

# Estimating occupancy dynamics and encounter rates with species misclassification: a semi-supervised individual-level approach

Anna I. Spiers <sup>a,b</sup>, J. Andrew Royle <sup>c</sup>, Christa L. Torrens <sup>d</sup>, Maxwell B. Joseph <sup>a</sup>

<sup>a</sup> Earth Lab, Cooperative Institute for Research in Environmental Sciences, University of Colorado Boulder, CO, USA

<sup>b</sup> Department of Ecology and Evolutionary Biology, University of Colorado Boulder, CO, USA

<sup>c</sup> US Geological Survey Eastern Ecological Science Center

<sup>d</sup> Institute of Arctic and Alpine Research, University of Colorado, Boulder, CO, USA

## Correspondence

Anna Spiers

Email: [anna.spiers@colorado.edu](mailto:anna.spiers@colorado.edu)

## Running headline

Species classification-occupancy model

## Abstract

1  
2 1. Large-scale, long-term biodiversity monitoring is essential to meeting conservation and land  
3 management goals and identifying threats to biodiversity. However, multispecies surveys are prone  
4 to various types of observation error, including false positive/negative detection, and misclassifi-  
5 cation, where a species is encountered but its species identity is not correctly identified. Previous  
6 methods assume an imperfect classifier produces species-level classifications, but in practice, partic-  
7 ularly with human observers, we may end up with extraspecific classifications including "unknown",  
8 morphospecies designations, and taxonomic identifications coarser than species. Disregarding these  
9 types of species misclassification in biodiversity monitoring datasets can bias estimates of ecologi-  
10 cally important quantities such as demographic rates, occurrence, and species richness.

11 2. Here we develop an occupancy model that accounts for species non-detection and misclassifica-  
12 tion. Our framework accommodates extinction and colonization dynamics, allows for additional  
13 uncertain 'morphospecies' designations in the imperfect species classifications, and makes use of  
14 individual specimen with known species identities in a semi-supervised setting. We compare the  
15 performance of our joint classification-occupancy model to a reduced classification model that dis-  
16 cards information about occupancy and encounter rate on a withheld test set. We illustrate our  
17 model with an empirical case study of the carabid beetle (Carabidae) community at the National  
18 Ecological Observatory Network Niwot Ridge Mountain Research Station, west of Boulder, CO,  
19 USA, and quantify taxonomist identification error by accounting for classification probabilities.

20 3. Species occupancy varied through time and across sites and species. The model yielded high  
21 probabilities (30 to 92% medians) of classification where the imperfect classifier matched the true  
22 species. The classification model informed by occupancy and encounter rates outperformed the  
23 classification that was not, and these differences were most pronounced for abundant species.

24 4. Our probabilistic framework can be applied to datasets with imperfect species detection and clas-  
25 sification. This model can identify commonly misclassified species, helping biodiversity monitoring  
26 organizations systematically prioritize which samples need validation by an expert. Our Bayesian  
27 approach propagates classification uncertainty to offer an alternative to making conservation deci-  
28 sions based on point estimates

29 **Keywords** — carabid, imperfect classifier, morphospecies, NEON, observation error, occupancy models, semi-  
30 supervised, species misclassification

## 31 1 Introduction

32 Large-scale, long-term biodiversity monitoring is essential to meeting conservation and land management goals and  
33 identifying threats to biodiversity. Such comprehensive datasets increasingly include multispecies surveys that capture

34 information-rich co-occurrence data, enabling community-level analyses (Iknayan et al., 2014; Ovaskainen et al., 2017).  
35 However, multispecies surveys are prone to various types of imperfect detection, including false absences where a  
36 species is present but not detected (Dorazio and Royle, 2005), and misidentification, where a species is encountered  
37 but its species identity is not correctly recorded (Miller et al., 2011).

38 Occupancy models account for observation error in biodiversity surveys that seek to understand species dis-  
39 tributions, track population changes, and describe mechanisms underlying population and community dynamics  
40 (MacKenzie et al., 2002). Latent presence/absence states are modeled explicitly, with an observation model that  
41 accounts for the details of the detection process, including the potential for false negatives (non-detections at occu-  
42 pied sites) and false positives (detections at unoccupied sites) (Royle and Link, 2006; Miller et al., 2012; Chambert  
43 et al., 2015; Wright et al., 2020). Disregarding false positives in biodiversity monitoring datasets can bias estimates  
44 of ecologically important quantities such as demographic rates, occurrence, and species richness (McClintock et al.,  
45 2010; Chambert et al., 2015, 2018).

46 Multi-species surveys are also subject to errors in species identifications by imperfect classifiers. Imperfect classi-  
47 fiers include citizen scientists (e.g., North American Breeding Bird Survey (Sauer et al., 2017)), technicians trained in  
48 local taxonomy (e.g., invertebrate trapping by NEON (Hoekman et al., 2017)), automated methods (e.g., bat acoustic  
49 recording software (Wright et al., 2020) or convolutional neural networks used with camera trap data (Tabak et al.,  
50 2019)). Previous methods assume an imperfect classifier produces species-level classifications, but in practice, partic-  
51 ularly with human observers, we may end up with extraspecific classifications including "unknown", morphospecies  
52 designations, and taxonomic identifications coarser than species.

53 If species are prone to misclassification, then samples with known species identities might be used to correct  
54 estimates of occupancy parameters. However, using these data presents a methodological challenge. We refer to  
55 this situation as "semi-supervised": true species identities are known for some but not all individuals. Previous  
56 multi-species occupancy models that accommodate misclassification have used multinomial models that sum over all  
57 individuals (Wright et al., 2020), or site-level validation data where the occupancy state of a species is known only  
58 at a site- or plot-level but not at an individual-level (Chambert et al., 2018). Using individual-level validation data  
59 requires a different approach.

60 Misclassified species identities can be dealt with using one of two contrasting approaches. A simple two step  
61 approach 1) uses a classifier to assign species IDs to each individual (creating one complete synthetic dataset from  
62 classifier output, for which species identities are treated as known), then 2) analyzes the constructed dataset using  
63 a downstream model (e.g., an occupancy model). This two step approach does not propagate uncertainty in species  
64 identity to the downstream model, and the assignment of species identities in the first stage does not use any  
65 information about occupancy or encounter rates. In contrast, a joint model directly uses classifier output as data,  
66 relating the observation process to underlying ecological states in one step. Such an approach can simultaneously  
67 account for uncertainty in species identities, and leverage information about occupancy and encounter rates to inform  
68 species identity estimates (Wright et al., 2020). However, there remains the practical question of how much value

69 is added by a joint model vs. a two-stage approach. *A priori*, we expect that a joint model should produce better  
70 estimates of true species identities by using information on occupancy and encounter rates, but this has not yet been  
71 tested.

72 Here we develop an individual-level, semi-supervised, dynamic occupancy model that accounts for species non-  
73 detection and misclassification. Our Bayesian approach propagates classification uncertainty to offer an alternative to  
74 making conservation decisions based on point estimates. Our framework extends the classification occupancy model  
75 of Wright et al. (2020) to 1) accommodate extinction and colonization dynamics, 2) allow for additional uncertain  
76 "morphospecies" designations in the imperfect species classifications, and 3) make use of labeled samples with known  
77 species identities in a semi-supervised setting. Further, we compare the performance of a classification occupancy  
78 model to a reduced classification model that discards information about occupancy and encounter rate on a withheld  
79 test set. We demonstrate our model with an empirical case study of the carabid beetle (Carabidae) community at  
80 the National Ecological Observatory Network (NEON) Niwot Ridge Mountain Research Station (NIWO), west of  
81 Boulder, CO, USA, and quantify taxonomist identification error by accounting for classification probabilities.

