390 two of the ZI models which differ only in the inclusion (ZI6) or exclusion (ZI2) of the bias predictor
391 from the zero component (Tab. 2). Although theoretically a ZI model that has only the biological

392 predictor in the count component, but both the biological and bias predictor in the zero component (as
393 with ZI3), would be the most obvious choice, in the real world the bias predictor may also have some
394 biological influence on the species distribution, and the researcher may not be sure whether it is a better
395 predictor of bias or biology. We therefore chose to use ZI6 rather than ZI3, to simulate better a real
396 world modelling scenario in which the causes of bias are unknown.

397

398 Model performance (D) was calculated for each simulation run with a particular combination of sample
399 size and altitude threshold. Finally, in order to evaluate the improvement in model performance created
400 by adding predictors of zero inflation, the difference in ' D ' was calculated between each model (GLM4
401 and ZI2, GLM4 and ZI6, and ZI2 and ZI6). This was repeated using both count abundance and sampling
402 abundance predictions for the ZI models. Again, model fitting was repeated 10 times each with two
403 sampling strategies (random and biased). Therefore, there were 200 simulation runs for the random
404 species (10 repeats, two sampling strategies and 10 levels of sampling zero-inflation), and 1,400
405 simulation runs for the altitude species (10 repeats, two sampling strategies, 10 levels of sampling zero-
406 inflation and seven altitude thresholds (levels of biological zero-inflation)).

407

408 ***Simulation 3: Comparing abundance versus presence-absence when aggregating spatial data***

409 Often when fitting distribution models the only data available are presence-only, and multiple species
410 occurrences within a grid cell are usually classified as a single presence. Often the predictors are only
411 available at a coarser spatial scale than the species occurrence data, forcing the modeller to aggregate
412 occurrences into coarser scale presence-only or presence-absence estimates. The coarser the resolution
413 at which the distribution is modelled, the more information is lost about both the precise location of
414 species occurrences, and species abundance (or occurrence density). However, if occurrences are
415 instead aggregated into count data, information about abundance or occurrence density is retained at all
416 scales, which may be more beneficial for conservation purposes. Therefore, even if only presence-only
417 data are available, ZI models fitted at a larger spatial scale using the summed counts of occurrence may
418 provide a better modelling method than traditional presence-only SDM that aggregate multiple
419 occurrences into presence-absence data. This effect is likely to be more pronounced when the species
420 data are biased, because ZI models attempt to model the excess zeroes from sampling bias, whereas
421 other methods, unless they explicitly incorporate bias correction, make no attempt to model or remove
422 the bias.

423

424 Our final simulation study addressed this question by comparing the performance of Poisson GLM and
425 ZI models predicting abundance of the altitude species (as was carried out in Simulation 1) with two
426 commonly used modelling methods that predict presence-absence: presence-absence binomial GLMs,
427 and presence-only MaxEnt models. This represents a scenario where the raw species occurrences

428 (simulated at a 100m resolution) are available at a greater resolution than the predictors (at a 1-km
429 resolution), so the modeller is required to make a decision on how to aggregate the data.

430

431 To fit the binomial GLM presence-absence models, the source data for which need to be in the form of
432 presence-absence rather than abundance, simulated 1-km cells that received an abundance count of zero
433 based on either the random or biased sampling strategy for the ZI models in Simulation 1 (i.e. 80% of
434 cells that were not considered to have been sampled) were classified automatically as an absence, and
435 any cell with species occurrences that was sampled was classified as a presence. All binomial GLMs
436 were fitted using the package ‘stats’ in R. As with Simulation 1, two GLMs were fitted, one with only
437 the biological predictor (‘Binomia-GLM1’ equivalent to GLM3) and one with the biological and bias
438 predictor (‘Binomial-GLM2’ equivalent to GLM4). Binomial occurrence predictions (i.e. predicted
439 probability of presence) were estimated across the study area from each model using 10-fold cross-
440 validation.

441

442 Two MaxEnt presence-only models were also fitted to the altitude species occurrence data, one with
443 altitude as the only predictor (‘Maxent1’), and one with both altitude and distance from nearest town as
444 predictors (‘Maxent2’). To produce presence-only data collected under a random or biased sampling
445 strategy, only occurrence points at a 100m resolution that fell within a 1-km cell that had been sampled
446 for the ZI models in Simulation 1 were retained; only these cells would be classified by MaxEnt as a
447 presence. Each model was fitted using the ‘dismo’ package (Hijmans, Phillips, Leathwick & Elith,
448 2017) in R, at a 1-km resolution with 10,000 randomly selected background ‘pseudo-absences’ and 10
449 repetitions across each set of town centres.

450 Comparing the performance of count/abundance models (Poisson GLM and ZI models) and
451 presence/presence-absence models (MaxEnt and binomial GLMs) required evaluation metrics which
452 could work with both types of model. As it is less feasible to convert presence-absence predictions to
453 abundance to use ‘D’, two other evaluation metrics were selected: Area Under the Curve (AUC) and
454 the Spearman’s Rank correlation coefficient (r_s) between the model predictors (‘altitude’ and/ or
455 ‘distance from town’) and each of the model predictions of count/abundance (GLM/ ZI) or habitat
456 suitability (MaxEnt/ binomial GLM). In order to calculate AUC for the ZI and GLM models, abundance
457 predictions were converted to binary presence-absence predictions, using an abundance threshold above
458 which the species was considered to be predicted to be present. Because some models produced
459 predicted abundances that all fell below one, the threshold for conversion was chosen to be the mean
460 abundance prediction across all grid cells for each individual model i.e. the threshold varied across each
461 GLM or ZI model. Mean AUC was calculated across the 10 repetitions for each model based on the
462 presence-absence predictions for all models compared to the true presence-absence based on all
463 occurrence locations across the study area. It should be noted that neither of these metrics offer a perfect

464 measure of model performance. AUC causes a loss of information from the Poisson GLMs and ZI
465 models, which are designed to predict abundance, while Spearman's rank retains more of the
466 information in the predictions of both types of model, but is necessarily relatively crude.

467

468 Finally, in order to assess the impact of the scale of data aggregation on the performance of abundance
469 and presence-absence models, additional models were fitted and compared across two other scales of
470 increasing coarseness: 2-km and 5-km. The larger the grid cell, the larger the mean count of occurrences
471 per cell, and hence the more data potentially lost by converting to presence-absence. ZI count abundance
472 predictions at a 2-km and 5-km scale were obtained following the methodology of Simulation 1 using
473 the ZI6 model structure and again converted to presence-absence predictions. MaxEnt and binomial
474 GLM presence-absence predictions at a 2-km and 5-km scale were obtained following the methodology
475 outlined previously in Simulation 3. Model predictors (altitude and distance from town) were converted
476 to coarser scales by calculating the mean values of each predictor at a 1-km resolution for each 2-km or
477 5-km cell. As before, all predictions were evaluated using AUC and Spearman's Rank correlation
478 coefficient (r_s).

479

480 **Results**

481 ***Simulation 1: Investigating the accuracy of species distribution maps from ZI models***

482 The results from Simulation 1 confirm that count abundance predictions from the ZI models provide
483 the most accurate estimates (according to the metric D) of true species abundance (Fig. 2 & S1.2).
484 Estimating true abundance based purely on the biology of the species rather than sampling processes is
485 usually the aim of ecological research, and these results suggest the count abundance predictions are
486 most likely able to fulfil these aims. In contrast, all GLMs are poor at predicting true abundance because
487 they do not separately model the excess (false) zeros generated by grid cells that have not been sampled.
488 The problem is exaggerated when sampling is not just incomplete, but is also biased; if the GLM
489 includes a predictor which is correlated with sampling effort (distance from nearest town), the model
490 performs even less well (compare pink and blue bars for GLM3 (without bias predictor) and GLM4
491 (with bias predictor) in Fig. 2) because it detects a spurious negative association between this predictor
492 and abundance (top panels, Fig. S1.3). Similarly, ZI sampling abundance predictions (predictions from
493 the whole model that potentially include the influence of sampling bias) perform poorly; rather than
494 estimating true abundance, reflecting the species niche, they predict abundance as it would appear to
495 observers employing each sampling strategy (Fig. 2 & S1.2). Again, these predictions are particularly
496 poor when sampling is biased (compare pink and blue bars for ZI2 and ZI6 in Fig. 2). These findings
497 hold true for all three species (altitude, altitude_randomised and random) (Fig. S1.2 & S1.3).

498

499 The ability to model excess zeros separately led to dramatically improved predictive power of true
500 abundance for all ZI models (see count abundance predictions in Fig. 2 and Fig. S1.2), although one

501 (ZI2) performed relatively less well than the others when sampling was biased (Fig. 2 & S1.2). In ZI2,
502 the bias predictor was included in the count component but not the zero component, meaning that like
503 the GLMs it detected a spurious negative association between abundance and distance from the nearest
504 town (middle panels, Fig. S1.3); if they included the bias predictor, the other ZI models (e.g. ZI3 or
505 ZI6) correctly detected that it was positively associated with the probability of an excess zero being
506 recorded (lower panels, Fig. S1.3).

