

1 Journal name: Ecological Applications (ESA)

2 Manuscript type: Article

3 Manuscript title: Advancing single species abundance models: robust models for predicting abundance  
4 using co-occurrence from communities

5 Author names: Aliénor Stahl<sup>1</sup>, Eric J. Pedersen<sup>1</sup> and Pedro R. Peres-Neto<sup>1</sup>

6 Affiliations: 1: Concordia University, Department of Biology, Montreal, Canada

7 Corresponding author: Aliénor Stahl, [alienor.stahl@gmail.com](mailto:alienor.stahl@gmail.com)

8 Open research statement: Data and code to replicate the presented analyses are available in Zenodo at  
9 <https://doi.org/10.5281/zenodo.11150229>

10 Key words: abundance; co-distribution; co-occurrence; latent; prediction; simulation

11

## 12 **Abstract**

13 Accurate estimates of abundance are crucial for successful conservation and management.  
14 However, gathering abundance data is costly. Species Abundance Models (SAMs) are  
15 increasingly used to predict variation in abundance for resource management for single  
16 species, but collecting enough relevant environmental information to build effective SAMs  
17 can often be challenging. Species co-occurrence patterns may provide additional information  
18 on missing environmental predictors, and data on presence-absence species co-occurrence are  
19 typically easier to collect than abundance or detailed environmental data. However, it is still  
20 not clear when supplementing abiotic data with co-occurrence data should improve abundance  
21 predictions, as co-occurrence data itself represents a noisy indicator of the local environment.  
22 Using simulated data where we manipulated the strength of relevant environmental predictors  
23 across multiple species, we assessed the conditions that improve model predictions of a target  
24 species by using co-occurrence data on the remaining species as a proxy for missing  
25 environmental predictors. Because species often share environmental preferences in nature, an  
26 aspect simulated in our data, latent variables are expected to summarize important  
27 environmental gradients across co-occurring species. We employed Gaussian copulas to  
28 generate presence-absence co-occurrence-based latent variables as proxies. These latent  
29 variables, along with various combinations of environmental predictors, were subsequently  
30 used as predictors in SAMs. We evaluated the accuracy of these models in predicting the  
31 presence and abundance of target species through model validation exercises. Our results  
32 showed that incorporating presence-absence latent predictors generally improved model  
33 performance when compared to models lacking relevant environmental predictors, although  
34 there was considerable variation in performance across simulations. All models tended to have  
35 greater error rates when predicting abundant species compared to rare species. The goal of our

36 proposed framework is to offer a novel and easy to implement method for accurately  
37 predicting abundance from both biotic and environmental information.

38

## 39 **Introduction**

40 Community ecology has grown increasingly quantitative in response to the demand for a  
41 deeper understanding and more accurate predictions regarding how ecological factors and  
42 processes influence abundance, biomass, and interactions among both coexisting and non-  
43 coexisting species (Flecker and Matthews 1999; Persson 2008). Abundance serves as a critical  
44 indicator for individual species, their communities, and/or the state of the environment,  
45 enabling us to quantify ecosystem functioning (e.g., predation pressure, densities of preys  
46 available, the probability of reproductive encounters) (Degnbol and Jarre 2004). However,  
47 abundance data is generally difficult to collect across many different locations in  
48 heterogeneous landscapes (e.g., across many lakes in a landscape) whereas data on the  
49 presence or absence of communities of species can be easier to collect at landscape scales  
50 (Jackson and Harvey 1997). As such, it would be useful for landscape-scale management to  
51 be able to predict the local abundance of specific species based on easier-to-sample data such  
52 as the presence or absence of other species.

53 Many conventional models used to predict abundance rely on local (e.g., lake temperature)  
54 and regional (e.g., number of growing degree days) environmental variables (Lek et al. 1996;  
55 Brosse et al. 1999; VanDerWal et al. 2009; Boyce et al. 2016; Bradley 2016; Sobrino et al.  
56 2020). While environmental variables are relatively easy to gather through sampling or  
57 existing datasets, they are unlikely to encompass the multitude of sources of variation  
58 necessary for accurately predicting the abundances of target species of interest and other  
59 responses related to their communities, such as species composition. This limitation arises

60 because it is not often possible to measure all relevant environmental variables, and many  
61 species and community responses depend on factors beyond just environmental ones.

62 Additional factors, such as species interactions and history of introducing exotic species,  
63 among many others, also play important roles in shaping species patterns of species  
64 distributions, including abundance, and biodiversity (richness and species composition) in  
65 local communities and regionally (i.e., large scale variation).