## 82 2 Materials and Methods

### 83 2.1 Modeling occupancy dynamics with misclassification

84 Consider data collected at sites  $i = 1, \dots, N$ , according to a robust design (Hoekman et al., 2017) where each site is  
85 visited  $J$  times within primary periods  $t = 1, \dots, T$ , where the occupancy states are assumed to be constant within  
86 primary periods.

#### 87 2.1.1 State model

88 We are interested in occupancy states and encounter rates for species  $k = 1, \dots, K$ . Sites are either occupied ( $z_{i,k,t} = 1$ )  
89 or not ( $z_{i,k,t} = 0$ ). We assume that the occupancy states arise as Bernoulli random variables:

$$z_{i,k,t} \sim \text{Bernoulli}(\psi_{i,k,t}).$$

90 The probability of occupancy in the initial primary period is  $\psi_{i,k,1}$ . Subsequent occupancy dynamics depend on  
91 the probability of colonization  $\gamma_{i,k,t}$  and persistence  $\phi_{i,k,t}$ , such that for  $t > 1$ :

$$\psi_{i,k,t} = z_{i,k,t-1}\phi_{i,k,t-1} + (1 - z_{i,k,t-1})\gamma_{i,k,t-1}$$

#### 92 2.1.2 Encounter model

93 On any particular sampling occasion  $j$ , we encounter  $L_{i,j,k,t}$  individuals with encounter rate  $\lambda_{i,j,k,t}$ . We assume  
94 that the number of encounters is a Poisson random variable:  $L_{i,j,k,t} \sim \text{Poisson}(z_{i,k,t}\lambda_{i,j,k,t})$ . In a setting with  
95 misclassification, the number of encountered individuals  $L_{i,j,k,t}$  is not observed directly because of uncertainty in the

96 true species identities of encountered individuals. We do however observe the total number of individuals encountered  
 97 on any particular occasion:  $L_{i,j,..,t} = \sum_{k=1}^K L_{i,j,k,t}$ . The properties of sums of Poisson random variables allow us to  
 98 model these observed totals as:

$$L_{i,j,..,t} \sim \text{Poisson}\left(\sum_{k=1}^K z_{i,k,t} \lambda_{i,j,k,t}\right).$$

### 99 2.1.3 Observation model

100 In addition to observing the total number of encountered individuals on an occasion  $L_{i,j,..,t}$ , we assume that we  
 101 also obtain imperfect species classifications for each encountered individual. In cases where individuals have been  
 102 encountered ( $L_{i,j,..,t} > 0$ ), we obtain imperfect classifications of individuals  $l = 1, \dots, L_{i,j,..,t}$  and model these as arising  
 103 from a categorical distribution with a species-specific probability vector:

$$y_{i,j,l,t} \sim \text{Categorical}(\boldsymbol{\theta}_{k[i,j,l,t]}),$$

104 where  $y_{i,j,l,t}$  is the imperfect classification, and  $\boldsymbol{\theta}_{k[i,j,l,t]}$  is a probability vector associated with the true species of  
 105 individual  $l$ , which we denote  $k[i, j, l, t]$ . Element  $k'$  in the vector  $\boldsymbol{\theta}_{k[i,j,l,t]}$  represents the probability that an individual  
 106 is classified into category  $k'$ , conditional on the true species identity  $k[i, j, l, t]$ , such that  $\theta_{k[i,j,l,t],k'} = \Pr(y_{i,j,l,t} =$   
 107  $k' \mid k[i, j, l, t])$ . If species are always misclassified as other species, then  $\theta_k$  will be a vector of length  $K$  (Wright et al.,  
 108 2020). If there are extraspecific classes (e.g. morphospecies),  $\theta_k$  may have more than  $K$  elements.

109 True species identities are modeled as:

$$k[i, j, l, t] \sim \text{Categorical}\left(\frac{z_{i,k,t} \lambda_{i,j,k,t}}{\sum_k z_{i,k,t} \lambda_{i,j,k,t}}\right).$$

110 If ground truth species identity data are available for some individuals, then  $k[i, j, l, t]$  is partly observed and this  
 111 model can be trained in a semi-supervised setting. In the unsupervised setting, this individual-level formulation is a  
 112 disaggregated version of the single-season multinomial model of Wright et al. (2020) (Appendix S1).

### 113 2.1.4 Incorporating morphospecies designations

114 In some settings the imperfect classifier might assign more classes than there are unique species, so that the vector  
 115  $\boldsymbol{\theta}_k$  has more than  $K$  elements. For example, in the NEON beetle data, if a parataxonomist is unable to identify a set  
 116 of similar individuals, they will classify those individuals as a unique morphospecies associated with that sampling  
 117 occasion. Thus, it is possible for individuals to be classified into  $\tilde{K} \geq K$  classes, where  $\tilde{K}$  is sum of the number  
 118 of species and the total number of morphospecies designations. In such cases, the matrix  $\Theta = (\boldsymbol{\theta}'_1, \dots, \boldsymbol{\theta}'_K)$  can be  
 119 rectangular, with the first  $K$  columns corresponding to the classification probabilities for species 1, ...,  $K$ , and the  
 120 remaining columns corresponding to classification probabilities for non-species classes, e.g., morphospecies:

$$\Theta = \left[ \begin{array}{cc} \text{Species} & \text{Morphospecies} \\ \text{classifications} & \text{classifications} \\ \hline \theta_{1,1} & \dots & \theta_{1,K} & \dots & \theta_{1,K'} \\ \vdots & \ddots & & \ddots & \\ \theta_{K,1} & & \theta_{K,K} & & \theta_{K,K'} \end{array} \right] \left. \vphantom{\begin{array}{c} \theta_{1,1} \\ \vdots \\ \theta_{K,1} \end{array}} \right\} \text{True species identities}$$

## 121 2.2 Case study

### 122 2.2.1 Application to NEON carabid data

123 We fit our model to the carabid pitfall trap sampling data collected by NEON at NIWO in 2015-2018 (National Eco-  
124 logical Observatory Network, 2021). Carabids are a ubiquitous and speciose family of ground-dwelling invertebrates  
125 that are commonly collected by passive sampling methods, like pitfall traps, as described in (Hoekman et al., 2017).  
126 Carabids are a well-studied sentinel group that make an excellent study system for assessing community occupancy  
127 rates and classification accuracy. Collecting and identifying carabids is resource-intensive, but NEON lowers this  
128 barrier to entry by providing a public carabid dataset with three levels of classification (parataxonomist, expert  
129 taxonomist, then DNA barcoding). Although NEON processes carabid samples at the domain level (sites within the  
130 same ecoregion) (Hoekman et al., 2017), we focus our analysis on one NEON site, NIWO, to assess occupancy across  
131 co-occurring species. We use the 2015-2018 dataset at NIWO since carabid sampling started in 2015 and expert  
132 classification data were not yet fully available for 2019 at the time of analysis in 2020 due to data latency (National  
133 Ecological Observatory Network, 2021). NIWO is a site in the southern Rocky Mountains, spanning subalpine conifer  
134 forest and alpine tundra.

135 We outline the relevant data collection protocol here, but (Hoekman et al., 2017) offer more detail regarding  
136 NEON's carabid pitfall trap data product. The sampling design at every NEON site consists of ten permanent plots  
137 across the site with four pitfall traps per plot. Traps are sampled and reset biweekly during the growing season, with  
138 a range of 5-7 collections per year at NIWO. Our model runs at the plot-level. In 2018 one plot was permanently  
139 relocated to ensure sampling was allocated proportionally to the NLCD cover types represented (NEON help desk,  
140 personal communication).

