



HAL
open science

Use of ambiguous detections to improve estimates from species distribution models

Julie Louvrier, Anja Molinari-jobin, Marc Kéry, Thierry Chambert, David Miller, Fridolin Zimmermann, Eric Marboutin, Paolo Molinari, Oliver Müller, Rok Černe, et al.

► To cite this version:

Julie Louvrier, Anja Molinari-jobin, Marc Kéry, Thierry Chambert, David Miller, et al.. Use of ambiguous detections to improve estimates from species distribution models. *Conservation Biology*, 2018, 33, pp.185 - 195. 10.1111/cobi.13191 . hal-03502432

HAL Id: hal-03502432

<https://hal.science/hal-03502432>

Submitted on 4 Jan 2022

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 Use of ambiguous detections to improve estimates from species distribution models.

2 Julie Louvrier^{1,2}, Anja Molinari-Jobin³, Marc Kéry⁴, Thierry Chambert¹, David Miller⁵, Fridolin
3 Zimmermann³, Eric Marboutin⁶, Paolo Molinari⁷, Oliver Müller⁸, Rok Černe⁹, Olivier
4 Gimenez¹

5
6 ¹CEFE, Univ Montpellier, CNRS, Univ Paul Valéry Montpellier 3, EPHE, IRD, Montpellier, France.

7 ²Office National de la Chasse et de la Faune Sauvage, CNERA prédateurs et animaux déprédateurs, Parc
8 Micropolis, 05000, Gap, France.

9 ³KORA, Thunstrasse 31, 3074, Muri, Switzerland.

10 ⁴Swiss Ornithological Institute, 6204, Sempach, Switzerland.

11 ⁵Department of Ecosystem Science and Management, Pennsylvania State University, University Park, PA, 16802,
12 U.S.A.

13 ⁶ONCFS, Gières, France.

14 ⁷Italian Lynx Project, 33018, Tarvisio, Italy.

15 ⁸National Office of Forests, Vaduz, Liechtenstein.

16 ⁹Slovenia Forest Service, Ljubljana, Slovenia.

17

18 Abstract

19 As large carnivores recover throughout Europe, there is a need to study their distribution to
20 determine their conservation status and assess the potential for conflicts with human
21 activities. However, efficient monitoring of many large carnivore species is challenging due
22 to their rarity, elusive behavior and large home range size. In Europe, most current monitoring
23 protocols rely on multiple detection methods, which can include opportunistic sightings from
24 citizens in addition to designed surveys. Two types of detection errors may occur in such
25 monitoring schemes; false negatives and false positives. When not accounted for, both can
26 bias estimates from species distribution models (SDMs). False negative detections can be
27 accounted for in SDMs that deal with imperfect detection. In contrast, false positive
28 detections, due to species misidentification, have only rarely been accounted for in SDMs.
29 Generally, researchers use *ad hoc* methods to avoid false positives through data filtering to
30 discard ambiguous observations prior to analysis. These practices may discard valuable
31 ecological information on the distribution of a species. Here, we investigated the costs and
32 benefits of including data types that might include false positives rather than discard them for
33 SDMs of large carnivores. We showcase a dynamic occupancy model that simultaneously

34 accounts for false negatives and positives to jointly analyze data that include both
35 unambiguous detections and ambiguous detections. Using simulations, we show that the
36 addition of ambiguous detections increases the precision of parameter estimates. The analysis
37 of data on the Eurasian lynx (*Lynx lynx*) suggested that incorporating ambiguous detections
38 produced more precise estimates of the ecological parameters and revealed additional
39 occupied sites in areas where the species is likely expanding. Overall, our work shows that
40 ambiguous data should be considered when studying the distribution of large carnivores,
41 through the use of dynamic occupancy models accounting for misidentification.

42

43 **Introduction**

44 The distribution and abundance of large carnivores in many parts of the world has been
45 declining for centuries because of habitat loss and human persecution (Ripple et al. 2014).
46 Thanks to active conservation measures, several species of large carnivores have recently
47 been expanding their ranges substantially in Europe. As a result, most European countries
48 currently host at least one viable population of large predators (Chapron et al. 2014). This
49 recent expansion led to the emergence of conflicts with humans (Ripple et al. 2014). In this
50 context, accurate distribution mapping, i.e., species distribution models (SDMs; Elith &
51 Leathwick 2009), is essential to determine the conservation status and recovery success
52 (IUCN, 2012), to target potential areas of occurrence and understand large carnivores range
53 dynamics, identify the possible areas where they might be recovering in the future (Chapron
54 et al., 2014) and mitigate conflicts often associated with the recovery of large carnivores
55 (Guillera-Arroita et al. 2015) like, e.g., livestock depredation related to wolves' recolonization
56 (Marucco & Mcintire, 2010). However, their rarity, elusive behavior and low density render
57 efficient monitoring of large carnivores difficult (Ripple et al. 2014).

58 The monitoring of large carnivores in Europe relies on several survey methods that are
59 implemented by professionals and members of the public ("citizens"). Citizens in particular
60 add to the ability to survey large spatial coverage over extended periods in time, which would
61 be costly if done by professionals only (Molinari-Jobin et al. 2017). The practice of engaging
62 the public in a project that produces reliable data and information usable by scientists and/or
63 decision-makers is a primary goal of citizen science (CS; McKinley et al. 2017). CS is
64 becoming an important tool in ecology to study the distribution, abundance and species
65 richness of plants and animals (Silvertown 2009; Dickinson et al. 2012). However, CS-
66 generated data present potential quality issues especially when the goal is to build SDMs.