66 In many cases, however, the environmental data gathered and used for predicting abundance  
67 variation in space (e.g., across sites) may stand as the primary source of low predictive  
68 accuracy, rather than other additional factors. For instance, relevant environmental variables  
69 may be missing or subject to measurement errors, or there could be time lags in  
70 environmental fluctuations and related changes in abundances (Myers 1998; Dornelas et al.  
71 2013; Bengtsson, Baillie, and Lawton 1997); and these lags may vary spatially and temporally  
72 (i.e., non-stationarity in lag-responses) even for the same species. If an unmeasured driver  
73 affects the abundance of at least two species, whether positively, negatively, or even in  
74 opposite directions between the species, one can expect that information about the distribution  
75 of one of these two species would improve the prediction of the other. This is especially  
76 expected when the probability of a species' presence or absence is related to its abundances,  
77 and when the presence or absence of other species act as proxies for unmeasured quantitative  
78 factors (e.g., low versus high values), or qualitative factors (e.g., presence or absence of the  
79 missing factor). Indeed, several studies have shown that, for certain species, the most accurate  
80 predictor of abundance was information regarding the presences and absences of other species  
81 (González-Salazar, Stephens, and Marquet 2013; Lewis et al. 2017; Öglü et al. 2019; Olkeba  
82 et al. 2020). While pairwise comparisons can be somewhat effective when studying single  
83 species, the interactions among multiple species can be complex and may not be adequately  
84 captured by pairwise comparisons alone.

85 It is generally not feasible to include the presence of all species in a regional species pool as  
86 predictors in a model targeting even the abundance of a single species. This is because even a  
87 moderately sized regional species pool may result in tens or hundreds of additional predictors  
88 in any abundance model. As such, incorporating the presence of other species into abundance  
89 models requires some form of dimension reduction of the species pool prior to analysis. In  
90 addition, many dimension reduction methods can borrow information across species and  
91 characterize their patterns of co-occurrence in a much-reduced number of axes, thereby  
92 improving predictive power based on these axes rather than considering all species separately  
93 (Carreira-Perpinán 1997; Cunningham 2008).

94 A solution to incorporating complex co-occurrence data while retaining a low dimensionality  
95 is to employ latent variable models (Walker and Jackson 2011). Latent variables are  
96 unobservable variables or factors that are not directly measured but rather estimated based on  
97 the associations (covariation) among species. These latent variables aim to estimate the joint  
98 model probability distribution of species presences-absences and represent the underlying  
99 structure or patterns in the data by specifying how data points (e.g., species composition  
100 across local communities or sites) are likely to be generated. Several methods exist to estimate  
101 latent variables from abundance or presence-absence data, including non-model-based (e.g.,  
102 classic ordination methods such as principal component analysis) and model-based (e.g.,  
103 mixed-model ordinations) methods (Walker and Jackson 2011; Popovic et al. 2019; Popovic,  
104 Hui, and Warton 2022). The power of latent variable methods stems from their ability to  
105 capture hidden variation in a dataset in low dimensionality (ter Braak and Prentice 1988; ter  
106 Braak 1985). Our contribution here is to demonstrate the robustness of modeling the  
107 abundances of single target species as function of latent variables that model the co-  
108 occurrence (presence-absence patterns) of the other species. This aspect is particularly  
109 important for the management and conservation programs tailored to specific species. We

110 introduce this general modeling framework and evaluate its ability to represent sources of  
111 predictive error caused by unmeasured drivers through detailed simulations.

112 The goal of this study is to assess the robustness of our proposed framework for advancing  
113 single species abundance distribution models using species co-occurrence data of other  
114 species in their communities. We used detailed simulations to contrast the performance of  
115 models containing various levels of information on the environment and community  
116 composition. Moreover, because we generate abundance distributions for all species in our  
117 simulations, we can contrast our model performance between abundance-based and species-  
118 co-occurrence based. Specifically, using comprehensive simulations, we set out to assess the  
119 performance of our proposed species-abundance framework by: (1) deriving rules for  
120 determining the number of latent variables used in modeling single species abundances, (2)  
121 contrasting model performance containing varying levels of information about the true  
122 underlying drivers (environment) versus latents (i.e., environmental proxies based on co-  
123 occurrence patterns of species sharing variable levels of environmental affinities; Figure 1),  
124 and (3) assessing how predictive performance varies as a function of sample size (i.e., number  
125 of sites or local communities used as input into the model). In this study, we focused on  
126 scenarios in which species and their communities are influenced solely by environmental  
127 variation, without considering the impact of species interactions or dispersal, which can either  
128 enhance or diminish model performance (i.e., increase or decrease predictive accuracy,  
129 respectively).