141 All carabid samples are classified by a parataxonomist, and a subset are sent to an expert taxonomist for validation  
142 (Figure 1) (Hoekman et al., 2017). Species classification by parataxonomists is considered imperfect due to the brief  
143 taxonomic training of parataxonomists. Identification by an expert taxonomist is treated as confirmation data and  
144 is limited due to budget constraints. We confirmed the accuracy of the expert taxonomist classifications in finding  
145 that all individuals sent for DNA barcoding by NEON match the expert taxonomist's identification for the samples  
146 we used. In the few cases where the expert taxonomist could not identify a specimen to species-level, we use their  
147 genus-level classification for the validation dataset.

148 Our dataset contains 4772 individuals, 1764 of which were identified by an expert taxonomist, and 62 species  
149 classified by the parataxonomist, 23 of which are morphospecies. Morphospecies identifications are unique to each

150 NEON site and year. We fit our model using all individuals and used no environmental covariates. A hurdle for  
151 the NEON community in using the carabid pitfall trap data is reconciling parataxonomist and expert taxonomist  
152 classifications (Figure 1). Only one study to date has been published using the NEON carabid pitfall trap data (Egli  
153 et al., 2020), but Egli et al. (2020) analyze only the subset of individuals that have expert taxonomist classifications.  
154 Our model lowers the barrier of entry for using all imperfectly classified individuals while leveraging the available  
155 validation data.

## 156 2.2.2 Model specification

157 We used informative priors for the species classification probability vectors  $\theta_1, \dots, \theta_K$  that placed higher probability  
158 density on the correct species classification. In the case of NEON beetle data, this is reasonable given the training  
159 that parataxonomists receive in beetle identification. Because all elements of each  $\theta_k$  vector need to sum to one,  
160 and each element is bounded between 0 and 1, we used a Dirichlet prior:  $\theta_k \sim \text{Dirichlet}(\alpha_k)$ . We chose the  
161 Dirichlet concentration values  $\alpha_k$  by comparing draws from the Dirichlet prior distribution to our prior intuition  
162 about parataxonomist skill, using 200 along the diagonal (the element corresponding to the correct species identity),  
163 and 2 elsewhere.

164 We used multivariate normal priors at the species and site level, which allowed for correlations among param-  
165 eters. These priors share information among initial occupancy, persistence, colonization, and encounter rates. The  
166 motivation for this stemmed from a prior expectation that these parameters could be related. For example, species  
167 that are more abundant might be more likely to occur initially, persist, or colonize new sites. Similar arguments  
168 could be made about relationships among parameters at a site level.

169 Each species is associated with a vector  $\alpha_k$  of length 4, where  $\alpha_{k,1}$ ,  $\alpha_{k,2}$ ,  $\alpha_{k,3}$ , and  $\alpha_{k,4}$  are species-specific  
170 adjustments on initial occupancy, persistence, colonization, and encounter rates respectively. We assume that the  
171 species-specific adjustments are drawn from a multivariate normal prior with mean equal to zero, and an unknown  
172 covariance matrix:  $\alpha_k \sim \text{Normal}(\mathbf{0}, \Sigma^{(\alpha)})$ . Similarly, site-specific adjustments  $\epsilon_i$  were drawn from a different multi-  
173 variate normal prior. These adjustments were added together on a transformed scale to compute initial occupancy,  
174 persistence, colonization, and encounter rates, e.g.,  $\text{logit}(\psi_{i,k,1}) = \epsilon_{i,1} + \alpha_{k,1}$ . A full model specification for the case  
175 study is available in Appendix S2 (Plummer et al., 2003).

176 To evaluate how the occupancy and encounter rate components of the full model informed classification probability  
177 estimates, we also developed a reduced model that discards all information about occupancy and abundance, using  
178 just the expert and para-taxonomist species classifications to estimate the classification matrix  $\Theta$ . This comparison  
179 reveals the extent to which occupancy and encounter rates inform classification probabilities. If there are no differences  
180 in the estimates of classification probabilities, then a two-stage model which first models misclassification and then  
181 passes the posterior on as a prior for an occupancy/encounter model should perform as well as the joint model in  
182 which the classification model is integrated with the occupancy model. In addition to comparing posterior estimates  
183 for  $\Theta$ , we withhold a randomly selected 20% of the imperfect classifications to evaluate which model (full or reduced)

184 better predicts the data generated by the parataxonomist. All models were fit using JAGS (Appendix S3, dclone, and  
185 R v4.0.2 (Plummer et al., 2003; Sólymos, 2010; R Core Team, 2020) and visualized with ggplot (Wickham, 2016).

## 186 3 Results

### 187 3.1 Dynamic occupancy model

188 The occupancy model was designed to allow correlation between parameters across sites and species. Occupancy,  
189 growth, and turnover rates also varied through time. Sites with high encounter rates tended to have low initial  
190 occupancy and colonization probabilities and high persistence probabilities (Figure 2). Further, sites with high  
191 colonization rates tended to have high initial occupancy probabilities and low persistence probabilities. At the species  
192 level, we saw positive correlations among many of the model components, but in particular, species' encounter rate  
193 was positively correlated with species' initial occupancy, persistence, and colonization rates (Figure 2). Species varied  
194 in their detection success by the imperfect classifier, from ones that were common and consistently identified correctly  
195 (e.g. *Calathus advena*) to ones that were not identified at all (e.g., *Dicheirotichus mannerheimii*) but were caught  
196 by the expert taxonomist.

### 197 3.2 Classification model

198 The model yielded high probabilities of classification along the diagonal of the  $\theta$  confusion matrix where the expert  
199 and para-taxonomist identifications match (Figure 3). The model favors the parataxonomist's skill by giving more  
200 weight in the theta prior to diagonal values, making morphospecies classifications less probable. Individuals with  
201 morphospecies classifications make up a sizeable portion of the community, 811 out of the total 4772 total individuals  
202 identified by the parataxonomist. Despite the dirichlet priors favoring parataxonomist accuracy, some species had  
203 nontrivial probabilities of being classified as morphospecies than as the true species by the parataxonomist. For  
204 example, the parataxonomist was more likely to classify *Pterostichus (Hypherpes) sp.* as D13.2016.MorphBT than as  
205 the true species (Figure 3). However, no species had more than 3% probability (median) of being classified as another  
206 species (i.e., our model results indicate that the parataxonomist is most likely to identify a species either correctly  
207 or as a morphospecies).

208 To evaluate the value added by informing the classification model with occupancy and encounter rates, we com-  
209 pared the full model to a reduced classification model that discards all information about occupancy and abundance.  
210 Most  $\theta_k$  probability vectors do not differ between the full and reduced model results. However, we see differences for  
211 a few species where there are non-overlapping  $\theta$  posterior density distributions between the full and reduced models  
212 (i.e., Theta[46,46] and Theta [48,48], Figure 4). These differences are found most notably for the abundant species.  
213 The full model yielded higher classification probabilities for the abundant species. Further, the reduced model has  
214 wider 95% credible intervals compared to the full model for many theta indices (Figure 5). Thus, we find that a joint  
215 occupancy-classification model outperforms a two-stage model (classification, then occupancy).



216 We evaluated the performance of the full and reduced models by withholding a randomly selected 20% (352  
217 individuals) of the imperfect classifications that have an expert identification (1764 individuals). For the withheld  
218 individuals, the validation metric macro-averages are listed in Table 1. For every validation metric, the full model  
219 yields better results than the reduced model. The validation metrics calculated to the species level highlighted  
220 substantial differences between the models for common species. This confirms the result that occupancy dynamics  
221 improve classification model performance compared to using classification data alone.