67 Difficulty detecting large carnivores means that animals can be missed at sites where
68 they are present (i.e., producing false negative observations). Occupancy models were
69 developed to deal with false negative errors (Guillera-Arroita 2017) and are recommended for
70 analyzing CS data (Isaac et al. 2014). While datasets produced by CS have been proven
71 valuable (Kosmala et al. 2016), professionals may present a better expertise than citizens to
72 detect or identify the species of interest, diminishing the risk of identification errors
73 (Fitzpatrick et al. 2009). False positives can occur when the species of interest is "detected" at
74 a site where it does not occur, resulting from misidentification (Miller et al. 2011). Recent
75 studies have demonstrated the importance of accounting for misidentification for SDMs
76 (Miller et al. 2011, 2013; Chambert et al. 2015). Ignoring misidentification may lead to
77 overestimating species range (Royle & Link 2006; McClintock et al. 2010).

78 Methods of observations typically used to survey large carnivores are based on
79 indirect observation methods through signs of presence such as tracks, prey remains, camera-
80 trap photos, or scats (Molinari-Jobin et al. 2017). Observations then go through a filtering
81 process performed by experts to assess the reliability of evidence of presence. Recent studies
82 of the distribution of European large carnivores were based only on the reliable observations,

83 i.e., those remaining after discarding ambiguous detections and validated by experts
84 (Molinari-Jobin et al. 2017). This means that part of the observations may end up being
85 discarded, even though they may contain relevant ecological information on the species
86 distribution. This raises the question whether this information can somehow be extracted and
87 made useful in the context of SDMs?

88 Here, we investigated the pros and cons of removing ambiguous detections in SDMs
89 of large carnivores versus keeping all records and formally accounting for misidentification.
90 We showcase a dynamic occupancy model accounting for both false negative and false
91 positive errors (Miller et al. 2011, 2013) to jointly analyze unambiguous and ambiguous
92 detections. To assess the performance of this approach, we performed a simulation study
93 which compares the analysis of unambiguous and ambiguous detections *vs.* using
94 unambiguous detections only.

95 We illustrate these methods in a case study with a SDM of the Eurasian lynx (*Lynx*
96 *lynx*) throughout the European Alps (Molinari-Jobin et al. 2017). Observations differ in their
97 reliability in terms of the likely incidence of false positives. Ambiguous detections, which are
98 usually discarded, represent almost a third of all observations in the dataset and have a larger
99 geographic range than unambiguous detections. We expected improved precision in
100 ecological parameter estimates when all data were included in an analysis, despite having to
101 accommodate additional nuisance parameters to deal with misidentification.

102

103 **Material and Methods**

104 1- Occupancy model accounting for misidentification

105 Dynamic occupancy models allow the estimation of occupancy and its temporal dynamics as
106 a function of local extinction and colonization probabilities, while accounting for imperfect

107 species detection (MacKenzie et al. 2003). These models can be formulated as state-space
 108 models to separate the state process, whether a species is present or not at a site and how that
 109 changes through time, from the observational process, whether the species is observed at a
 110 site during a given period depending on whether or not it was actually present (Royle & Kéry
 111 2007). We define $z_{i,1}$ as the initial latent occurrence state of site i (with $z = 1$ denoting
 112 presence and $z = 0$ absence), and $z_{i,t}$ the latent state for of site i at time t . The state process is
 113 initiated by the initial occupancy probability $\psi_{i,1}$ for site i , then governed by colonization
 114 probability $\gamma_{i,t}$ (the probability that a site i that is not occupied at time t will become occupied
 115 at time $t+1$), and extinction probability $\varepsilon_{i,t}$ (the probability that an occupied site i at time t will
 116 become unoccupied at time $t+1$). We model $z_{i,1}$ as a draw from a Bernoulli distribution with
 117 probability $\psi_{i,1}$. All subsequent latent states $z_{i,t}$ for $t > 1$ are draws from another Bernoulli
 118 distribution that combines both possible extinction and colonization events:

$$119 \quad z_{i,t+1}|z_{i,t} \sim \text{Bernoulli}(z_{i,t}(1 - \varepsilon_{i,t}) + (1 - z_{i,t})\gamma_{i,t}).$$

120 If a site is occupied in year (or season) t it will still be occupied with probability $1 - \varepsilon_{i,t}$ or if it
 121 is unoccupied it will become occupied with probability $\gamma_{i,t}$. Each site is surveyed during
 122 secondary occasions (or survey) j within year (or season) t . Site occupancy models rely
 123 satisfying the site closure assumption, whereby the latent occurrence state of a site does not
 124 change within a sampling season, whereas occupancy dynamics (colonization, extinction)
 125 happen between years (or seasons).

126 In addition to the state process, the observation process leads to the data $y_{i,j,t}$: the
 127 observed state of site i during secondary occasion (or survey) j within year (or season) t .
 128 Hereafter, we drop the indices when possible to ease the reading. In our study, $y = 0$ denoted
 129 no detection, $y = 1$ an unambiguous detection and $y = 2$ an ambiguous detection. To account
 130 for unambiguous and ambiguous detections, we followed the formulation of Miller et al.