130

## 131 **Material and method**

132 The simulations to test our framework followed the subsequent steps (see Figure 1 for an  
133 illustration of how this general workflow for a single simulated landscape):

- 134 1. Use stochastic simulations to generate landscape-scale environmental variation for  
 135 each site in a landscape, and to generate coefficients for each species determining how  
 136 average species abundance should vary as a function of environmental variables.
- 137 2. Simulate the abundance of species in each site, based on the environmental variables  
 138 and coefficients generated in step 1.
- 139 3. Calculate latent variables from the presence-absence data of the previously generated  
 140 abundance using Gaussian Copulas.
- 141 4. Using a subset of the data generated, train a set of statistical models for each species to  
 142 predict local abundance. Trained models varied in the number of included  
 143 environmental variables and whether the model included latent variables.
- 144 5. Use a suite of metrics to evaluate the ability of each model to predict patterns of  
 145 presence-absence and abundance for the sites that were not used to estimate the  
 146 models.

147 **Steps 1 and 2: simulating communities**

148 We used a Poisson model to simulate species abundances across different landscapes  
 149 representing communities spread across  $E$  environmental gradients, assuming that the values  
 150 of the environmental gradients were uncorrelated from one another, and that the log of the  
 151 mean abundance of each species was equal to the sum of linearly dependent functions of each  
 152 of the environmental gradients plus a species-specific intercept:

$$A_{s,j,u} \sim \text{Poisson}(\mu_{s,j,u}) \quad 1(a)$$

$$\mu_{s,j,u} = \exp(b_{0,s,u} + b_{1,s,u}X_{1,j,u} + b_{2,s,u}X_{2,j,u} + \dots + b_{E,s,u}X_{E,j,u}) \quad 1(b)$$

153 Here  $\mu_{s,j,u}$  is the expected number of individuals (abundance) of a species at a site,  
 154 conditional on the environmental covariates included in the model. The abundance values

155 were drawn from a Poisson distribution with mean  $\mu_{s,j,u}$ .  $s$  denotes species,  $j$  sites, and  $u$  the  
156 landscape.  $A_{s,j,u}$  is the abundance of the  $s^{\text{th}}$  species in site  $j$  of landscape  $u$ ,  $X_{1,j,u}$  to  $X_{E,j,u}$  are  
157 the  $E$  environmental covariates that vary for each site  $j$  of each landscape  $u$ ,  $b_{0,s,u}$  the  
158 intercept that vary for each species  $s$  and landscape  $u$ , and  $b_{1,s,u}$  to  $b_{E,s,u}$  fixed coefficients  
159 relative to environmental variables  $1$  to  $E$  for species  $s$  in landscape  $u$ .

160 We simulated environmental covariates by drawing  $J$  independent, normally distributed  
161 values for each of the  $E$  environmental variables for each landscape (step 1). Thus, values for  
162 each covariate were statistically independent, with each environmental covariate having a  
163 mean of 0 and a variance of 1 across sites. These environmental covariates can be interpreted  
164 as environmental gradients given that they were generated independently. The coefficients  
165 ( $b_{0,s,u}, b_{1,s,u}, \dots, b_{E,s,u}$ ) for each species were drawn from a uniform distribution with a range of  
166 -2.4 to 1.2 for the intercept, and -0.8 to 0.8 for the slopes. The ranges for the coefficients were  
167 determined through simulation trials where we identified the minimum and maximum  
168 coefficients that allowed for all species to be present in at least 10% of sites and at most in  
169 90% of sites. The selected parameters allowed to generate species with different levels of  
170 strength between abundance and environment variables (e.g., narrow versus broad niche  
171 breadths; step 2). Table 1 summarizes how each variable in eq. 1 was generated. The  
172 distribution across species of spatially averaged species abundance within each landscape was  
173 approximately log-normally distributed (Figure 2), resembling common patterns found in  
174 natural communities.

### 175 **Step 3: Latent variables generation and their abilities to represent missing** 176 **environmental variation**

177 Different methods are available for incorporating presence-absence information into a latent  
178 model (Popovic et al. 2019; Zou and Zhang 2009; Blanchet, Cazelles, and Gravel 2020). The



179 copula approach used here is a model-based latent approach to estimate latent variables from  
180 multivariate data sets, as implemented in the ecoCopula R package (Popovic et al. 2019). This  
181 Gaussian Copula graphical model approach combines a multivariate distribution (e.g.,  
182 multivariate Gaussian) with a set of marginal distributions (e.g., binomial, Poisson). Due to its  
183 high versatility (i.e., allowing for the selection of the multivariate distribution as well as the  
184 modeling of the appropriate discrete marginal distributions), it holds significant potential for  
185 applications in ecology (Anderson et al. 2019). Additionally, it has been shown to be one of  
186 the most accurate latent estimation methods in heterogenous environments (i.e., varying with  
187 a binary environmental covariate) (Popovic et al. 2019) and has been identified as the fastest  
188 and most robust latent variable quantification method for count and binomial (presence-  
189 absence) data (Popovic et al. 2022).