## 222 4 Discussion

223 We developed a statistical approach to improve classification in multispecies datasets that leverages occupancy dy-  
224 namics. Our probabilistic framework can be applied to datasets with imperfect species detection and also errors  
225 in classification of samples. This approach builds on recent work on classification in occupancy models (Devarajan  
226 et al., 2020, and references therein) by evaluating the advantage of a joint occupancy-classification model, allowing  
227 imperfect classifications to outnumber species, and leveraging individual-level confirmation data in a semi-supervised  
228 setting. While analyses targeting species richness may be shielded to a certain extent from imperfect classification  
229 (Egli et al., 2020), any population- or community-level analyses with taxonomic specificity require an understanding  
230 of classification uncertainty in the data. Whereas imperfect classifiers offer classification point-estimates, our model  
231 provides a vector of probabilities for every species.

232 This is the first model to consider how occupancy and encounter rates contribute to improving classification.  
233 We found that our joint occupancy-classification model outperformed a reduced model that disregarded occupancy  
234 dynamics in estimating imperfect classification (Figure 6, Table 1). When looking at the validation metrics at the  
235 species level, the joint model surpasses the reduced model even more for abundant species. In Figure 5, we see  
236 that the full model posteriors have smaller CI widths than the reduced model for many species, which visualizes the  
237 superior precision of the full model.

238 False positive and negative species classifications are inevitable in any field collection, due to time and money  
239 constraints or imperfect classifiers (Royle and Link, 2006; Miller et al., 2012; McClintock et al., 2010; Hoekman  
240 et al., 2017). Accounting for false identifications is important to reduce bias in occupancy dynamics estimated from  
241 multispecies biodiversity monitoring datasets (McClintock et al., 2010; Miller et al., 2011; Chambert et al., 2015;  
242 Miller et al., 2015). Alternative models that account for false positives may consider data from only the focal species  
243 (Chambert et al., 2015) or from binary observations (Chambert et al., 2017). Like Wright et al. (2020), we use  
244 available counts from an imperfect classifier (Figure 1). However, we use all species detected, no matter how rare.  
245 By allowing taxonomic uncertainty propagation for multispecies datasets where imperfect classifications outnumber  
246 species (e.g., unknown, morphospecies, to the family level) by using a rectangular classification matrix  $\theta$  (Figure  
247 3), we remove an assumption that previous occupancy modeling methods have used (Chambert et al., 2018; Wright  
248 et al., 2020).

249 Our model is semi-supervised and makes use of data at the individual-level. Whereas alternative models use

250 data pooled at the site- or occasion-level (Chambert et al., 2018; Wright et al., 2020), our model leverages the rich  
251 information at the individual-level to reveal which species are commonly mistaken by the imperfect classifier and how  
252 often. Individuals identified by the expert taxonomist we know as true positives so can be used as partially-observed  
253 occupancy data in our semi-supervised model. In contrast, the model by Wright et al. (2020) is unsupervised, but  
254 also at the individual-level. Because our analysis is done at the individual-level, we can use species counts to inform  
255 classification (Chambert et al., 2017). Although the model priors favor parataxonomist accuracy, the model found  
256 high probability of classification for a couple of morphospecies that were abundant in the data (D13.2015.MorphO  
257 and D13.2016.MorphBT) (Figure 3).

258 Despite its innovations, our model has limitations. For NEON's carabid data from NIWO, which we used to fit our  
259 model, the parataxonomists were skilled and had high identification agreement with expert taxonomist classifications.  
260 The model may yield unexpected results when applied to a NEON site with lower parataxonomist accuracy. This  
261 raises the question of how much validation data is necessary to fit the model for varying degrees of imperfect classifier  
262 accuracy. This could be answered with simulations to identify what percentage or which type of samples should be  
263 prioritized for expert taxonomist classification to yield desired results.

264 We tried various iterations of the model before coming to the final disaggregated, semi-supervised, individual-  
265 level model. Using an aggregated data approach, we found the model either would not converge or struggled with  
266 identifiability, yielding multimodal posteriors for  $\theta$ . Changing the Dirichlet priors to favor parataxonomist accuracy  
267 helped but did not eliminate the problem. Future work could more explicitly incorporate false positives by informing  
268 the  $\theta$  priors with a list of species commonly misidentified by the imperfect classifier.

269 A model assumption is that samples were selected at random for verification. In reality, NEON prioritizes samples  
270 that were not classified to the species level by the parataxonomist. The model outlined here could be extended to  
271 represent processes such as this, by which samples are selected for verification by a more accurate classifier. Such  
272 extensions could take advantage of information contained in whether or not samples are selected for verification (e.g.,  
273 the fact that a sample was not chosen for verification is informative, as it may indicate higher confidence in the initial  
274 classification).

275 Large-scale, long-term biodiversity surveys are critical to inform land management and conservation policy  
276 (Hughes et al., 2017) and require affordable and efficient species classifications to stay on track and within budget.  
277 Egli et al. (2020) make a case for better training of parataxonomists to improve classification error rates. However,  
278 training people is expensive and staff can be transient, calling for a more systematic solution. This probabilistic  
279 approach can model species occupancy while accounting for imperfect classification, without additional training.  
280 Innovations in occupancy models in general, are rapidly being made to consider an expanding variety of study sys-  
281 tems and experimental designs (Bailey et al., 2014). Our results support the concept that ecological dynamics (i.e.  
282 occupancy and encounter rates) inform classification probabilities and lays a foundation for future work to build  
283 upon.

284

## 285 **Acknowledgements**

286 The National Ecological Observatory Network is a program sponsored by the National Science Foundation and  
287 operated under cooperative agreement by Battelle Memorial Institute. This material is based in part upon work  
288 supported by the National Science Foundation through the NEON Program. We thank G. Vagle for taking part in  
289 the conception of this project and J. Coulombe for her graphical design assistance. The work was supported by the  
290 CU Boulder Grand Challenge investment in Earth Lab. AIS was supported as a GRA at Earth Lab for work on this  
291 project. Any use of trade, product, or firm names is for descriptive purposes only and does not imply endorsement  
292 by the U.S. Government.

## 293 **Authors' contributions**

294 AIS, CLT, and MBJ conceived the project idea; AIS, JAR, and MBJ designed methodology; AIS curated the data;  
295 AIS and MBJ analysed the data; AIS and MBJ led the writing of the manuscript. All authors contributed critically  
296 to the drafts and gave final approval for publication.

## 297 **Data availability**

298 Carabid data accessible through the NEON data portal at [https://data.neonscience.org/data-products/DP1.](https://data.neonscience.org/data-products/DP1.10022.001)  
299 [10022.001](https://data.neonscience.org/data-products/DP1.10022.001) (National Ecological Observatory Network, 2021). Code for data cleaning and analysis is available at  
300 <https://github.com/annaspiera/NEON-NIW0-misclass>

## 301 **References**

- 302 Bailey, L. L., MacKenzie, D. I., and Nichols, J. D. (2014). Advances and applications of occupancy models. *Methods*  
303 *in Ecology and Evolution*, 5(12):1269–1279. <https://doi.org/10.1111/2041-210X.12100>.
- 304 Chambert, T., Grant, E. H. C., Miller, D. A., Nichols, J. D., Mulder, K. P., and Brand, A. B. (2018). Two-species  
305 occupancy modelling accounting for species misidentification and non-detection. *Methods in Ecology and Evolution*,  
306 9(6):1468–1477. <https://doi.org/10.1111/2041-210X.12985>.
- 307 Chambert, T., Miller, D. A., and Nichols, J. D. (2015). Modeling false positive detections in species occurrence data  
308 under different study designs. *Ecology*, 96(2):332–339. <https://doi.org/10.1890/14-1507.1>.
- 309 Chambert, T., Waddle, J. H., Miller, D. A., Walls, S. C., and Nichols, J. D. (2017). A new framework for analysing  
310 automated acoustic species detection data: Occupancy estimation and optimization of recordings post-processing.  
311 *Methods in Ecology and Evolution*, 9(3):560–570. <https://doi.org/10.1111/2041-210X.12910>.