131 (2013). We defined an additional parameter $d_{i,j,t}$ which took the value of 1 if any detection
132 (ambiguous or unambiguous) was made at site i during survey j within year t , and 0 if not. For
133 occupied sites, by definition $d = 1$ denoted a true detection while for unoccupied sites, $d = 1$
134 denoted a false positive detection. For both occupied and unoccupied sites, $d = 0$ meant no
135 detection was made hence, $y = 0$. At an occupied site, the possible observations are: no
136 detection ($y = 0$), unambiguous detection ($y = 1 | d = 1$) or ambiguous detection ($y = 2 | d = 1$).
137 For occupied sites, the probability of a true detection (i.e. $d = 1$) during a secondary sampling
138 occasion (or survey) is defined as $P(d = 1 | z = 1)$, hereafter written as p_{11} . The probability that
139 a true detection will be classified as unambiguous is given by the probability $P(y = 1 | d = 1)$
140 hereafter written as b . The probability of an unambiguous detection is $p_{11}b$ and the probability
141 for an ambiguous detection (i.e. $y=2$) for an occupied site is $p_{11}(1-b)$. For unoccupied sites (i.e.
142 $z = 0$), by definition, unambiguous detections ($y = 1 | d = 1$) do not occur thus, only two
143 possible observations can be made: an ambiguous detection ($y = 2 | d = 1$), which in this case
144 is a false positive, or no detection ($y = 0$). The probability of a false positive detection (i.e. $d =$
145 1) occurring at an unoccupied site i during a secondary sampling occasion (or survey) j is
146 $P(d=1 | z=0)$, hereafter written as p_{10} . Then the probabilities, unconditional on state z of a site,
147 of recording the three possible observed states (y) are:

$$148 \quad P(y = 0) = P(z = 1)P(d = 0 | z = 1) + P(z = 0)P(d = 0 | z = 0)$$

$$149 \quad = \psi(1 - p_{11}) + (1 - \psi)(1 - p_{10}) \text{ for no detection;}$$

$$150 \quad P(y = 1) = P(z = 1)P(d = 1 | z = 1)P(y = 1 | d = 1)$$

$$151 \quad = \psi p_{11} b \text{ for unambiguous detection;}$$

$$152 \quad P(y = 2) = P(z = 1)P(d = 1 | z = 1)P(y = 2 | d = 1) + P(z = 0)P(d = 1 | z = 0)$$

$$153 \quad = \psi p_{11}(1-b) + (1 - \psi)p_{10} \text{ for ambiguous detection.}$$

154

155 2- Simulations

156 We conducted a simulation study to examine the performance of a dynamic occupancy model
157 that also accounts for possible false positives (MUA – “Model Unambiguous/Ambiguous”) in
158 comparison with the dynamic occupancy model that only accounts for false-negatives, i.e.
159 fitted with unambiguous data only (MU – “Model Unambiguous”). To assess the ability of
160 both models to estimate ecological parameters, we defined four scenarios in which parameters
161 which control false positive detections and true detections varied (Table 1).

162 [Table 1 about here]

163 First, because the ecological parameters have an influence on the amount of detections
164 produced, we chose two main situations in which the occupancy probability is either “high”
165 or “low”. In the “high” occupancy scenario, we set the initial occupancy probability ψ_1 at 0.8,
166 the colonization probability γ at 0.4 and extinction probability at 0.1 to maintain a high
167 occupancy probability. This scenario would correspond to a fairly well-established species
168 reflected by its high occupancy probability across time. In the “low” occupancy scenario, we
169 set the initial occupancy probability ψ_1 at 0.1, the colonization probability γ at 0.1 and
170 extinction probability at 0.1 to maintain a low occupancy probability. This scenario would
171 correspond to a rare species with a low occupancy probability across time.

172 The detection parameters also have an influence on the amount of false positive and true
173 positive detections. First, true detections are controlled by the true detection probability p_{11}
174 and the probability to classify a true detection as unambiguous b . Therefore, in both “high”
175 and “low” occupancy scenarios, we consider two situations in which b is either “high” (i.e. set
176 at 0.8) or “low” (i.e. set at 0.5), leading to four scenarios. For all scenarios, we set p_{11} at 0.4.
177 When b is equal to 0.8, most of the true detections are classified as unambiguous. This

178 scenario would correspond to the monitoring of a species that is not easily mistaken for
179 another one or done by people trained to recognize accurately the presence signs of the
180 species. When b is equal to 0.5, a larger part of the true detections is classified as ambiguous.
181 This scenario would correspond to the monitoring of a species that can easily be mistaken or
182 done by untrained people, for instance from the general public. Second, the amount of false
183 positive detections is controlled by the false positive detection probability p_{10} . In all four
184 scenarios, we looked at how the models performed under seven different values of p_{10} ,
185 varying from 0.01 to 0.3, leading to twenty-eight different simulation scenarios. Finally,
186 because our main objective was to assess the effect of accounting for ambiguous data,
187 environmental variation was not included into our simulation study.

188 For ease of reading, the “high” occupancy “high” b scenario will be referred to as HH; the
189 “high” occupancy “low” b will be HL; the “low” occupancy “high” b will be LH; and the
190 “low” occupancy “low” b will be LL.

191 In our simulations, we generated data for 100 sites over 5 years and 3 surveys. To
192 remain realistic in the simulations, the number of surveys were chosen to mimic the case
193 study characteristics. For each scenario, we simulated $S = 500$ datasets and we fitted both
194 models to each dataset. For the initial occupancy probability ψ_1 , the colonization probability γ
195 and the extinction probability ε in both models in each scenario, we calculated the relative
196 bias and mean squared error (MSE).