507

508 Predicted distribution maps based on both the count abundance predictions and sampling abundance
509 predictions also support these findings (Fig. 3 & S1.4). Maps produced using ZI count abundance
510 predictions that account for bias where necessary (i.e. including predictors of bias in the zero component
511 when sampling is biased), correlate strongly with the biological predictor layer (altitude) ($r_s > 0.9$) and
512 show little influence of bias (distance from towns) (Fig. S1.5). When sampling is biased, both neglecting
513 to account for the bias in the zero component, or using the sampling abundance predictions, results in
514 low accuracy distribution maps that correlate more strongly with the bias predictor (r_s value between -
515 0.64 to -0.71) and less strongly with the biological predictor (r_s values between 0.60 to 0.74) (Fig. S1.5).
516 Distribution maps produced by the GLMs were also less accurate when sampling was biased and
517 predictors correlating with bias were included (Fig. 3 & S1.4). Maps from the GLMs which include the
518 bias predictor (GLM4) show a strong influence of sampling bias similar to that seen in the ZI sampling
519 abundance predictions. These maps show relatively weak correlations to the altitude predictor ($r_s = 0.60$)
520 compared to their counterpart GLMs that do not include the bias predictor (GLM3) ($r_s = 0.99$) (Fig.
521 S1.5). The prediction map from the GLM including both the biological and bias predictors (GLM4)
522 with biased sampling also shows a strong correlation to the bias predictor ($r_s = -0.72$).

523

524 Additional maps that depict the probability of each grid cell being an excess zero (i.e. predictions from
525 the zero component of a ZI model) further highlight the ability of ZI models to model separately the
526 biological and sampling processes, as well as providing insight into the nature of bias in the species
527 data (Fig. 3 & S1.4). This means that in real studies in which the sources of sampling bias are unknown,
528 inclusion of predictors that may correlate with sampling bias (e.g. distance to towns or roads,
529 accessibility, land-use etc.) in both the count and zero components of ZI models can help to both model
530 and identify likely causes of bias. This is a unique feature of the ZI models, and is something which the
531 GLMs are unable to reproduce; these models cannot provide insight into the bias or prediction maps
532 that eliminate sampling effects within the data.

533

534 ***Simulation 2: Examining the impact of the extent of zero-inflation in the data***

535 Real species occurrence or abundance data will suffer from variable levels of zero inflation resulting
536 from both biological and sampling processes. Therefore, the better performance of ZI models compared
537 with GLMs described in Simulation 1 may not occur in all circumstances, so exploring this issue was

538 our aim of Simulation 2. As anticipated, ZI count abundance predictions and GLM abundance
539 predictions have similar accuracy when the data are not zero-inflated: when the whole study area is
540 surveyed, all absences are ‘true absences’, the species is randomly distributed with no biological zero-
541 inflation, and the difference in performance is zero (Fig. 4, see random species (R) in left and middle
542 panels). When considering the random species only (i.e. with no biological zero-inflation), as less of
543 the study area is surveyed, zero-inflation as a result of sampling increases, and therefore the
544 effectiveness of ZI model count abundance predictions improves in comparison to GLMs. Although
545 this phenomenon occurs under both sampling strategies, it is most noticeable when both sampling is
546 biased and that bias is accounted for in the model (by including the bias predictors in the ZI zero
547 component as in ZI6 for example).

548

549 As with the random species, when there are high levels of incomplete sampling for the altitude species
550 (e.g. ~20% or fewer cells are sampled), ZI model count abundance predictions are consistently better
551 than GLM predictions, regardless of biological zero-inflation (Fig. 4, left and middle panels). However
552 as more of the area is surveyed (> 20%), the difference in performance decreases. At low levels of
553 biological zero-inflation, this difference tends towards zero. However, at higher levels of biological
554 zero-inflation, GLM predictions are actually more accurate than the ZI model count abundance
555 predictions under both random and biased sampling scenarios. This can best be understood by looking
556 at Fig. S1.6 showing the results based on sampling abundance predictions from the ZI model, rather
557 than count abundance predictions: in contrast to the count abundance predictions, as biological zero-
558 inflation increases, ZI sampling abundance predictions increasingly outperform those of the GLM. This
559 is because the zero component, which is combined with the count component to create the sampling
560 abundance prediction, is able to predict the excess zeroes caused by the biological driver, while the
561 GLM cannot. Therefore, if high levels of biological zero-inflation are suspected in the data, both the
562 count and sampling abundance predictions should be considered and evaluated before choosing the best
563 predictions of species abundance.

564

565 Reiterating our results from Simulation 1, when sampling is random there is no benefit of including the
566 bias predictor in the zero component under any levels of sampling or biological zero-inflation (Fig. 4
567 & S1.6, top right panels). Under biased sampling scenarios, models accounting for bias (by including
568 the bias predictor in the zero component as in ZI6 for example) are most effective when there are high
569 levels of sampling-related zero-inflation and low levels of biological zero-inflation. As either the area
570 surveyed or biological zero-inflation increases, the effectiveness of these models reduces compared to
571 models that fail to account for bias (Fig. 4, bottom right panel). Nevertheless, the majority of differences
572 seen between ZI models are relatively small compared to those between the ZI models and GLMs.

573

574

575 ***Simulation 3: Comparing abundance versus presence-absence data across multiple spatial scales***

576 The results from Simulation 3 support our hypothesis that, when dealing with biased species data,
577 modelling aggregated count data using ZI models is a better choice than modelling aggregated presence-
578 absence or presence-only data, as is commonly done in traditional SDM studies, using approaches such
579 as binomial GLMs or MaxEnt (Fig. 5). The only model to perform consistently well across all spatial
580 scales when dealing with the biased species data was the ZI model, which maintained strong correlations
581 to the biological predictor ($r_s > 0.9$) and low correlations to the bias predictor ($-0.12 < r_s < 0.07$) across
582 all scales (Fig. 5). Predicted maps of the altitude species distribution also show that the ZI model count
583 abundance predictions provide the most accurate reflection of the true species distribution as the scale
584 of data aggregation increases (Fig. S1.7). Binomial-GLM2 and MaxEnt2 models, which incorporate
585 the bias predictor, produced predictions that are heavily influenced by sampling bias at a 1-km scale,
586 with strong correlations to the bias predictor ($r_s < -0.75$) (Fig. 5 & S1.7). These increase in strength as
587 scale increases to 2-km and 5-km, so that both model predictions produce correlations to the bias
588 predictor close to one ($r_s < -0.92$). Both MaxEnt1 and binomial-GLM1 (which do not include the bias
589 predictor) were able to produce accurate predictions with the biased data at a 1-km resolution, although
590 performance declined as the scale became coarser. Even when the species data was collected using a
591 random sampling strategy, the performance of the presence-absence models declined as the scale
592 became coarser and more information was lost with data aggregation (Fig. 5); this phenomenon was not
593 seen in the ZI models and performance remained high as scale increased.

594

595 Model evaluation using mean AUC based on the presence-absence predictions also supports these
596 findings (Fig. 6 & S1.8). Across all three scales, the ZI model was best suited to model the biased
597 species data compared to the MaxEnt and binomial GLM models that were fitted using the bias predictor
598 (Fig. 6). The presence-absence models have a much larger variance in performance than the ZI
599 abundance models, especially at coarser scales, with some repetitions producing AUC values below 0.5
600 and above 0.9 (Fig. 6 & S1.8). The ZI model also outperformed several of the MaxEnt and binomial
601 GLMs fitted without the bias predictor, including the MaxEnt1 model at a 2-km scale and the binomial-
602 GLM1 at a 5-km scale (Fig. S1.8), although it produced slightly lower mean AUC values than some of
603 the presence-absence models when the bias predictor was excluded. Nevertheless, if the sampling bias
604 source is unknown, it might be difficult to exclude completely predictors correlating with the bias, so
605 choosing a ZI model is still likely to be the safest option to produce the best, most robust predictions
606 least affected by sampling bias.

607

608

609

610

611

612 **Discussion**

613 Sampling bias in species data is problematic for SDM, and many researchers call for greater awareness
614 and development of correction methods to deal with this issue (Araújo & Guisan, 2006; Bystriakova,
615 Peregrym, Erkens, Bezsmertna, & Schneider, 2012; Kramer-Schadt et al., 2013). Our simulations using
616 ZI models highlight a novel approach for dealing with sampling bias and zero-inflation in SDM, which
617 we believe can be applied to a wide variety of ecological and conservation research questions that use
618 large databases of species records. Our results reveal that ZI models have the potential both to reduce
619 the impact of bias on predictions which are used for biological inference, and to provide insights into
620 previously unknown causes or correlates of sampling bias. This method can be used with both raw
621 abundance data, and with abundance data created by summing occurrences from presence-only data
622 across a larger spatial scale, and therefore offers an alternative to traditional presence-only SDM
623 methods. As spatial occurrence data is often present at a finer scale than the environmental predictors,
624 decisions about data aggregation have to be made when fitting distribution models. We found that even
625 though information about the precise location of species occurrences is sacrificed, aggregating species
626 occurrences into counts of abundance and fitting ZI models produces better estimates of a species
627 distribution, especially when the species data is biased by sampling methods, than aggregating
628 occurrences into presence-absence form at a coarser spatial scale, as is common with traditional SDM
629 methods such as binomial GLMs or MaxEnt.