- 312 Devarajan, K., Morelli, T. L., and Tenan, S. (2020). Multi-species occupancy models: review, roadmap, and recom-  
313 mendations. *Ecography*. <https://doi.org/10.1111/ecog.04957>.
- 314 Dorazio, R. M. and Royle, J. A. (2005). Estimating size and composition of biological communities by modeling the  
315 occurrence of species. *Journal of the American Statistical Association*, 100(470):389–398. <https://doi.org/10.1198/016214505000000015>.
- 316
- 317 Egli, L., LeVan, K. E., and Work, T. T. (2020). Taxonomic error rates affect interpretations of a national-scale  
318 ground beetle monitoring program at national ecological observatory network. *Ecosphere*, 11(4):e03035. <https://doi.org/10.1002/ecs2.3035>.
- 319
- 320 Hoekman, D., LeVan, K. E., Ball, G. E., Browne, R. A., Davidson, R. L., Erwin, T. L., Knisley, C. B., LaBonte, J. R.,  
321 Lundgren, J., Maddison, D. R., et al. (2017). Design for ground beetle abundance and diversity sampling within  
322 the national ecological observatory network. *Ecosphere*, 8(4):e01744. <https://doi.org/10.1002/ecs2.1744>.
- 323 Hughes, B. B., Beas-Luna, R., Barner, A. K., Brewitt, K., Brumbaugh, D. R., Cerny-Chipman, E. B., Close, S. L.,  
324 Coblenz, K. E., De Nesnera, K. L., Drobniitch, S. T., et al. (2017). Long-term studies contribute disproportionately  
325 to ecology and policy. *BioScience*, 67(3):271–281. <https://doi.org/10.1002/ecs2.1744>.
- 326 Iknayan, K. J., Tingley, M. W., Furnas, B. J., and Beissinger, S. R. (2014). Detecting diversity: emerging methods to  
327 estimate species diversity. *Trends in ecology & evolution*, 29(2):97–106. <https://doi.org/10.1016/j.tree.2013.10.012>.
- 328
- 329 MacKenzie, D. I., Nichols, J. D., Lachman, G. B., Droege, S., Andrew Royle, J., and Langtimm, C. A. (2002).  
330 Estimating site occupancy rates when detection probabilities are less than one. *Ecology*, 83(8):2248–2255. [https://doi.org/10.1890/0012-9658\(2002\)083\[2248:ESORWD\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2002)083[2248:ESORWD]2.0.CO;2).
- 331 McClintock, B. T., Bailey, L. L., Pollock, K. H., and Simons, T. R. (2010). Unmodeled observation error induces  
332 bias when inferring patterns and dynamics of species occurrence via aural detections. *Ecology*, 91(8):2446–2454.  
333 [https://doi.org/10.1890/0012-9658\(2002\)083\[2248:ESORWD\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2002)083[2248:ESORWD]2.0.CO;2).
- 334 Miller, D. A., Bailey, L. L., Grant, E. H. C., McClintock, B. T., Weir, L. A., and Simons, T. R. (2015). Performance  
335 of species occurrence estimators when basic assumptions are not met: a test using field data where true occupancy  
336 status is known. *Methods in Ecology and Evolution*, 6(5):557–565. <https://doi.org/10.1111/2041-210X.12342>.
- 337 Miller, D. A., Nichols, J. D., McClintock, B. T., Grant, E. H. C., Bailey, L. L., and Weir, L. A. (2011). Improving  
338 occupancy estimation when two types of observational error occur: Non-detection and species misidentification.  
339 *Ecology*, 92(7):1422–1428. <https://doi.org/10.1890/10.1890/10-1396.1>.
- 340 Miller, D. A., Weir, L. A., McClintock, B. T., Grant, E. H. C., Bailey, L. L., and Simons, T. R. (2012). Experimental  
341 investigation of false positive errors in auditory species occurrence surveys. *Ecological Applications*, 22(5):1665–  
342 1674. <https://doi.org/10.1890/11-2129.1>.

343 National Ecological Observatory Network (2021). Data products: Neon.dp1.10022.001. Provisional data downloaded  
344 from <https://data.neonscience.org/data-products/DP1.10022.001>.

345 Ovaskainen, O., Tikhonov, G., Norberg, A., Guillaume Blanchet, F., Duan, L., Dunson, D., Roslin, T., and Abrego,  
346 N. (2017). How to make more out of community data? a conceptual framework and its implementation as models  
347 and software. *Ecology letters*, 20(5):561–576. <https://doi.org/10.1111/ele.12757>.

348 Plummer, M. et al. (2003). Jags: A program for analysis of bayesian graphical models using gibbs sampling. In  
349 *Proceedings of the 3rd international workshop on distributed statistical computing*, volume 124, pages 1–10. Vienna,  
350 Austria.

351 R Core Team (2020). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical  
352 Computing, Vienna, Austria. <https://www.R-project.org/>.

353 Royle, J. A. and Link, W. A. (2006). Generalized site occupancy models allowing for false positive and false negative  
354 errors. *Ecology*, 87(4):835–841. [https://doi.org/10.1890/0012-9658\(2006\)87\[835:GSOMAF\]2.0.CO;2](https://doi.org/10.1890/0012-9658(2006)87[835:GSOMAF]2.0.CO;2)

355 Sauer, J. R., Niven, D. K., Hines, J. E., D. J. Ziolkowski, J., Pardieck, K. L., Fallon, J. E., and Link, W. A. (2017).  
356 The north american breeding bird survey, results and analysis 1966 - 2015. [https://www.mbr-pwrc.usgs.gov/  
357 bbs/bbs.html](https://www.mbr-pwrc.usgs.gov/bbs/bbs.html).

358 Sólymos, P. (2010). dclone: Data cloning in R. *The R Journal*, 2(2):29–37. <https://doi.org/10.32614/RJ-2010-011>,

359 Tabak, M. A., Norouzzadeh, M. S., Wolfson, D. W., Sweeney, S. J., VerCauteren, K. C., Snow, N. P., Halseth, J. M.,  
360 Di Salvo, P. A., Lewis, J. S., White, M. D., et al. (2019). Machine learning to classify animal species in camera  
361 trap images: Applications in ecology. *Methods in Ecology and Evolution*, 10(4):585–590. [https://doi.org/10.  
362 1111/2041-210X.13120](https://doi.org/10.1111/2041-210X.13120).

363 Talts, S., Betancourt, M., Simpson, D., Vehtari, A., and Gelman, A. (2018). Validating bayesian inference algorithms  
364 with simulation-based calibration. *arXiv preprint arXiv:1804.06788*. <https://arxiv.org/abs/1804.06788>.

365 Wickham, H. (2016). *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. [https://ggplot2.  
366 tidyverse.org](https://ggplot2.tidyverse.org).