197

198 3- Case study: Eurasian lynx in the Alps 1995–2014

199 After its total eradication in the Alps by around 1930, the Eurasian lynx (*Lynx lynx*) has been
200 reintroduced multiple times between 1970 and today in Switzerland, Italy, Austria and
201 Slovenia (Molinari-Jobin et al. 2017). In the 1990s, experts from the seven Alpine countries

202 set up the international lynx monitoring program SCALP (Status and Conservation of the
203 Alpine Lynx Population). The monitoring of the elusive lynx, relies on a network of > 1300
204 trained experts (game wardens, hunters, and naturalists) covering seven Alpine countries.
205 Signs of presence were classified into three reliability categories: C1 included “hard facts”
206 data, e.g. dead lynx, lynx removed from the wild as young orphans and put into captivity,
207 lynx photos and a few genetic samples, C2 are detections that were confirmed by a lynx
208 expert, (all livestock killed by lynx that was compensated, verified wild prey remains, and
209 tracks) and C3 are data that could not be verified by experts (unverified tracks and wild prey
210 remains) and unverifiable data such as any sighting, scats and vocalizations. We treated C1
211 and C2 data as unambiguous detections, assuming there were no false-positive detections in
212 these data, while the C3 data were treated as ambiguous detections. From 1995 to 2014, 8415
213 observations (67%) were classified as unambiguous detections and 3991 (33%) as ambiguous.
214 If unambiguous and ambiguous detections occurred at a site, we accounted for the
215 unambiguous detections only. Non-detections were obtained on the sites that were sampled
216 but where no lynx presence was reported during a survey within a year. In Molinari-Jobin et
217 al. (2017), a dynamic occupancy model was fitted using unambiguous detections only (i.e.,
218 using our model MU) to assess the effects of environmental covariates on different
219 parameters of the model and to assess distribution-based population trends. A 10 x 10 km grid
220 was used to define the distribution units which correspond to the approximate size of female
221 lynx’ home-range in the Alps (Molinari-Jobin et al. 2017). Surveys were defined as three
222 replicated two-month periods: November-December; January-February; and March-April.
223 Here, we used the same data set as did Molinari-Jobin et al. (2017), but in addition we also
224 used the C3 data and fitted a dynamic occupancy model that combined both unambiguous and
225 ambiguous detections (MUA). In addition, we used the same covariates for the parameters
226 that are in common in the models MU and MUA. We considered the effects of forest cover

227 and distance to the release site on ψ_1 ; the effects of year, forest cover, and number of observed
228 occupied contiguous neighbors on ε ; and the same effects plus that of human density and
229 elevation on γ .

230 For the new parameters in MUA, p_{11} and p_{10} , we used the effect of elevation and forest
231 cover and a random site-by-winter effect to accommodate unmodeled spatial heterogeneity in
232 detection rates in every combination of site and winter. A “network” covariate was also
233 included to account for heterogeneity in sampling effort in time and space. This covariate
234 took the following values based on the amount of effort for the location and time period – 0:
235 no information was available regarding the sampling effort in which case we assumed that it
236 was small but never exactly null, owing to the large number of observers and organizations
237 that collaborate in the Alpine lynx monitoring (Molinari-Jobin *et al.*, 2012); 1: trained lynx
238 monitoring network were present on the site; and 2: experienced lynx monitoring network
239 members were actively searching for lynx signs. We also considered a linear year effect, i.e.,
240 an annual trend, on p_{10} to investigate whether this probability decreased as observers gained
241 experience over time. Finally, we kept the probability b to classify a true positive detection as
242 unambiguous constant. We considered the effect of a covariate as “significant” if its 95%
243 credible interval (CRI) did not overlap 0.

244 To evaluate the added value of incorporating the C3 data (ambiguous detections) into
245 the analysis, we compared the maps of occupancy produced by the two models by calculating
246 and mapping the difference in the site- and year-specific estimates of realized occurrence $\hat{z}_{i,t}$
247 $(MU) - \hat{z}_{i,t}(MUA)$.

248 We provide the codes to run the simulations and fit the models described above in Appendix
249 S1 and Appendix S2.

250

251 Results

252 1- Simulations

253 When looking at the MSE, MUA performed better than MU in all 4 scenarios when the
254 probability of false positive detection p_{10} was below or equal to 0.15 (Appendix S3). Above
255 this value of p_{10} both models performed equally well except in one scenario and for one
256 parameter when estimating the ecological parameters: MUA estimated the colonization
257 probability γ less precisely than MU only in the HL scenario for values of p_{10} between 0.20
258 and 0.30. MSE was at its highest value, varying between 0.04 and 0.25 in the HL scenario,
259 then between 0.04 and 0.14 in the HH scenario. MSE was at its lowest value in the LH
260 scenario, varying between 0.02 and 0.06, then between 0.02 and 0.11 in the LL scenario.

261 Both models estimated the initial occupancy probability ψ_1 and γ with biases below or equal
262 to 5% in the three scenarios HH, HL, LH (Appendix S3). In the LL scenario, MU estimated
263 ψ_1 with a bias above 5% (up to 8%) and MUA had a lower bias than MU. Finally, for the
264 extinction probability ε , MUA performed better or equivalently above 5% in terms of bias in
265 the scenarios HH and HL, and worse or equivalently above 5% in the LH and LL scenarios.