630

631 Species distribution maps are an important resource for conservation planners (Rodríguez, Brotons,
632 Bustamante & Seoane, 2007), yet there is often little consideration of inaccuracies or uncertainty in
633 these maps or associated models (Elith, Burgman & Regan, 2002; Zuquim et al., 2014). Our results
634 show how the biological information value of maps based on GLM, MaxEnt and ZI sampling abundance
635 predictions can be reduced by sampling bias. In contrast, the distribution maps produced from the
636 predictions from the count component of ZI models are accurate reflections of the species niche and
637 true abundance, even when species data are spatially biased, providing that the bias influence is
638 accounted for in the model by included all predictors suspected of capturing or correlating with the bias
639 in both ZI count and zero components. If in doubt about whether a predictor is likely to be a source of
640 bias, inclusion in both parts will not only alleviate the problem of bias, but will also provide insight into
641 whether it actually is a introducing a large number of excess ('false') zeros. Additionally, ZI model
642 coefficients allowed examination of potential causes of bias; in ZI6 (the model including both the bias
643 and biological predictor in the zero component) from Simulation 1, 'distance from nearest town' was
644 influential only in the zero component, and was not spuriously identified as influencing true abundance.
645 Currently, there are few statistical models that allow post-modelling identification of bias sources.
646 Many SDM techniques rely on prior understanding and some form of quantification of the bias in order

647 to remove it (Phillips, 2008), so ZI models provide an advantage over these traditional bias correction
648 methods in their ability to shed light on potential causes of bias.

649

650 If all excess zeros are false zeros, count abundance predictions from ZI models should always reflect
651 the true species niche, and the zero component will be modelling only excess zeros from non-biological,
652 sampling processes. However, this scenario is unlikely in ecological systems. In reality, as in our
653 simulations with the altitude and altitude_randomised species, the excess zeros will result from a
654 combination of biological zero-inflation and sampling zero-inflation. Therefore, the count abundance
655 prediction may not always be predicting true abundance, and the zero component may actually be
656 dominated by biological processes, as we suggest is the case for the results from Simulation 2. In this
657 case, the sampling abundance prediction will actually be a more accurate reflection of true species
658 abundance. Nevertheless, by examining the significance and influence of predictors in both
659 components, their plausibility as causes of bias can be inspected: biological predictors of abundance are
660 likely to be significant in both parts of the ZI model, whereas sampling predictors are unlikely to appear
661 influential in the count component.

662

663 After identifying potential bias predictors, modellers can make more informed choices about whether
664 to eliminate these predictors from either ZI component, whether the zero component is more heavily
665 dominated by biological or sampling processes, and if the count abundance or sampling abundance is
666 more likely to reflect true species abundance. A good understanding of the biology of the species being
667 modelled is therefore key. Additionally, despite the post-model-fitting ability of ZI models to
668 distinguish bias, beginning any analysis of a zero-inflated dataset, it is important also to try and identify
669 the source of excess zeros as either from biology or sampling processes (Martin et al., 2005).
670 Consequently, although one benefit of ZI models is the ability to use different sets of covariates in the
671 count and zero components (Lambert, 1992; Zuur et al., 2009), it is important only to include
672 appropriate, relevant predictors in each part where possible.

673

674 The collection of species data varies widely in its scale and standardisation, from single museum
675 specimens collected by natural history experts, to more local, standardised recording schemes (Pocock
676 & Evans, 2014) and to international, opportunistic recording schemes such as eBird (Sullivan et al.,
677 2009). The more standardised and directed the protocols, the lower the likelihood of sampling bias and
678 'false zeros' in the data. In these cases, a simple Poisson or negative binomial GLM may suffice rather
679 than a ZI model; at very low levels of zero-inflation the performance of the GLMs was shown to be
680 equal to that of the ZI models in Simulation 2. Nevertheless, our findings from Simulation 2 suggest
681 that, regardless of biological zero-inflation, when sampling is suspected to be very incomplete
682 (estimated coverage of total study area < ~20%), ZI models will always be the optimum choice. At low
683 levels of biological zero-inflation, we found ZI models to be more effective than GLMs even when

684 sampling coverage approached levels as high as 90%, as might be the case for species with broad ranges
685 that have been extensively documented, such as important or conspicuous species in countries with long
686 histories of species record keeping.

687

688 As well as the Poisson distribution, the negative binomial distribution is also often used for count data,
689 which can also be applied within a zero-inflated modelling framework (Ridout, Hinde & Demétrio,
690 2001; Minami, Lennert-Cody, Gao, & Román-Verdesoto, 2007; Zuur et al., 2009). The negative
691 binomial distribution is able to model an extra proportion of the excess zeros compared to the Poisson
692 distribution through the use of an extra model parameter (Θ) (Fisher, 1941) and can therefore account
693 for biological aggregation and overdispersion in ecological data (Lindén & Mäntyniemi, 2011). We
694 chose not to investigate a ZI negative binomial model in these simulations to remove confusion when
695 communicating our main message, although we acknowledge that under high levels of biological zero-
696 inflation (as in Simulation 2), such models may well be more effective than the ZI Poisson models.
697 Therefore, when analysing presence-only species data suffering from high levels of sampling bias, a ZI
698 Poisson model will usually be effective, but it is valuable to know that there are different ZI model
699 types that can be used to address ecological or statistical issues that may arise in species data.

700

701 The majority of SDM research to date has focused on producing presence-absence or presence-only
702 distribution maps of species or communities (Brotons, Thuiller, Araújo, & Hirzel, 2004; Phillips,
703 Anderson & Schapire, 2006; Lyashavska, Brus & van der Meer, 2016). Species abundance maps are
704 produced more infrequently, often due to the practical difficulty of measuring absolute abundance
705 (Lyashavska et al., 2016). However, their ability to display extra information about density means they
706 are often more informative and preferred (Pearce & Ferrier, 2001; Barry & Welsh, 2002; Johnston et
707 al., 2015).

708

709 Although count data are known commonly to suffer from zero-inflation, ZI models have been used to
710 produce accurate species abundance maps from systematically collected species data in very few studies
711 (Bouyer et al., 2015; Lyashavska et al., 2016), and none have acknowledged or explored bias in their
712 data. It is also not recommended to use SDM to predict species abundance from presence-only or
713 presence-absence data (Jiménez-Valverde, Aragón & Loboet, 2020), so ZI models that fit abundance
714 by default should be able to cover this methodology gap in the field of SDM. Additionally, scale is
715 hugely important in SDM. Species distributions are often modelled at coarse resolutions across national
716 or international scales due to the availability of predictors, even though occurrences relate more to
717 localised environmental factors (Guisan et al., 2007; Kuemmerlen et al., 2014). The coarser the grain
718 size used in presence-absence or presence-only SDM, the more the raw occurrences are aggregated into
719 a binary variable and density information is lost. Therefore, it is likely that at coarse resolutions, using

720 abundance rather than occurrence data preserves more information and will produce more accurate
721 maps of habitat suitability.

722

723 Our findings from Simulation 3 suggest that when having to decide how to aggregate data to match the
724 coarser resolution of the environmental predictors, the best method is to aggregate species occurrences
725 into counts of abundance and fit using a ZI model, rather than aggregate into presence-absence data and
726 fit using a traditional SDM method such as MaxEnt. This provides two main benefits over presence-
727 absence methods in that a) ZI models are able to identify and account for bias without prior knowledge
728 of the bias sources and b) extra information about species abundance is retained and modelled. We
729 found that as scale became increasingly coarser, only the ZI models retained a high level of predictive
730 power and were an accurate reflection of species niche compared to MaxEnt or binomial GLMs,
731 especially when the data suffered from sampling bias. We believe that ZI models have an advantage
732 over other statistical methods in that they can be used with either presence-absence data or abundance
733 data collected from citizen science projects: presence-absence data can just be aggregated into a count
734 at a particular resolution. Furthermore, scale was shown to have little influence on the predictive power
735 of ZI models providing bias was accounted for. Nevertheless, this was only simulated across relatively
736 small resolutions (up to 5-km) due to the limitations of the study area and requirement for zero-inflated
737 data, whereas many studies map distributions at larger scales (> 10-km) (Thuiller, Lavorel, Sykes, &
738 Araújo, 2006; Luoto, Virkkala, & Heikkinen, 2007). It is uncertain therefore whether this pattern holds
739 true across more coarse scales of analysis.

740

741 In this paper, we have investigated the performance of ZI models under a relatively restricted set of
742 scenarios. We acknowledge that our findings may therefore be case-specific and we are addressing this
743 with on-going research (Nolan et al., unpubl.). For example, we chose to use a simple scenario in which
744 only two predictors, a biological predictor and a bias predictor, generate patterns in the species
745 distribution. The altitude species was assigned a simple preference for high altitudes, when in fact, there
746 are likely several different environmental influences on the species niche. Furthermore, some of these
747 biological predictors of species presence will also predict sampling bias. Therefore, it is important that
748 prior consideration is given to the possible influences of any predictor included in the model on both
749 ecological processes and sampling behaviour before it is decided whether to include it in either part of
750 the ZI model.