367 Wright, W. J., Irvine, K. M., AlMBERG, E. S., and Litt, A. R. (2020). Modelling misclassification in multi-species  
368 acoustic data when estimating occupancy and relative activity. *Methods in Ecology and Evolution*, 11(1):71–81.  
369 <https://doi.org/10.1111/2041-210X.13315>

370 **5 Figures**

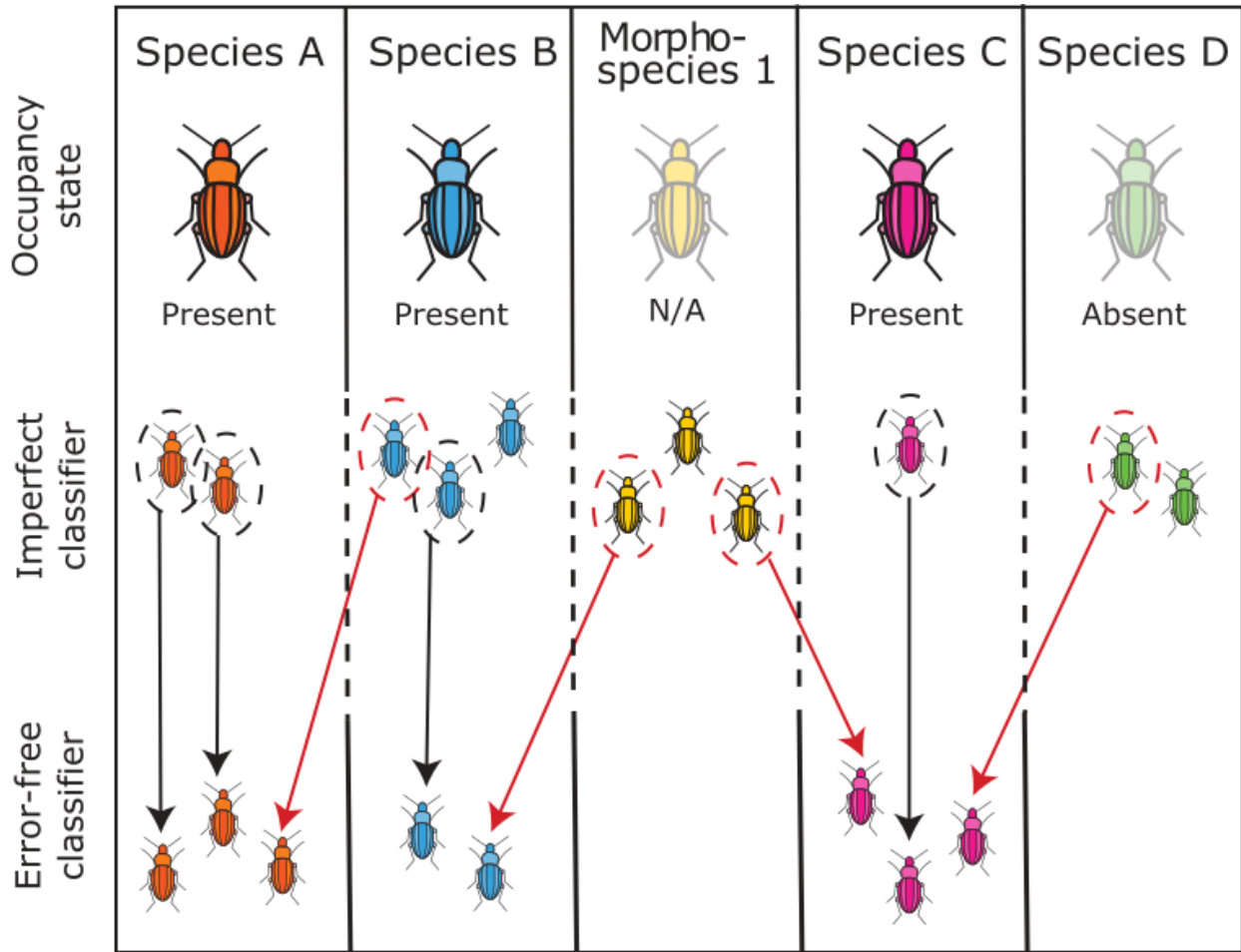


Figure 1: Classification scenarios in NEON carabid data. Each column is an imperfect classifier label. Each species is either present or absent, and morphospecies don't have an occupancy state. In some cases, the imperfect classifier (parataxonomist) matches the error-free classifier (expert taxonomist) (black arrow), in other cases the imperfect classifier was wrong (red arrow), while in other cases still, the error-free classification is unknown due to lack of validation data. For example, the *Morphospecies 1* individual with no error-free classification must belong to a different column, but this species identity is unknown in the raw data.