266

267 2- Lynx case study

268 When we fitted the MUA with both unambiguous and ambiguous detections (i.e., for C1,
269 C2 and C3 data), the true detection probability, p_{11} , was higher on sites with a high forest
270 cover, and appeared to vary according to the season and network (Table 2). Elevation had no
271 effect on p_{11} . The false positive detection probability, p_{10} , was higher on sites with a high
272 forest cover and varied according to network (Table 2). While elevation and season had no
273 significant effect on p_{10} , we found that this probability decreased with time (Table 2). Both

274 models gave similar estimates for ψ_1 , ε and γ but MUA produced more precise estimates than
275 MU (Appendix S4).

276 The probability b of classifying a true detection as unambiguous was estimated at 0.81
277 with high precision (CRI 0.79 - 0.83). At the beginning of the study period, in the winter
278 1995/1996, we estimated the mean occupancy probability ψ over all sites at 0.04 (CRI 0.03-
279 0.07), p_{11} was estimated on average at 0.11 (CRI 0.10 - 0.25) and p_{10} was estimated at 0.006
280 (CRI 0.004- 0.01). For the end of the study period, the winter 2013/2014, we estimated the
281 mean ψ at 0.1 (0.0899; 0.11), p_{11} was estimated on average at 0.17 (0.09; 0.24) and p_{10} at
282 0.007 (0.003; 0.010). MUA estimated a few more occupied sites than MU for both winters
283 1995/1996 and 2013/2014 (between 4 in 1995/1996 to 13 in 2013/2014, see Fig. 1, middle
284 and bottom panel) and estimated occupied sites that were estimated occupied by MU too. The
285 additional sites that were estimated occupied from MUA were sites where ambiguous
286 detections had occurred (Fig. 1, top panel).

287

288 [Table 2 about here]

289 [Figure 1 about here]

290

291 **Discussion**

292 Assessing the distribution of large carnivores at large scales is a central information for
293 assessing their conservation status, and abundance (IUCN, 2012; Jedrzejewski et al., 2018),
294 target potential conflict areas (Marucco & Mcintire, 2010) and understand the mechanism of
295 the distribution's dynamics for successful management (Eriksson & Dalerum, 2018).
296 Producing more precise and less biased estimates by adding ambiguous data with a model
297 accounting for false positive detections can bring new insights into species' distribution in
298 places where getting unambiguous data is challenging. Due to the large areas involved, the

299 monitoring of large carnivores in Europe relies on a large network of both professional and
300 non-professional observers (Louvrier et al. 2018; Molinari-Jobin et al. 2017). While false-
301 negative detections have received much attention in the species distribution modeling
302 literature with the rise of occupancy models (MacKenzie et al. 2003; Bailey et al. 2014),
303 dealing with ambiguous detections has been studied much less (Miller et al. 2011; Chambert
304 et al. 2015). Here, using simulations we demonstrate that jointly analyzing unambiguous and
305 ambiguous detections with the appropriate dynamic occupancy models led to increased
306 precision in the estimates of ecological parameters when p_{10} was low. When this probability
307 was above 0.20, both models estimated ecological parameters with almost equivalent
308 precision which varied between its highest values in the “high” occupancy scenarios and its
309 lowest values in the “low” occupancy scenarios. Both models produced estimates of
310 ecological parameters with low bias except for one ecological parameter in one specific
311 scenario.

312 When looking at the results of the lynx analysis, we found that adding ambiguous data
313 helped produce more precise estimates and provided additional spatial information that
314 improved inference in areas where the species likely occurred at very low density (e.g., at a
315 colonization front).

316

317 **What did we learn from the simulation study?**

318 MUA performed better than MU in most of the scenarios. Two factors seemed to have an
319 influence on models’ performances: the false-positive probability p_{10} and the occupancy
320 probability. In terms of precision, MUA performed better when p_{10} was low and performed
321 equivalently when p_{10} was high. In the case of a low occupancy probability, the estimates of
322 extinction probability were found to be more biased positively under the MUA than the MU

323 leading to an overestimation of ϵ and the distribution. For the other parameters and the other
324 scenarios, MUA produced estimates with low biases. Whether a species is occurring at “high”
325 or “low” occupancy probability can often be evaluated prior to the analyses based on the
326 knowledge of the species ecology or on previous studies. Overall, we recommend always
327 including ambiguous data, as in most of the scenarios MUA performed better than or
328 equivalently to MU in terms of both precision and bias for the ecological parameter estimates.

329

330 **Shall we account for ambiguous data when studying the distribution of large**
331 **carnivores?**

332 Using a model incorporating both unambiguous and ambiguous data, we estimated the effect
333 of several covariates on the dynamics of Lynx occupancy in the entire range of the Alps. This
334 SDM exercise allowed assessing trends in the distribution of the species, informing its
335 conservation status (Guisan & Thuiller 2005). We found covariate effects to be similar in
336 direction and magnitude to those estimated by Molinari-Jobin et al. (2017) who fitted the
337 simpler MU to the lynx data with unambiguous detections only (Table 2). We refer the reader
338 to their study for a detailed description of these effects and their possible biological
339 interpretation. Our results showed that the probability to make a false positive detection
340 decreased over time. This could be due to observers remaining in the network becoming less
341 likely to make false positive detections with time as they became more experienced in
342 recognizing the species (Jordan et al. 2012). This was corroborated by the fact that the
343 number of ambiguous detections decreased over the duration of the study period (Molinari-
344 Jobin et al. 2012). Additionally, the use of camera trapping has increased over time, leading to
345 an increasing amount of C1 detections and therefore diminishing the proportion of C3 in the
346 datasets (Molinari-Jobin et al. 2017). The learning process of citizens in scientific monitoring