751

752 GLMs, and by extension ZI models, have been criticised for their inability to capture the complex, non-
753 linear relationships which may often characterise species responses to the environment, in contrast with
754 more modern methods such as MaxEnt or other machine learning techniques which are more flexible
755 (Austin, 2002). Nevertheless, GLMs and ZIs also have some clear benefits, such as the ease with which
756 they can be applied, and the transparency of their design. Here, we have shown an additional benefit of

757 ZI models not yet available with any other modelling approach: the ability to simultaneously account
758 for bias and to make inferences about it, when predicting distributions from incomplete sampling. We
759 believe that our approach using ZI models has broad applicability to a variety of scenarios when bias is
760 present, and there are suspected predictors of bias available. ZI models should be especially valuable
761 when species abundance is of interest to the modeller, such as when modelling distributions of
762 individual large animals or trees. Although we acknowledge that GLMs and ZI models have limitations,
763 there is a range of options for more complex versions of these models, such as those incorporating
764 polynomial terms, interactions and LASSO variable selection (Hastie et al., 2009; Vollerling et al.,
765 2019), which might allow such models to capture non-linear/complex responses to the environment at
766 the same time as modelling the causes of excess zeroes.

767

768 In our simulations, we assume that all ‘false absences’ are due to sampling bias, but it is likely that in
769 many cases, particularly for rare or cryptic species, they are also generated by detection errors
770 (Fitzpatrick, Preisser, Ellison, & Elkinton, 2009; Dickinson et al., 2010; Kosmala, Wiggins, Swanson,
771 & Simmons, 2016). The species range size and the scale of detectability of the individuals is likely to
772 influence the interpretation of the model “abundance” predictions. For example under-estimation of
773 true abundance could occur when modelling small organisms which appear frequently during the
774 survey, and will be more representative of the likelihood of successfully sampling the species. On the
775 other hand, over-estimation could occur when modelling large, mobile organisms that cover multiple
776 sampling locations, so prediction abundance might be a proxy of the probability of encountering one of
777 a small number of individuals. Hence, there may be three sources of excess zeros: true zeros from
778 unsuitable habitat, false zeros from lack of sampling and false zeros from detection error. When
779 detection errors are significant, ZI models will not be able to distinguish between the different types of
780 false zeros; but by including predictors in both the count and zero components of the model that capture
781 the processes generating all types of zeros, we believe that ZI models will still be able (mostly) to
782 account for these excess ‘false’ zeros, and combined with expert knowledge can provide some
783 information about their sources.

784

785 **Conclusion**

786 Large collections of species data are extremely useful for SDM and conservation, and yet are limited
787 by issues associated with the recording processes, including sampling bias and zero-inflation. Our
788 simulations show that ZI models can fit biased data and identify sources of bias. Most importantly for
789 conservation, by using only predictions from the count component of the ZI model (i.e. the count
790 abundance predictions), biased species data can be used to produce distribution maps comparable to
791 those using unbiased data. We also highlight the importance of considering the use of abundance data
792 in SDM, especially at large spatial scales, when valuable ecological information about density is lost if
793 data in each cell are converted to presence-absence or presence-only. ZI models are advantageous

794 compared to other commonly used SDM techniques such as MaxEnt owing to their ability to retain
795 information about abundance and also to identify and remove bias without prior knowledge of the bias
796 sources. We believe ZI models have been largely overlooked in ecological research, even though they
797 have a huge potential to be useful in SDM, and could have great benefits for conservation and our
798 environment.

799

800 **Acknowledgements**

801 This research was made possible by funding support from the University of Nottingham, UK and the
802 Woodland Trust, Grantham, UK. We gratefully acknowledge Tim Newbold and Richard Field for their
803 helpful comments on a draft version of this manuscript.

804

805

806 **Data availability**

807 R code and data used to produce the analyses and results reported within this manuscript can be found
808 freely available at the following online location: <https://doi.org/10.6084/m9.figshare.13118417.v2>.
809 Elevation and temperature data used in the study are publicly available for download from WorldClim
810 DEM (<https://www.worldclim.org/>).

811

812 **Biosketch**

813 **Victoria Nolan** is an ecologist with interests in distribution modelling and conservation at a
814 macroecological scale, alongside an interest in large data and novel methods of statistical analysis.

815 **Tom Reader** is an Associate Professor with interests in ecology, behaviour and statistics, with
816 particular focus on the evolution of animal signals such as Batesian mimicry and aposematism.

817 **Francis Gilbert** is a Professor of Ecology with interests in distribution modelling and conservation,
818 together with the evolution of ecological and behavioural attributes of organisms, with a specific focus
819 on hoverflies.

820

821

822

823

824

825

826 **References**

- 827 Araújo, M. B., & Guisan, A. (2006). Five (or so) challenges for species distribution modelling. *Journal of*
828 *Biogeography*, 33(10), 1677–1688.
- 829 Austin, M. P. (2002). Spatial prediction of species distribution: an interface between ecological theory and
830 statistical modelling. *Ecological Modelling*, 157, 101–118.
- 831 Barry, S. C., & Welsh, A. H. (2002). Generalized additive modelling and zero inflated count data. *Ecological*
832 *Modelling*, 157(2), 179–188.
- 833 Bird, T. J., Bates, A. E., Lefcheck, J. S., Hill, N. A., Thomson, R. J., Edgar, G. J., ... Frusher, S. (2014).
834 Statistical solutions for error and bias in global citizen science datasets. *Biological Conservation*, 173,
835 144–154.
- 836 Boakes, E. H., McGowan, P. J. K., Fuller, R. A., Chang-qing, D., Clark, N. E., O'Connor, K., & Mace, G. M.
837 (2010). Distorted views of biodiversity: spatial and temporal bias in species occurrence data. *PLoS*
838 *Biology*, 8(6).
- 839 Boria, R. A., Olson, L. E., Goodman, S. M., & Anderson, R. P. (2014). Spatial filtering to reduce sampling bias
840 can improve the performance of ecological niche models. *Ecological Modelling*, 275(Supplement C),
841 73–77.
- 842 Bouyer, Y., Rigot, T., Panzacchi, M., Moorter, B. V., Poncin, P., Beudels-Jamar, R., ... Linnell, J. D. C. (2015).
843 Using Zero-Inflated models to predict the relative distribution and abundance of roe deer over very
844 large spatial scales. *Annales Zoologici Fennici*, 52(1–2), 66–76.
- 845 Brotons, L., Thuiller, W., Araújo, M. B., & Hirzel, A. H. (2004). Presence-absence versus presence-only
846 modelling methods for predicting bird habitat suitability. *Ecography*, 27(4), 437–448.
- 847 Brunson, C., Fotheringham, S., & Charlton, M. (1998). Geographically Weighted Regression. *Journal of the*
848 *Royal Statistical Society: Series D (The Statistician)*, 47(3), 431–443.
- 849 Bystriakova, N., Peregrym, M., Erkens, R., Bezsmertna, O., & Schneider, H. (2012). Sampling bias in
850 geographic and environmental space and its effect on the predictive power of species distribution
851 models. *Systematics and Biodiversity*, 10.
- 852 Cunningham, R. B., & Lindenmayer, D. B. (2005). Modeling count data of rare species: some statistical issues.
853 *Ecology*, 86(5), 1135–1142.
- 854 Dénes, F. V., Silveira, L. F., & Beissinger, S. R. (2015). Estimating abundance of unmarked animal populations:
855 Accounting for imperfect detection and other sources of zero inflation. *Methods in Ecology and*
856 *Evolution*, 6(5), 543–556.
- 857 Dickinson, J. L., Zuckerberg, B., & Bonter, D. N. (2010). Citizen Science as an ecological research tool:
858 challenges and benefits. *Annual Review of Ecology, Evolution, and Systematics*, 41(1), 149–172.
- 859 Dormann, C., M. McPherson, J., B. Araújo, M., Bivand, R., Bolliger, J., Carl, G., ... Wilson, R. (2007).
860 Methods to account for spatial autocorrelation in the analysis of species distributional data: A review.
861 *Ecography*, 30(5), 609–628.
- 862 Dudík, M., E. Schapire, R., & J. Phillips, S. (2005). Correcting sample selection bias in maximum entropy
863 density estimation. *Advances in Neural Information Processing Systems*, 17.

864 Dwyer, R. G., Carpenter-Bundhoo, L., Franklin, C. E., & Campbell, H. A. (2016). Using citizen-collected
865 wildlife sightings to predict traffic strike hot spots for threatened species: A case study on the southern
866 cassowary. *Journal of Applied Ecology*, 53(4), 973–982.

867 Elith, J., Burgman, M. A., & Regan, H. M. (2002). Mapping epistemic uncertainties and vague concepts in
868 predictions of species distribution. *Ecological Modelling*, 157(2-3), 313–329.

869 Elith, J., & Leathwick, J. R. (2009). Species distribution models: Ecological explanation and prediction across
870 space and time. *Annual Review of Ecology, Evolution, and Systematics*, 40(1), 677–697.

871 Elith, J., Phillips, S.J., Hastie, T., Dudík, M., Chee, Y.E., & Yates, C. J. (2011). A statistical explanation of
872 MaxEnt for ecologists. *Diversity and Distributions*, 17, 43–57.