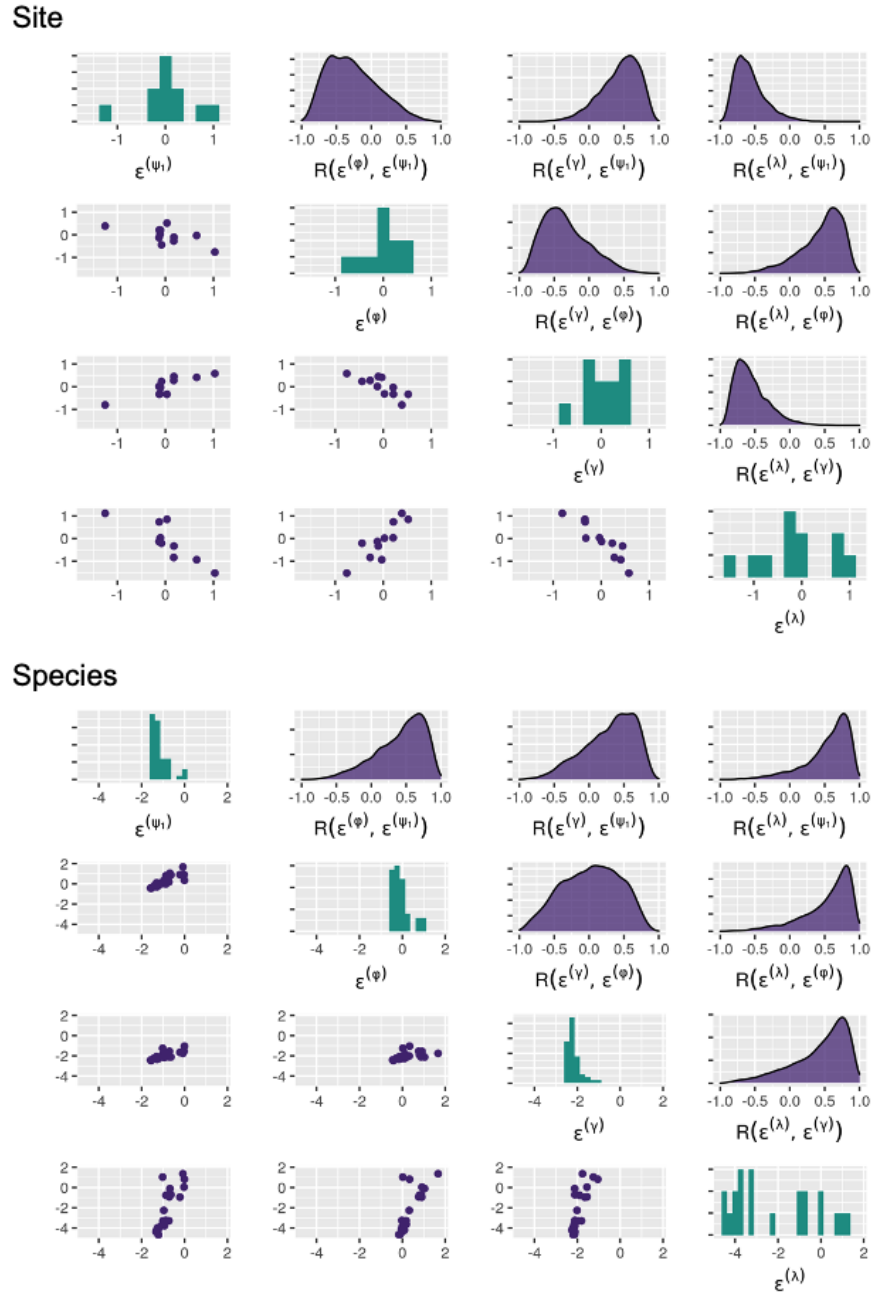


Figure 2: Random effects at the site and species levels. Rows correspond to the rows in the random effect covariance matrix: initial occupancy (row/col 1), persistence (row/col 2), colonization (row/col 3), and encounter rate (row/col 4). Along the diagonal are marginal histograms of posterior medians. Below the diagonal are pairwise scatterplots (each point is a site or species). Above the diagonal are posterior density plots of the pairwise correlation.

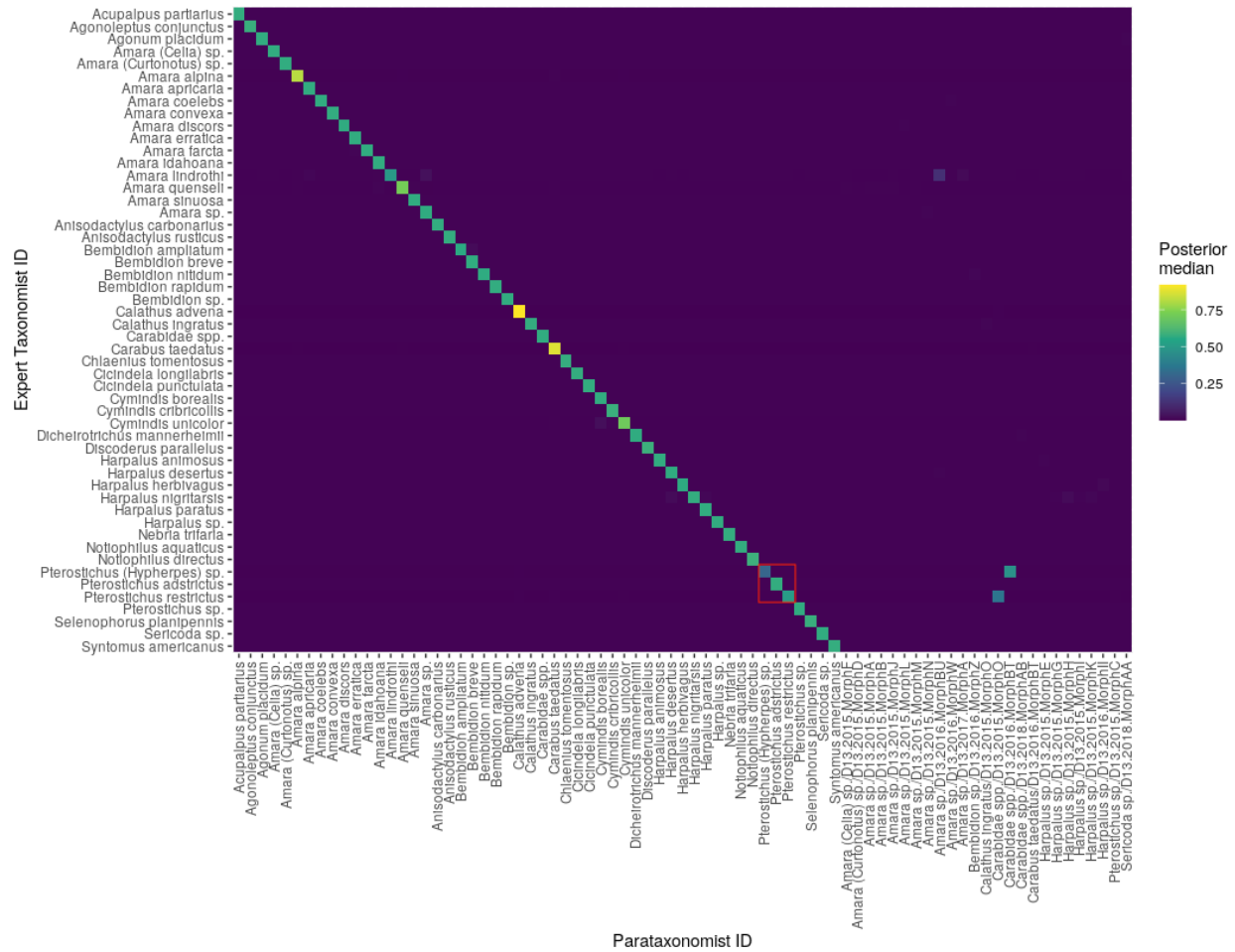


Figure 3:  $\Theta$  confusion matrix for the full model (classification with occupancy). Heat map of posterior median values where a value is interpreted as the probability that the species in that row is identified as a species in that column. Values along the diagonal are where the species is correctly identified. Values in each row sum to 1. The  $\theta$  posterior distributions for the 3x3 cells outlined in red are illustrated in Figure 4, along with the posteriors for the reduced (classification alone) and prior models.