347 programs has been studied in the past (Dickinson et al. 2012; Jordan et al. 2012) and it was
348 found that the general public not only learned through participation but also became more
349 aware of the general ecological issues and became more prone to understand scientific
350 research (Bonney et al. 2009). We found that the probability to make a true detection was
351 similar to the probability to detect the species in MU fitted by Molinari-Jobin et al. (2017).
352 This makes sense because the probability to detect the species in MU is equal to the
353 probability to make a true detection multiplied by the probability to classify a detection as
354 unambiguous. We also found that there was a probability of 0.8 to classify a true detection as
355 unambiguous. This may be due to the fact that observers in the network are highly competent
356 at detecting the species and produce reliable data. This could also reflect that it is relatively
357 easy to identify the signs of presence of lynx because there is almost no confusion possible
358 with other species present in the area. Whenever the focus species can be mistaken for another
359 one, if data quality is not sufficient (e.g. tracks in the snow for wolves which can be mistaken
360 for dogs), true detections can be classified as ambiguous. There can also be false positive
361 detections coming from misidentification when b is low. In this case, the amount of true
362 detections in ambiguous data will be non-negligible. In a case where b is low and only
363 unambiguous data is used, a large part of true presences can be missed and the resulting
364 distribution will be underestimated (Miller et al. 2011).

365 The occupancy estimates under both models agree to suggest that the lynx case study
366 corresponds to the LH simulation scenario (compare Table 1 to Appendix S4). When
367 inspecting the distribution maps produced by MUA, we saw that adding ambiguous detections
368 brought new and useful information. Some sites were estimated as occupied by MUA, while
369 these same sites were estimated as non-occupied by MU (Fig 1). Because of the low
370 occupancy of the lynx and its elusive behavior, the number of times the species was detected
371 was very low. Because the probability to classify a detection as unambiguous b was high,

372 only few true detections were classified as ambiguous, which might explain why adding them
373 did not change the parameter estimates but helped producing more precise estimates. In turn,
374 it provides new insights in the context of managing a protected species (Guillera-Arroita et al.
375 2015). The sites we found to be occupied thanks to the incorporation of ambiguous detections
376 could likely represent areas where the species is currently expanding. These same sites also
377 point to places where lynx have not occurred before and negative interactions might occur due
378 to the novelty of lynx presence. Sites that appeared occupied after including ambiguous data
379 can inform the prediction of location of potential conflicts. Finally, if the objective is mapping
380 the colonization front to, e.g., mitigate conflicts, ambiguous data should be included.

381

382 **Recommendations**

383 Dynamic occupancy models in general provide a powerful and natural analytical framework
384 for changing species distributions (Kéry et al. 2013). More specifically, dynamic occupancy
385 models accounting for misidentification represent a powerful method to deal with detections
386 that cannot be categorized as certain in species distribution modeling. We recommend careful
387 categorization of field observations into unambiguous or ambiguous detections, for instance
388 by using several experts to classify the detections and use a standardized filtering
389 classification process, to avoid false positive detections mistakenly classified as reliable data.
390 This filtering process also allows avoiding too many detections that cannot be verified by
391 rejecting some of them. If some detections cannot be checked by experts for instance and
392 cannot be classified as unambiguous, observers might need to visit the sites where these
393 detections were made to get more reliable detections. Even though occupancy models can
394 deal with ambiguity, efforts should be put in the survey design and data collection to avoid
395 the production of false positive detections or at least reduce their proportion. In the case of

396 analyzing data from citizen-science, models accounting for false-positive detections can be a
397 good tool to assess species distribution if a classification of detections is made (e.g.:
398 unambiguous *vs* ambiguous). In the case of a species occurring at low density such as the
399 Eurasian lynx, additional information can bring new insights into the species distribution and
400 help targeting specific sites where the species is likely to occur in the future.

401

402 **References**

403 Bailey, L.L., Mackenzie, D.I. & Nichols, J.D. (2014) Advances and applications of occupancy
404 models. *Methods in Ecology and Evolution*, **5**, 1269–1279.

405 Bonney, R., Cooper, C.B., Dickinson, J., Kelling, S., Phillips, T., Rosenberg, K. V & Shirk, J.
406 (2009) Citizen Science : A Developing Tool for Expanding Science Knowledge and
407 Scientific Literacy. **59**, 977–984.

408 Chambert, T.C., Miller, D.A.W. & Nichols, J.D. (2015) Modeling false positive detections in
409 species occurrence data under different study designs. *Ecology*, **96**, 332–339.

410 Chapron, G., Kaczensky, P., Linnell, J.D.C., von Arx, M., Huber, D., Andren, H., Lopez-Bao,
411 J. V., et al. (2014) Recovery of large carnivores in Europe's modern human-dominated
412 landscapes. *Science*, **346**, 1517–1519.

413 Dickinson, J.L., Shirk, J., Bonter, D., Bonney, R., Crain, R.L., Martin, J., Phillips, T., et al.
414 (2012) The current state of citizen science as a tool for ecological research and public
415 engagement. *Frontiers in Ecology and the Environment*, **10**, 291–297.