873 Fisher, R. A. (1941). The Negative Binomial Distribution. *Annals of Eugenics*, 11(1), 182–187.

874 Fitzpatrick, M. C., Preisser, E. L., Ellison, A. M., & Elkinton, J. S. (2009). Observer bias and the detection of
875 low-density populations. *Ecological Applications*, 19(7), 1673–1679.

876 Fitzpatrick, M. C., Gotelli, N. J., & Ellison, A. M. (2013). MaxEnt versus MaxLike: Empirical comparisons
877 with ant species distributions. *Ecosphere*, 4(5), art55.

878 Fourcade, Y., Engler, J. O., Rödder, D., & Secondi, J. (2014). Mapping species distributions with MAXENT
879 using a geographically biased sample of presence data: A performance assessment of methods for
880 correcting sampling bias. *PLoS ONE*, 9(5).

881 Gouraguine, A., Moranta, J., Ruiz-Frau, A., Hinz, H., Reñones, O., Ferse, S. C. A., ... Smith, D. J. (2019).
882 Citizen science in data and resource-limited areas: A tool to detect long-term ecosystem changes. *PLOS*
883 *ONE*, 14(1), e0210007.

884 Guisan, A., Zimmermann, N. E., Elith, J., Graham, C. H., Phillips, S., & Peterson, A. T. (2007). What matters
885 for predicting the occurrences of trees: Techniques, data, or species' characteristics? *Ecological*
886 *Monographs*, 77(4), 615–630.

887 Hastie, T., Tibshirani, R., & Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference,*
888 *and Prediction, Second Edition*. Springer Science & Business Media.

889 Hijmans, R., Phillips, S., Leathwick, J., & Elith, J. (2017). dismo: Species Distribution Modeling. R package
890 version 1.1-4. <https://CRAN.R-project.org/package=dismo>.

891 Howard, C., Stephens, P. A., Pearce-Higgins, J. W., Gregory, R. D., & Willis, S. G. (2014). Improving species
892 distribution models: The value of data on abundance. *Methods in Ecology and Evolution*, 5(6), 506–
893 513.

894 Isaac, N. J. B., van Strien, A. J., August, T. A., de Zeeuw, M. P., & Roy, D. B. (2014). Statistics for citizen
895 science: Extracting signals of change from noisy ecological data. *Methods in Ecology and Evolution*,
896 5(10), 1052–1060.

897 Jiménez-Valverde, A., Aragón, P., Lobo, J.M. (2021). Deconstructing the abundance–suitability relationship in
898 species distribution modelling. *Glob. Ecol. Biogeogr.* 30, 327–338.

899 Johnston, A., Fink, D., Reynolds, M. D., Hochachka, W. M., Sullivan, B. L., Bruns, N. E., ... Kelling, S. (2015).
900 Abundance models improve spatial and temporal prioritization of conservation resources. *Ecological*
901 *Applications*, 25(7), 1749–1756.

902 Kadmon, R., Farber, O., & Danin, A. (2004). Effect of roadside bias on the accuracy of predictive maps
903 produced by bioclimatic models. *Ecological Applications*, 14(2), 401–413.

904 Komori, O., Eguchi, S., Saigusa, Y., Kusumoto, B., & Kubota, Y. (2020). Sampling bias correction in species
905 distribution models by quasi-linear Poisson point process. *Ecological Informatics*, 55, 101015.

906 Kosmala, M., Wiggins, A., Swanson, A., & Simmons, B. (2016). Assessing data quality in citizen science.
907 *Frontiers in Ecology and the Environment*, 14(10), 551–560.

908 Kramer-Schadt, S., Niedballa, J., Pilgrim, J. D., Schröder, B., Lindenborn, J., Reinfelder, V., ... Wilting, A.
909 (2013). The importance of correcting for sampling bias in MaxEnt species distribution models.
910 *Diversity and Distributions*, 19(11), 1366–1379.

911 Kuemmerlen, M., Schmalz, B., Guse, B., Cai, Q., Fohrer, N., & Jähnig, S. (2014). Integrating catchment
912 properties in small scale species distribution models of stream macroinvertebrates. *Ecological*
913 *Modelling*, 277, 77–86.

914 Lambert, D. (1992). Zero-Inflated Poisson regression, with an application to defects in manufacturing.
915 *Technometrics*, 34(1), 1–14.

916 Lindén, A., & Mäntyniemi, S. (2011). Using the negative binomial distribution to model overdispersion in
917 ecological count data. *Ecology*, 92(7), 1414–1421.

918 Luoto, M., Virkkala, R., & Heikkinen, R. K. (2007). The role of land cover in bioclimatic models depends on
919 spatial resolution. *Global Ecology and Biogeography*, 16(1), 34–42.

920 Lyashevskaya, O., Brus, D. J., & van der Meer, J. (2016). Mapping species abundance by a spatial zero-inflated
921 Poisson model: A case study in the Wadden Sea, the Netherlands. *Ecology and Evolution*, 6(2), 532–
922 543.

923 Mair, L., & Ruete, A. (2016). Explaining spatial variation in the recording effort of citizen science data across
924 multiple taxa. *PLOS ONE*, 11(1), e0147796.

925 Martin, T. G., Wintle, B. A., Rhodes, J. R., Kuhnert, P. M., Field, S. A., Low-Choy, S. J., ... Possingham, H. P.
926 (2005). Zero tolerance ecology: Improving ecological inference by modelling the source of zero
927 observations. *Ecology Letters*, 8(11), 1235–1246.

928 Meynard, C. N., & Kaplan, D. M. (2013). Using virtual species to study species distributions and model
929 performance. *Journal of Biogeography*, 40, 1–8.

930 Meynard, C. N., Leroy, B., & Kaplan, D. M. (2019). Testing methods in species distribution modelling using
931 virtual species: what have we learnt and what are we missing? *Ecography*, 42, 2021–2036.

932 Minami, M., Lennert-Cody, C. E., Gao, W., & Román-Verdesoto, M. (2007). Modeling shark bycatch: The
933 zero-inflated negative binomial regression model with smoothing. *Fisheries Research*, 84(2), 210–221.

934 Parnell, J.A.N., Simpson, D.A., Moat, J., Kirkup, D.W., Chantaranonthai, P., Boyce, P.C., ... Wilkin, P. (2003)
935 Plant collecting spread and densities: their potential impact on biogeographical studies in Thailand.
936 *Journal of Biogeography*, 30, 193–209.

937 Pearce, J., & Ferrier, S. (2001). The practical value of modelling relative abundance of species for regional
938 conservation planning: A case study. *Biological Conservation*, 98(1), 33–43.

939 Pearce, J. L., & Boyce, M. S. (2006). Modelling distribution and abundance with presence-only data. *Journal of*
940 *Applied Ecology*, 43(3), 405–412.

941 Phillips, S. J., Anderson, R. P., & Schapire, R. E. (2006). Maximum entropy modeling of species geographic
942 distributions. *Ecological Modelling*, 190(3), 231–259.

943 Phillips, S. J. (2008). Transferability, sample selection bias and background data in presence-only modelling: A
944 response to Peterson et al. (2007). *Ecography*, 31(2), 272–278.

945 Phillips, S. J., & Dudík, M. (2008). Modeling of species distributions with Maxent: New extensions and a
946 comprehensive evaluation. *Ecography*, 31(2), 161–175.

947 Phillips, S. J., Dudík, M., Elith, J., Graham, C. H., Lehmann, A., Leathwick, J., & Ferrier, S. (2009). Sample
948 selection bias and presence-only distribution models: Implications for background and pseudo-absence
949 data. *Ecological Applications*, 19(1), 181–197.

950 Pocock, M. J. O., & Evans, D. M. (2014). The success of the Horse-Chestnut Leaf-Miner, *Cameraria ohridella*,
951 in the UK revealed with hypothesis-led citizen science. *PLoS ONE*, 9(1).

952 Reddy, S., & Dávalos, L. M. (2003). Geographical sampling bias and its implications for conservation priorities
953 in Africa. *Journal of Biogeography*, 30(11), 1719–1727.

954 Renner, I. W., Elith, J., Baddeley, A., Fithian, W., Hastie, T., Phillips, S.J., ... Warton, D. I. (2015). Point
955 process models for presence-only analysis. *Methods in Ecology and Evolution*, 6, 366–379.

956 Ridout, M., Hinde, J., & Demétrio, C. G. B. (2001). A score test for testing a Zero-Inflated Poisson regression
957 model Against Zero-Inflated negative binomial alternatives. *Biometrics*, 57(1), 219–223.

958 Rocchini, D., Hortal, J., Lengyel, S., Lobo, J. M., Jiménez-Valverde, A., Ricotta, C., ... Chiarucci, A. (2011).
959 Accounting for uncertainty when mapping species distributions: The need for maps of ignorance.
960 *Progress in Physical Geography*, 35(2), 211–226.

961 Rodríguez, J. P., Brotons, L., Bustamante, J., & Seoane, J. (2007). The application of predictive modelling of
962 species distribution to biodiversity conservation. *Diversity and Distributions*, 13(3), 243–251.

963 Schmeller, D. S., Henry, P.-Y., Julliard, R., Gruber, B., Clobert, J., Dziock, F., ... Henle, K. (2009). Advantages
964 of volunteer-based biodiversity monitoring in Europe. *Conservation Biology: The Journal of the*
965 *Society for Conservation Biology*, 23(2), 307–316.