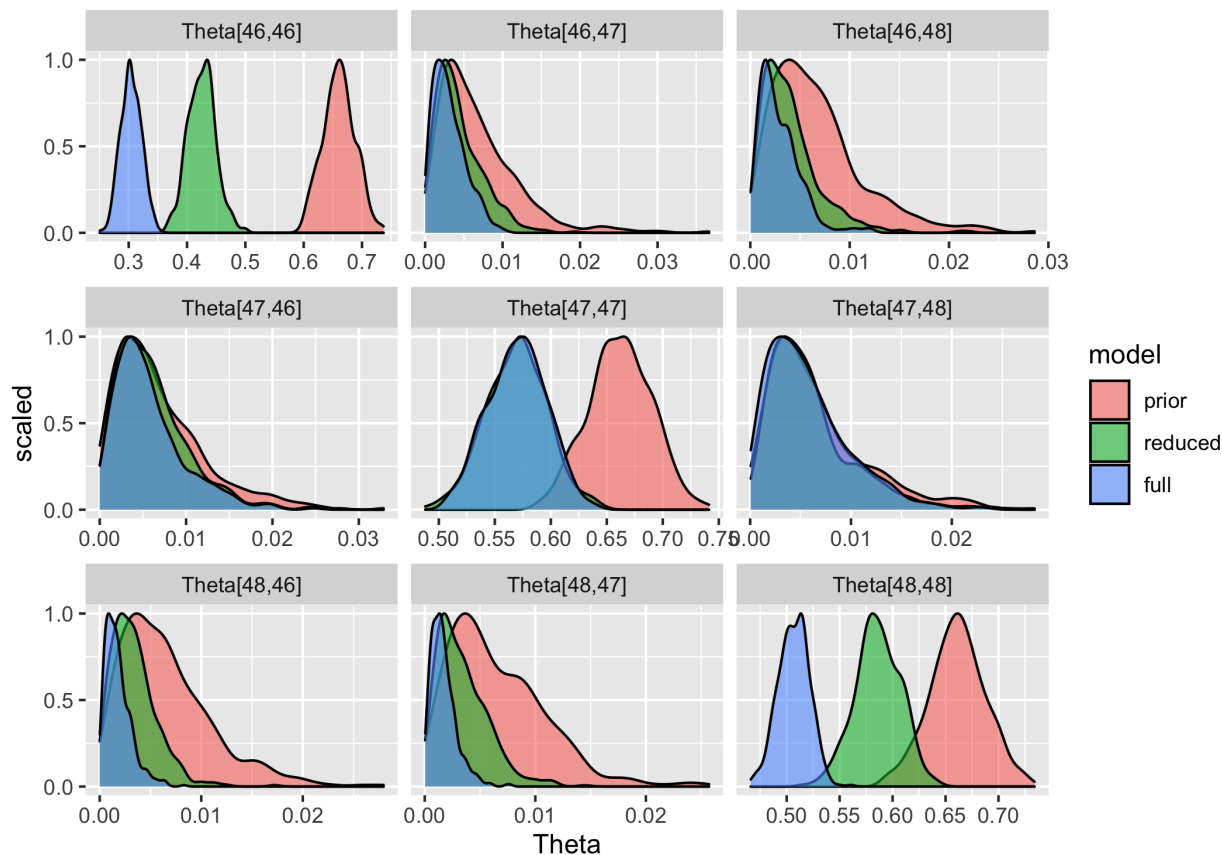


Figure 4: Comparison of  $\theta$  density distribution between prior, full model posterior (classification with occupancy), and reduced model posterior (classification alone) for select  $\theta$  confusion matrix indices. For non-abundant species,  $\theta_k$  probability distributions look like the second row. Here the three models' probability distributions overlap on the off-diagonal (Theta[47,46]/[47,48]) and the posterior distributions overlap but differ from the prior along the diagonal (Theta[47,47]). In contrast the top and bottom rows reflect probability distributions for abundant species. The posterior distributions overlap and are narrower and smaller than the prior on off-diagonal values (Theta[46, 47]/[46,48], Theta[48,46]/[48,47]). Along the diagonal, we see a difference in posterior probability distribution between the reduced and full models (Theta[46,46], Theta[48,48]), visualizing how the full and reduced models perform differently.

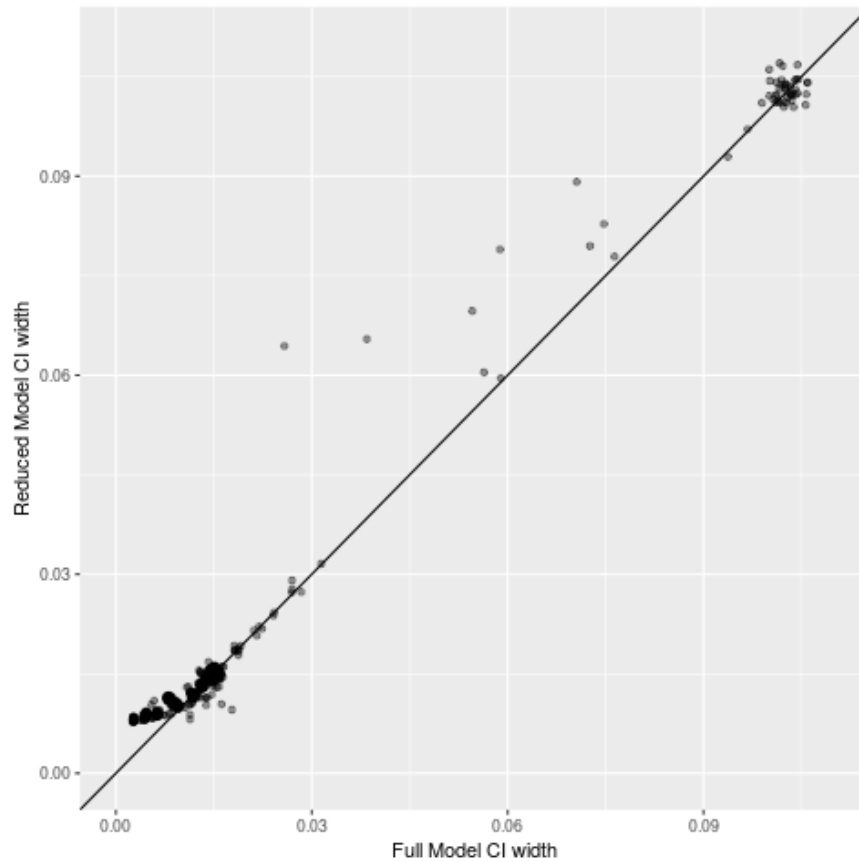


Figure 5: Comparison of precision between theta posterior 95% credible intervals (CI) of full model (occupancy with classification) and reduced model (classification alone). The full model is more precise for points above the line, indicating that the reduced model has a larger CI than the full model for that species.

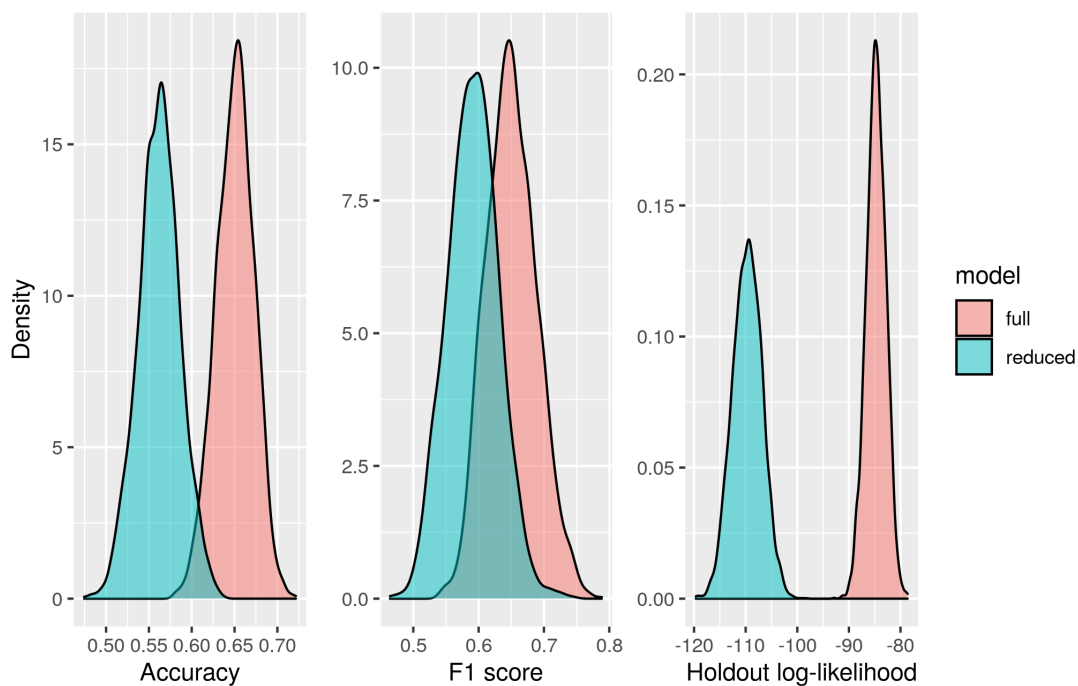


Figure 6: Validation metric macro-averages across posterior draws for accuracy, F1 score, and holdout log-likelihood.

|                        | Full model | Reduced model |
|------------------------|------------|---------------|
| Accuracy               | 0.65       | 0.56          |
| Precision              | 0.23       | 0.19          |
| Recall                 | 0.37       | 0.35          |
| F1 score               | 0.65       | 0.59          |
| Holdout log-likelihood | -84.7      | -110          |

Table 1: Validation metric macro-averages