416 Elith, J. & Leathwick, J.R. (2009) Species Distribution Models : Ecological Explanation and
417 Prediction Across Space and Time. *Annual review of ecology, evolution, and*
418 *systematics*, **40**, 677–697.

419 Eriksson, T. & Dalerum, F. (2018) Identifying potential areas for an expanding wolf
420 population in Sweden. *Biological Conservation*, **220**, 170–181.

421 Fitzpatrick, M.C., Preisser, E.L., Ellison, A.M. & Elkinton, J.S. (2009) Observer Bias and the
422 Detection of Low-Density Populations. *Ecological Applications*, **19**, 1673–1679.

423 Guillera-Arroita, G. (2017) Modelling of species distributions, range dynamics and
424 communities under imperfect detection: advances, challenges and opportunities.
425 *Ecography*, **40**, 281–295.

426 Guillera-arroita, G., Lahoz-monfort, J.J., Elith, J., Gordon, A., Kujala, H., Lentini, P.E.,
427 Mccarthy, M.A., et al. (2015) Is my species distribution model fit for purpose ? Matching
428 data and models to applications. *Global Ecology and Biogeography*, **24**, 276–292.

429 Guisan, A. & Thuiller, W. (2005) Predicting species distribution : offering more than simple
430 habitat models. *Ecology Letters*, **8**, 993–1009.

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

434 IUCN (2012) IUCN Red List Categories and Criteria: version 3.1, 2nd edn. *Gland,*
435 *Switzerland and Cambridge, UK.*

436 Jedrzejewski, W., Robinson, H.S., Abarca, M., Zeller, K.A., Velasquez, G., Paemelaere,
437 E.A.D., Goldberg, J.F., et al. (2018) Estimating large carnivore populations at global
438 scale based on spatial predictions of density and distribution – Application to the jaguar (
439 *Panthera onca*). *PLoS ONE*, 1–25.

440 Jordan, R.C., Ballard, H.L. & Phillips, T.B. (2012) Key issues and new approaches for
441 evaluating citizen-science learning outcomes. *Frontiers in Ecology and the Environment*,

442 **10**, 307–309.

443 Kéry, M., Guisera-arroita, G. & Lahoz-monfort, J.J. (2013) Analysing and mapping species
444 range dynamics using occupancy models. *Journal of Biogeography*, **40**, 1463–1474.

445 Kosmala, M., Wiggins, A., Swanson, A. & Simmons, B. (2016) Assessing data quality in
446 citizen science. *Frontiers in Ecology and the Environment*, **14**.

447 Louvrier, J., Duchamp, C., Lauret, V., Marboutin, E., Cubaynes, S., Choquet, R., Miquel, C.,
448 et al. (2018) Mapping and explaining wolf recolonization in France using dynamic
449 occupancy models and opportunistic data. *Ecography*, **41**, 647–660.

450 Mackenzie, D.I., Nichols, J.D., Hines, J.E., Knutson, M.G. & Franklin, A.B. (2003)
451 Estimating Site Occupancy , Colonization , and Local Extinction When a Species Is
452 Detected Imperfectly. *Ecology*, **84**, 2200–2207.

453 Marucco, F. & Mcintire, E.J.B. (2010) Predicting spatio-temporal recolonization of large
454 carnivore populations and livestock depredation risk : wolves in the Italian Alps. *Journal*
455 *of Applied Ecology*, **47**, 789–798.

456 McClintock, B.T., Bailey, L.L., Pollock, K.H. & Simons, T.R. (2010) Unmodeled observation
457 error induces bias when inferring patterns and dynamics of species occurrence via aural
458 detections. *Ecology*, **91**, 2446–2454.

459 Mckinley, D.C., Miller-rushing, A.J., Ballard, H.L., Bonney, R., Brown, H., Cook-patton,
460 S.C., Evans, D.M., et al. (2017) Citizen science can improve conservation science ,
461 natural resource management , and environmental protection. *Biological Conservation*,
462 **208**, 15–28.

463 Miller, D.A., Nichols, J.D., McClintock, B.T., Grant, E.H.C., Bailey, L.L. & Weir, L.A.
464 (2011) Improving occupancy estimation when two types of observational error occur:

465 Non-detection and species misidentification. *Ecology*, **92**, 1422–1428.

466 Miller, D.A.W., Nichols, J.D., Gude, J.A., Rich, L.N., Podruzny, K.M., Hines, J.E. &
467 Mitchell, M.S. (2013) Determining Occurrence Dynamics when False Positives Occur:
468 Estimating the Range Dynamics of Wolves from Public Survey Data. *PLoS ONE*, **8**.

469 Molinari-Jobin, A., Kéry, M., Marboutin, E., Marucco, F., Zimmermann, F., Molinari, P.,
470 Frick, H., et al. (2017) Mapping range dynamics from opportunistic data: Spatiotemporal
471 modelling of the lynx distribution in the Alps over 21 years. *Animal Conservation*, 1–13.

472 Molinari-Jobin, A., Kéry, M., Marboutin, E., Molinari, P., Koren, I., Fuxjäger, C.,
473 Breitenmoser-Würsten, C., et al. (2012) Monitoring in the presence of species
474 misidentification: The case of the Eurasian lynx in the Alps. *Animal Conservation*, **15**,
475 266–273.

476 Ripple, W.J., Beschta, R.L., Fortin, J.K. & Robbins, C.T. (2014) Trophic cascades from
477 wolves to grizzly bears in Yellowstone. *Journal of Animal Ecology*, **83**, 223–233.