966 Sileshi, G., Hailu, G., & Nyadzi, G. I. (2009). Traditional occupancy–abundance models are inadequate for
967 zero-inflated ecological count data. *Ecological Modelling*, 220(15), 1764–1775.

968 Smith, A. N. H., Anderson, M. J., & Millar, R. B. (2012). Incorporating the intraspecific occupancy–abundance
969 relationship into zero-inflated models. *Ecology*, 93(12), 2526–2532.

970 Sólymos, P., Lele, S., & Bayne, E. (2012). Conditional likelihood approach for analyzing single visit abundance
971 survey data in the presence of zero inflation and detection error. *Environmetrics*, 23(2), 197–205.

972 Sullivan, B. L., Wood, C. L., Iliff, M. J., Bonney, R. E., Fink, D., & Kelling, S. (2009). eBird: A citizen-based
973 bird observation network in the biological sciences. *Biological Conservation*, 142(10), 2282–2292.

974 Syfert, M. M., Smith, M. J., & Coomes, D. A. (2013). The effects of sampling bias and model complexity on the
975 predictive performance of MaxEnt species distribution models. *PLOS ONE*, 8(2), e55158.

976 Thuiller, W., Lavorel, S., Sykes, M. T., & Araújo, M. B. (2006). Using niche-based modelling to assess the
977 impact of climate change on tree functional diversity in Europe. *Diversity and Distributions*, 12(1), 49–
978 60.

979 Tiago, P., Ceia-Hasse, A., Marques, T. A., Capinha, C., & Pereira, H. M. (2017). Spatial distribution of citizen
980 science casuistic observations for different taxonomic groups. *Scientific Reports*, 7(1), 1–9.

981 Tiago, P., Pereira, H. M., & Capinha, C. (2017). Using citizen science data to estimate climatic niches and
982 species distributions. *Basic and Applied Ecology*, 20, 75–85.

983 Vollerling, J., Halvorsen, R., & Mazzoni, S. (2019). The MIAMaxent R package: Variable transformation and
984 model selection for species distribution models. *Ecology and Evolution*, 9, 12051–12068.

985 Welsh, A. H., Cunningham, R. B., Donnelly, C. F., & Lindenmayer, D. B. (1996). Modelling the abundance of
986 rare species: Statistical models for counts with extra zeros. *Ecological Modelling*, 88(1), 297–308.

987 Wenger, S. J., & Freeman, M. C. (2008). Estimating species occurrence, abundance, and detection probability
988 using Zero-Inflated distributions. *Ecology*, 89(10), 2953–2959. <https://doi.org/10.1890/07-1127.1>

989 Williams, M. R., Yates, C. J., Stock, W. D., Barrett, G. W., & Finn, H. C. (2016). Citizen science monitoring
990 reveals a significant, ongoing decline of the Endangered Carnaby’s black-cockatoo *Calyptorhynchus*
991 *latirostris*. *Oryx*, 50(4), 626–635.

992 Wisz, M. S., Hijmans, R. J., Li, J., Peterson, A. T., Graham, C. H., Guisan, A., & NCEAS Predicting Species
993 Distributions Working Group. (2008). Effects of sample size on the performance of species distribution
994 models. *Diversity and Distributions*, 14(5), 763–773.

995 Yackulic, C. B., Chandler, R., Zipkin, E. F., Royle, J. A., Nichols, J. D., Campbell Grant, E. H., & Veran, S.
996 (2013). Presence-only modelling using MAXENT: When can we trust the inferences? *Methods in*
997 *Ecology and Evolution*, 4(3), 236–243.

998 Zeileis, A., Kleiber, C. & Jackman, S. (2008). Regression Models for Count Data in R. *Journal of Statistical*
999 *Software*, 27, 1–25.

1000 Zuquim, G., Tuomisto, H., Jones, M. M., Prado, J., Figueiredo, F. O. G., Moulatlet, G. M., ... Emilio, T. (2014).
1001 Predicting environmental gradients with fern species composition in Brazilian Amazonia. *Journal of*
1002 *Vegetation Science*, 25(5), 1195–1207.

1003 Zuur, A. F., Ieno, E. N., Walker, N. J., Saveliev, A. A., & Smith, G. M. (2009). Zero-Truncated and Zero-
1004 Inflated models for count data. In *Mixed effects models and extensions in ecology with R* (pp. 261–
1005 293). Springer, New York, NY.

1006

1007

1008

1009

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019 **Tables**

1020

1021 *Table 1. Sources of zero-inflation in the simulated species occurrence data.*

Species	Source of zero-inflation		
	True abundance (before sampling)	Random sampling	Biased sampling
Random	No zero-inflation	Sampling	Sampling
Altitude	Biological	Biological and sampling	Biological and sampling
Altitude 2	Biological	Biological and sampling	Biological and sampling

1022

1023

1024 *Table 2. Ten predictor combinations were considered when modelling the simulated species distributions. Four*
 1025 *Generalised Linear Model (GLM) and six Zero-Inflated (ZI) model structures were considered using*
 1026 *combinations of the biological predictors (either altitude or altitude_randomised) and the bias predictor (distance*
 1027 *from nearest town), including different combinations in the count and zero components of the ZI models.*

Model	Predictors (GLM/ ZI Count component)	Predictors (ZI Zero component)
GLM1	Null (No predictors)	N/A
GLM2	Biased	N/A
GLM3	Biological	N/A
GLM4	Biological + bias	N/A
ZI1	Null (No predictors)	Null
ZI2	Biological + bias	Biological
ZI3	Biological	Biological + bias
ZI4	Biological	Biological
ZI5	Bias	Bias
ZI6	Biological + bias	Biological + bias