190 However, the copula model requires specifying the number of latent variables to estimate  
191 prior to model fitting. In general, at least  $E$  latent variables should be required to capture the  
192 variation in  $E$  independent environmental gradients, but it may be the case that more latent  
193 variables are needed to fully capture environmental variation. One frequently used method for  
194 determining the number of latent variables to retain is to compare AIC (Akaike Information  
195 Criterion) or BIC (Bayesian Information Criterion) for models with increasing numbers of  
196 latent variables until the chosen matrix reaches a minimum value (i.e., best predictive value of  
197 co-occurrence). However, initial testing on landscapes (simulated using the method in step 1)  
198 with varying numbers of latent variables consistently showed that, using the BIC method  
199 calculated in ecoCopula, the BIC score was always lowest for models with a single latent  
200 variable, regardless of the number of environmental predictors used to simulate species  
201 abundances. As such, we conducted a preliminary trial to evaluate the number of latent  
202 variables needed to best approximate the environmental gradients in our simulated  
203 landscapes.

204 Using eq. 1, we simulated  $U$  landscapes of size  $J$  (number of sites), containing  $S$  species and a  
205 varying  $E$  number of environmental predictors ( $U = 450, J \in \{100, 200, 300\}, S \in \{10, 20,$   
206  $30\}, E \in [1,5]$ ; Table 1). To evaluate the optimal number of latent parameters (axes) needed to  
207 best approximate the environmental gradients in our simulated landscapes and compare the  
208 impact of adding or removing latent variables, we generated several numbers of latent  
209 variables for each possible combination of parameter values. Therefore, for each possible  
210 combination of parameter values, we fitted the presence-absence data into a stacked species  
211 regression model before using a model-based ordination with Gaussian copulas by using the  
212 functions *stackedsdm* and *cord* from the package *ecoCopula* (Popovic et al. 2019, version 1.0-  
213 2) with  $L$  different numbers of latent factors to model them ( $L \in [1,5]$ ).

214 We extracted the BIC value of each of these models and subtracted from them the BIC of the  
215 best model from any given simulation set (i.e., lowest BIC for the species considered in the  
216 current landscape). To evaluate the effectiveness of the latent variables in representing (i.e.,  
217 serve as a proxy) environmental variation, we conducted a redundancy analysis (RDA) of the  
218 original environmental variables used to simulate species abundance regressed against the  
219 extracted latents using the function *rda* from the package *vegan* (Oksanen et al. 2022, version  
220 2.6-2). Ability of latents to represent environmental variation was measured via the RDA  
221 adjusted  $R^2$  (Peres-Neto et al. 2006). We determined from this trial that, regardless of the  
222 number of sites  $J$  or species  $S$  in the simulation, BIC was always lowest with a single latent  
223 variable (Appendix S1: Figure S1), but adjusted  $R^2$  did increase with the number of latent  
224 predictors, until the number of latents equalled  $E$ , after which the adjusted  $R^2$  did not increase  
225 with more latent variables (Appendix S1: Figure S2), so there is no reason to extract more  
226 than  $E$  latent variables for any given simulation.

#### 227 **Step 4: Contrasting the performance of abundance models**

228 We compared the models containing only the environmental variables used to generate  
229 species abundances (eq. 1) against the ones containing selected environmental variables and  
230 the latent variables (community composition). This allowed us to compare model  
231 performance under ideal conditions because we used the true environmental drivers used to  
232 simulate species abundances against models from which we removed various combinations of  
233 environmental variables (scenarios) and replaced them with latent variables (proxies) to  
234 represent the missing sources of variation. Note, however, that ideal conditions do not imply  
235 perfect model performance, as different species were simulated with varying degrees of  
236 strength and associated errors relative to environmental variables (e.g., narrow versus broad  
237 niche breadths).