478 Royle, J.A. & Kéry, M. (2007) A Bayesian state-space formulation of dynamic occupancy
479 models. *Ecology*, **88**, 1813–1823.

480 Royle, J.A. & Link, A.W. (2006) Generalized site occupancy models allowing for false
481 positive and false negative errors. *Ecology*, **87**, 835–841.

482 Silvertown, J. (2009) A new dawn for citizen science. *Trends in Ecology and Evolution*, **24**,
483 467–471.

484

485

486 Table 1: Parameters values for the simulation scenarios

Scenarios	<i>Initial occupancy probability ψ_1</i>	<i>Colonization probability γ</i>	<i>Probability to classify a true detection as unambiguous b</i>	<i>False positive detection probability p_{10}</i>
“high” occupancy	0.8	0.4	0.8	0.01
“high” b (HH)				0.5
				0.10
				0.15
				0.20
				0.25
				0.30
“high” occupancy	0.8	0.4	0.5	0.01
“low” b (HL)				0.5
				0.10
				0.15
				0.20
				0.25
				0.30
“low” occupancy	0.1	0.1	0.8	0.01
“high” b (LH)				0.5
				0.10
				0.15
				0.20
				0.25
				0.30
“low” occupancy	0.1	0.1	0.5	0.01
“low” b (LL)				0.5
				0.10
				0.15
				0.20
				0.25
				0.30

487

488

489

490

491 Table 2: Parameters estimates for the detection probabilities from both dynamic occupancy models accounting for unambiguous data only and
492 accounting for unambiguous and ambiguous data; the first column corresponds to the parameters estimates for the detection probability from the
493 model with unambiguous data only, the second column correspond to the parameters estimates for the probability of correctly detecting the
494 species given a site is occupied from the dynamic occupancy model accounting for unambiguous and ambiguous data, the last columns
495 correspond to the parameters estimates for the probability of incorrectly detecting the species given a site is unoccupied; posterior means,
496 standard deviation and the lower and upper bound of the 95% Bayesian credible interval are given. Effects with 95% Bayesian credible intervals
497 that do not contain zero are in bold.

498

Model with unambiguous data only (MU)					Model with unambiguous and ambiguous data (MUA)					Model with unambiguous and ambiguous data (MUA)				
Detection probability p	mean	sd	2.5%	97.5%	true detection probability p_{11}	mean	sd	2.5%	97.5%	false positive detection probability p_{10}	mean	sd	2.5%	97.5%
Intercept	-3.88	0.46	-4.88	-3.04	Intercept	-3.14	0.47	-4.15	-2.33	Intercept	-5.37	0.29	-5.96	-4.80
Effect of elevation	-0.11	0.05	-0.20	-0.01	Effect of elevation	-0.07	0.04	-0.16	0.01	Effect of elevation	-0.02	0.06	-0.13	0.09
Effect of forest	0.63	0.07	0.50	0.75	Effect of forest	0.67	0.06	0.55	0.79	Effect of forest	0.37	0.06	0.26	0.49
Effect of season 2	0.26	0.07	0.12	0.40	Effect of season 2	0.22	0.07	0.08	0.35	Effect of season 2	-0.14	0.12	-0.38	0.09
Effect of season 3	0.42	0.07	0.28	0.57	Effect of season 3	0.42	0.07	0.29	0.56	Effect of season 3	-0.04	0.12	-0.25	0.19
Effect of network 1	1.21	0.21	0.77	1.60	Effect of network 1	0.92	0.20	0.54	1.31	Effect of network 1	0.37	0.14	0.09	0.66
Effect of network 2	2.37	0.23	1.93	2.84	Effect of network 2	1.95	0.22	1.53	2.40	Effect of network 2	2.63	0.30	2.02	3.20
Residual effect	0.86	0.10	0.65	1.04	Residual effect	0.90	0.06	0.77	1.01	Residual effect	0.42	0.21	0.18	0.92
Effect of country France	-3.49	0.31	-4.09	-2.87	Effect of country France	-2.76	0.28	-3.32	-2.21	Effect of country France	-5.55	0.25	-6.05	-5.06
Effect of country Italy	-3.40	0.24	-3.86	-2.94	Effect of country Italy	-2.85	0.23	-3.31	-2.40	Effect of country Italy	-5.55	0.24	-6.01	-5.07
Effect of country Switzerland	-2.87	0.23	-3.32	-2.42	Effect of country Switzerland	-2.23	0.22	-2.68	-1.81	Effect of country Switzerland	-4.75	0.26	-5.25	-4.25
Effect of country	-4.00	0.25	-4.48	-3.51	Effect of country	-3.25	0.21	-3.67	-2.84	Effect of country	-5.62	0.23	-6.07	-5.16

Austria					Austria					Austria				
Effect of country					Effect of country					Effect of country				
	-4.54	0.35	-5.23	-3.86		-3.83	0.34	-4.51	-3.19		-4.62	0.28	-5.17	-4.05
Slovenia					Slovenia					Slovenia				
Effect of country					Effect of country					Effect of country				
	-4.99	1.39	-8.28	-2.65		-3.93	1.52	-7.41	-1.51		-6.11	0.49	-7.21	-5.23
Germany					Germany					Germany				
										Effect of time in				
										years	-0.04	0.02	-0.09	-0.01

500

501

502