1028

1029

1030

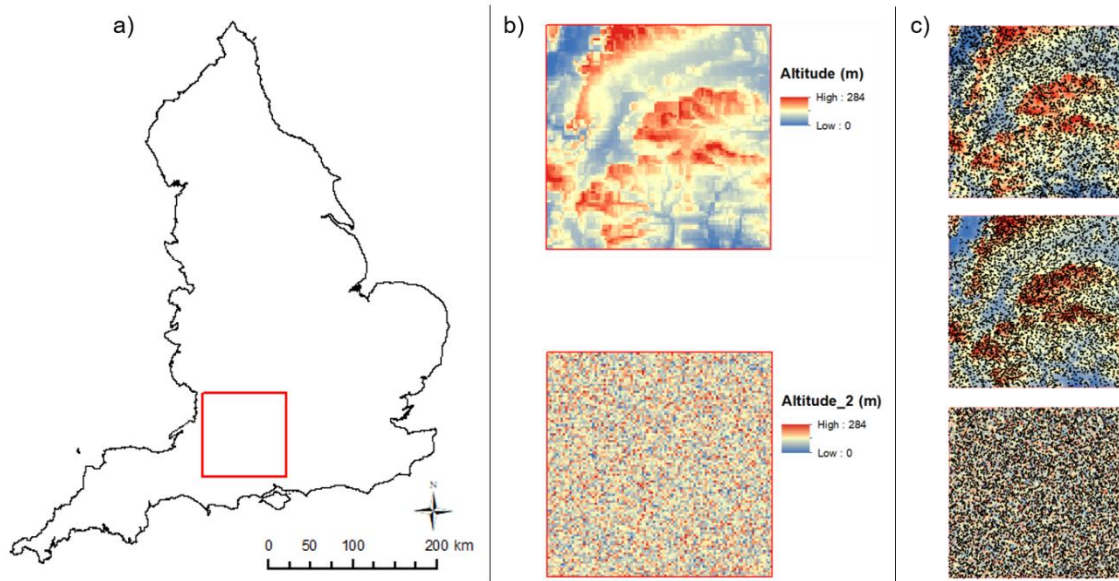
1031

1032

1033

1034

1035



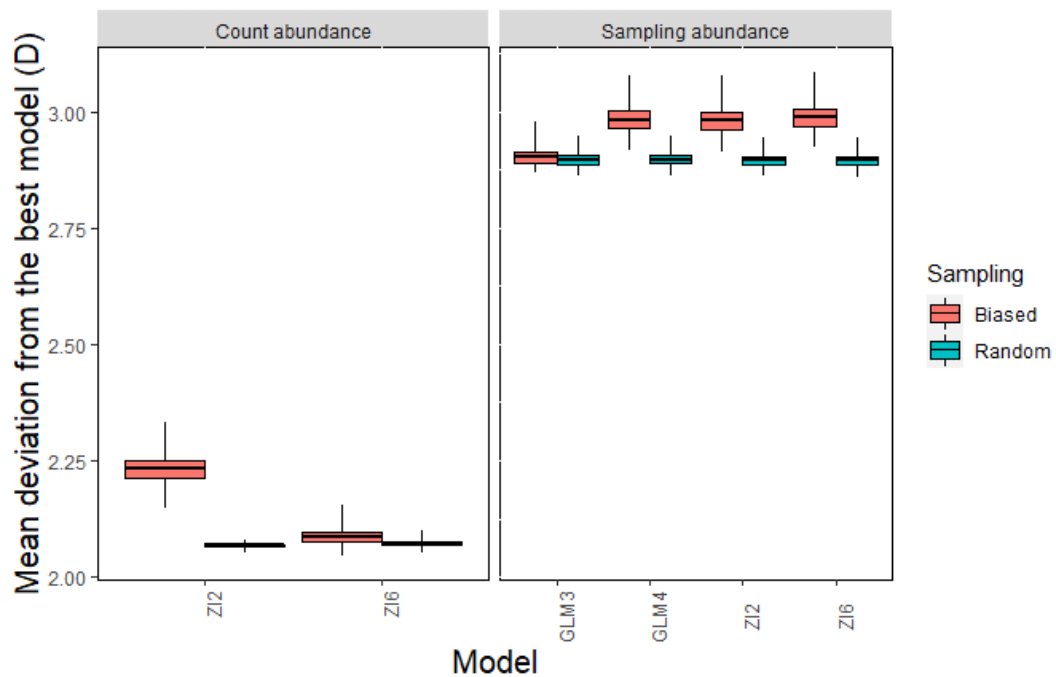
1037

1038 *Figure 1a) Simulation study area consisting of a group of 100 x 100 grid squares of 1 km² size randomly placed*
 1039 *within England covering a total area of 10,000 km² (outlined in red) (left). b) The biological predictors used to*
 1040 *fit the models: altitude (m) (top) and altitude_randomised (m) (randomised altitude layer with no spatial*
 1041 *autocorrelation, labelled here as 'altitude_2') (bottom) shown for the study area. c) A simulated species with*
 1042 *5,000 occurrence points showing no preference for altitude (random species) (top), a preference for high altitudes*
 1043 *based on a logarithmic scaler of altitude (altitude species) (middle), and a preference for high altitudes based on*
 1044 *a logarithmic scaler of altitude_randomised (altitude_randomised species) (bottom). Although the occurrence*
 1045 *points in this bottom map appear randomly distributed across the study area (similar to that of the 'random*
 1046 *species' in the top map), it is actually only the altitude values that are randomly distributed (in the 'altitude*
 1047 *randomised' layer): species occurrences are still placed based on preferring high altitude values, so are not*
 1048 *actually random with respect to the environmental predictor.*

1049

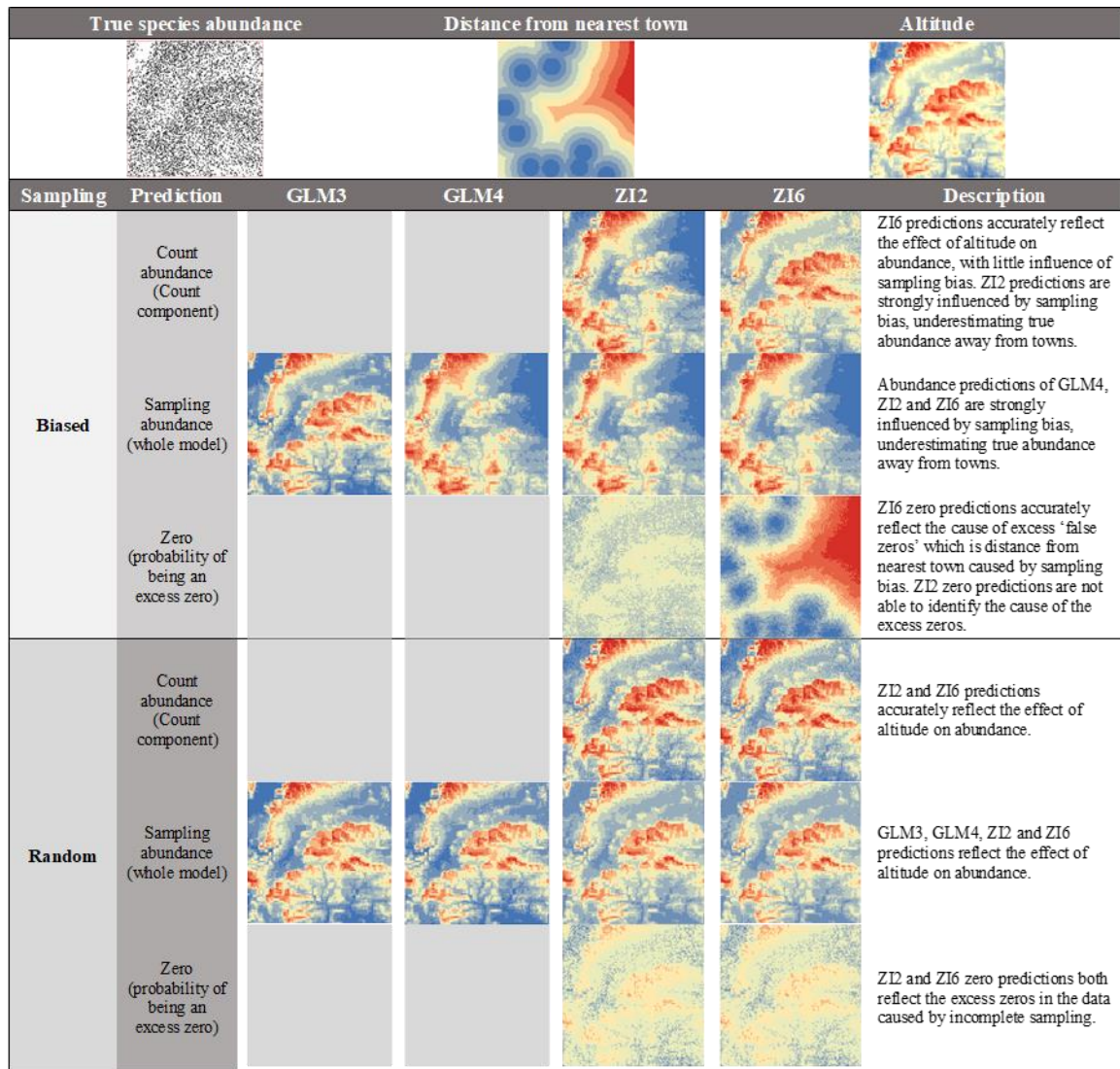
1050

1051



1052
 1053
 1054
 1055
 1056
 1057
 1058
 1059
 1060
 1061
 1062
 1063
 1064

Figure 2. Evaluation of abundance predictions (based on $D = \text{'deviation from the best model'}$) for a hypothetical organism with occurrences simulated based on a preference for high altitudes (altitude species). Mean D values ($\pm SE$ and data range) are shown for each sampling strategy (random or biased) across the 10 model repetitions. Panel a) (left) shows the evaluation of the count abundance predictions (from the zero-inflated (ZI) model count component only) from one model that accounts for sampling bias in the zero component (ZI6) and one model that does not (ZI2). Panel b) (right) shows the evaluation of the sampling abundance predictions (predictions from the whole model, and thus can be obtained from both ZI models and Generalised Linear Models (GLMs)) for four models: the same two ZI models as in Panel a), along with two GLMs: GLM3 including only the biological predictor and GLM4 including the biological and bias predictor. Only sampling abundance can be obtained from the GLMs, hence why Panel a) only shows results from the ZI models.



1065

1066

Figure 3. Example maps of abundance for a hypothetical species ('altitude species') whose occurrence is positively influenced by altitude, produced from two generalised linear models (GLMs) and two Zero-Inflated (ZI)

1067

models. Models were built with either data collected by randomly sampling grid cells (random) or with sampling bias (biased). Abundance maps from GLM3 (including the biological predictor only) and GLM4 (including both the biological and bias predictor) are produced using sampling abundance predictions (i.e. from the whole model).

1068

Both count abundance and sampling abundance predictions can be produced from the ZI models along with a map of the probability a cell is an excess zeros (zero). Both ZI models include a biological predictor (altitude) of both abundance and excess zeros, and bias predictor (distance from the nearest town) of abundance.

1069

ZI6 also includes 'distance from the nearest town' as a predictor of excess zeros. Individual cells are colour-coded based on abundance for the abundance predictions or on probability of being an excess zero for the zero predictions (high = red, low = blue).

1070

1071

1072

1073

1074

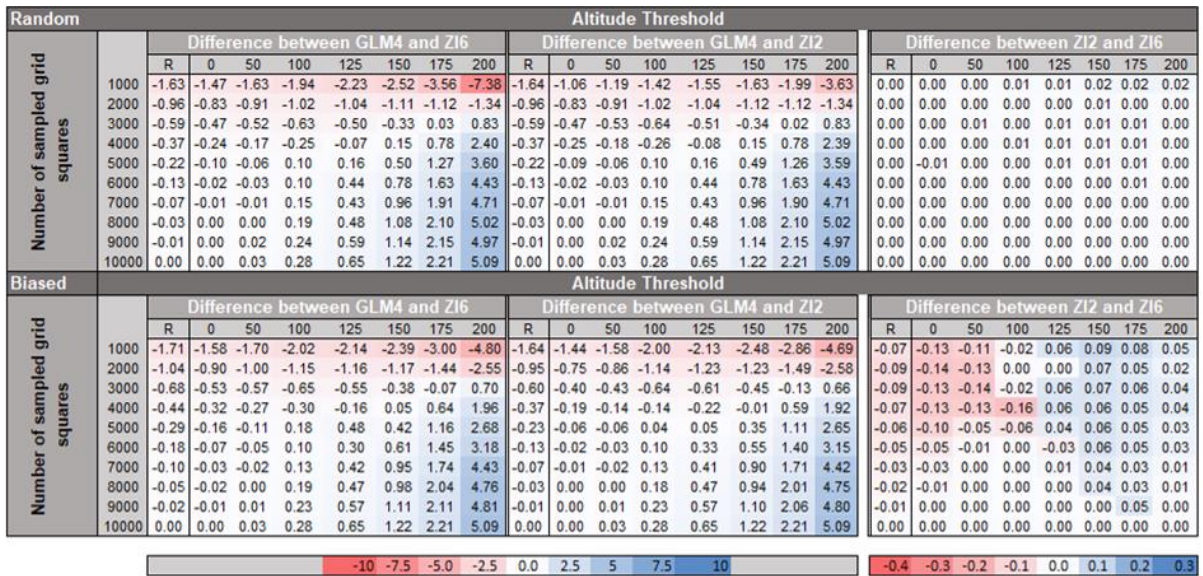
1075

1076

1077

1078

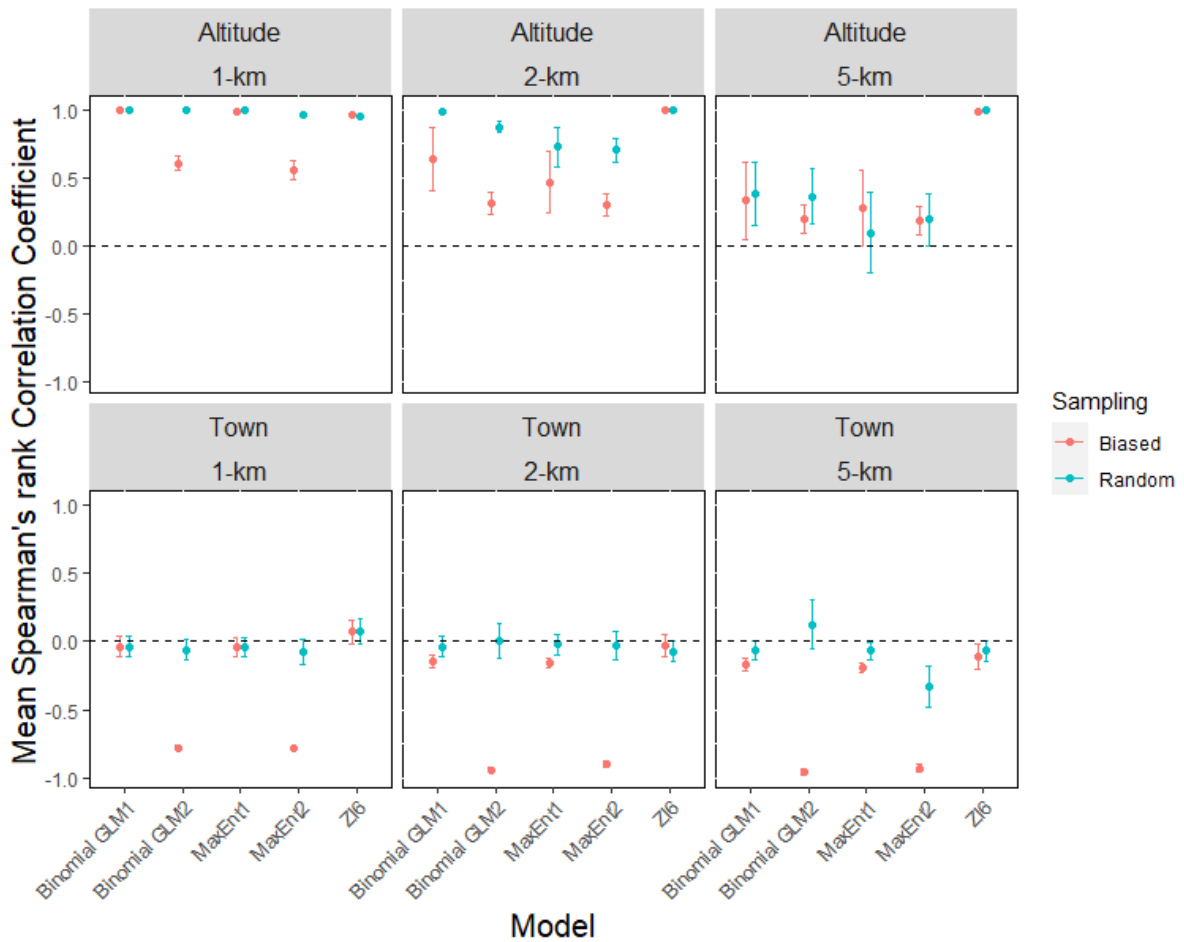
1079
 1080
 1081
 1082



1083

1084
 1085
 1086
 1087
 1088
 1089
 1090
 1091
 1092
 1093
 1094
 1095

Figure 4. Comparisons of model predictive power of true abundance between a generalised linear model (GLM) and two zero-inflated (ZI) models across varying levels of biological and sampling bias zero-inflation. Values represent the mean difference in D ('deviation from the best model') between GLM4 (containing both biological and bias predictor), ZI2 (excludes the bias predictor from the zero component) and ZI6 (includes the bias predictor in the zero component). Biological zero-inflation was increased by introducing a minimum altitude threshold below which the species cannot survive and therefore reducing its environmental niche. Sampling-related zero-inflation was increased by increasing the number of grid cells sampled across the study area in increments of 10%. Negative (red) values show scenarios where the ZI model performs better than the GLM (left and middle panels) or where ZI6 performs better than ZI2 (right panel), whereas positive (blue) values show scenarios where GLM4 outperforms the ZI models or ZI2 outperforms ZI6. 'R' represents the values for the random species whose occurrence is not related to altitude.



1096

1097

1098

1099

1100

1101

1102

1103

1104

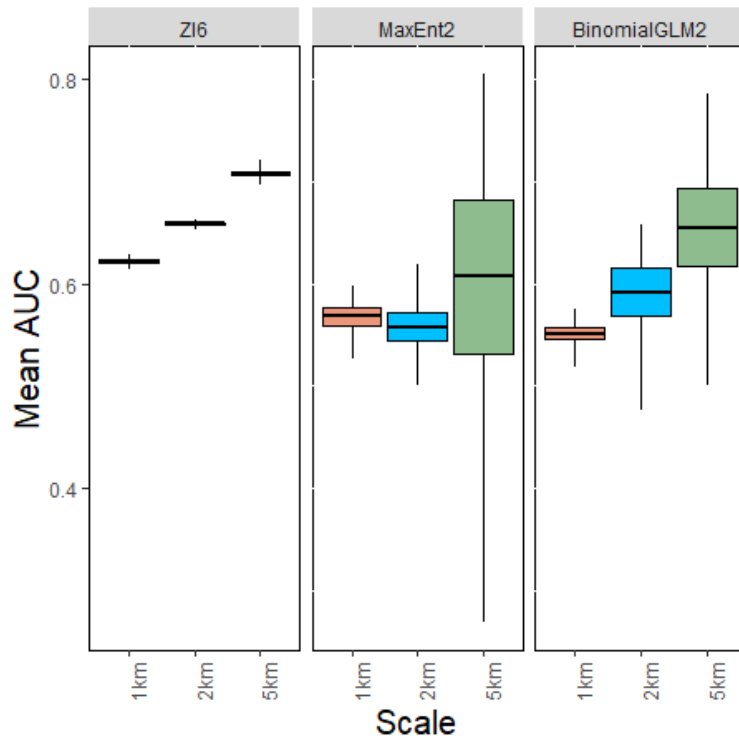
1105

1106

1107

1108

Figure 5. Mean spearman's Rank correlation coefficients (r_s) ($\pm SE$) between the model predictors (altitude and distance from nearest town) and model predictions for altitude species across three modelling scales: 1-km, 2-km and 5-km and two sampling strategies (random and biased). Three types of model are compared: 1) binomial generalised linear models (GLMs) that predict the probability of occurrence, 2) Maximum Entropy (MaxEnt) models that predict the probability of occurrence and 3) zero-inflated (ZI) models that predict the true (count) abundance of the species. Binomial-GLM1 and MaxEnt1 include only the biological predictor in the model, whereas Binomial-GLM2 and MaxEnt2 include both the biological and bias predictor. ZI6 model includes the bias and biological predictor in both the count and zero component.



1109

1110

1111 *Figure 6. Evaluation of MaxEnt, generalised linear model (GLM) and zero-inflated (ZI) model predictions of*
 1112 *altitude species presence-absence sampled using a biased strategy across the study area. Mean Area under the*
 1113 *Curve (AUC) (\pm SE and data range) across the 10 model repetitions is used to evaluate predictions across three*
 1114 *scales of data aggregation: 1-km, 2-km and 5-km. Three models are compared: 1) a ZI model able to account for*
 1115 *the bias in the zero component (ZI6) (see Methods for more information on the conversion of ZI abundance*
 1116 *predictions to presence-absence), 2) a MaxEnt model (MaxEnt2) that includes altitude and distance from town as*
 1117 *predictors and 3) a binomial GLM (Binomial-GLM2), also including altitude and distance from town as*
 1118 *predictors.*

1119