238 For this contrast, we created  $U$  landscapes, and for each landscape  $u$ , we generated  $K$   
239 replicates ( $U = 30$ ,  $K = 10$  replicates per landscape). For each replicate  $k$ , we simulated  
240 abundances for each  $s$  species in each site  $j$  using eq. 1, using three environmental variables  
241  $\mathbf{X}_1$ ,  $\mathbf{X}_2$  and  $\mathbf{X}_3$  per landscape containing multiple sites. We simulated 20 species and 1000  
242 sites per landscape. We fixed the number of latent factors to 3 as we had three environmental  
243 variables (see RDA results in previous section). Replicates (i.e., landscapes using the same  
244 coefficients but had varying values of environmental gradients) were used to allow a  
245 reasonable estimate of the metrics used to contrast model performances.

246 We randomly sampled 100 sites (out of the 1000 simulated) from each landscape  $u$  (referred  
247 here as to the training set), and for each training set we estimated abundance models with  
248 different combinations of environmental and latent predictors (step 4). Each model was  
249 estimated using a Generalized Linear Model (GLM), using a Poisson distribution with a log-  
250 link function (Kéry and Royle 2015). We used the *manyglm* function from the R package  
251 *mvabund* (Wang et al. 2022, version 4.2-1) to fit separate models for each replicate landscape  
252 simultaneously for all species separately.

253 We were interested in comparing models containing different combinations of environmental  
254 variables and latent variables. The complete list of model scenarios considered is described in  
255 Table 2. As each species had different strengths of relationship with each environmental  
256 variable (i.e., different coefficient values in eq. 1 were used to simulate each species), we  
257 ordered the models based on the decreasing values of the environmental coefficients used to  
258 simulate the species' abundance. For instance, if species *A* had the values of -0.5, 0 and 0.8 as  
259 coefficients for the environmental variables  $X_1$ ,  $X_2$ , and  $X_3$ , respectively,  $X_3$  had the largest  
260 influence on driving abundance values, followed by  $X_2$  (i.e., importance is given by  
261 decreasing coefficient values) and  $X_1$ . But if species *B* had values of 0.7, -0.5 and 0.3 as  
262 coefficients for the environmental variables  $X_1$ ,  $X_2$ , and  $X_3$  respectively, its abundance was  
263 mostly driven by variations of  $X_1$ , then  $X_3$  and finally  $X_2$ . When removing  $X_1$  from the  
264 predictors of a model, species *A* and *B* were not impacted in the same way due to the lesser  
265 influence  $X_1$  had on the abundance of species *A*. We predicted that including latent variables  
266 should increase predictive ability more when added to a model that only included  
267 environmental predictors that weakly predicted the abundance of an individual species. To test  
268 this, we compared model performance with and without latent variables for models including  
269 different combinations of strengths of environmental variables.

270 For models containing one environmental variable as predictor, we labeled the predictors as  
271 “high”, “intermediate”, and “low”, corresponding to the decreasing values of coefficients of  
272 the environmental variables. For models incorporating two environmental variables, we  
273 designated the model with the two highest coefficients as “high”, the model with the highest  
274 and lowest coefficient as “intermediate”, and the model with the two lowest coefficients as  
275 “low”.

## 276 **Step 5: comparison of model performance**

277 For each model estimated for each replicate within the same landscape, we generated  
278 predictions for species abundances at the remaining 900 sites in the landscape from which the  
279 sites were sampled from (the test set). To establish baselines for optimal model performance,  
280 we also calculated predicted abundances in the test set using the oracle model: i.e., the model  
281 employing the true coefficients used to simulate each species' abundances to predict the  
282 conditional expected abundance for each species in each site. The oracle model represents the  
283 best possible model for estimating the simulated abundances in each test set that could be  
284 derived using data from the training set. Two other models were singled out: (i) a benchmark  
285 model containing all three environmental variables, to identify in which scenarios having  
286 access to all environmental variables (drivers of the abundance) did not suffice to properly  
287 estimate the environmental coefficients (by comparing the performance of the benchmark  
288 model to that of the oracle model), and (ii) a latent model containing only the latent variables,  
289 to study how species co-occurrence patterns performed as predictors of their own. We  
290 assessed how effectively the different models, including the oracle model, predicted the  
291 pattern of presences and absences as well as the true abundances in the test set.

292 Although our primary focus was on predicting abundance, we evaluated the models for both  
293 presence-absence and abundance predictions. This approach was taken because, in many  
294 cases, the interest may lie in predicting presence or absence of a particular target species. It is  
295 important to note, however, that the latents used as predictors were always derived based on  
296 the presence-absence of other species.

### 297 *Metrics for evaluating presence-absence predictions*

298 The Poisson regression models estimated in step 4 can predict the probability of presence of  
299 each species in a given site, but to evaluate the effectiveness of the model for predicting  
300 presence, these probabilities need to be translated into concrete predictions for presence or  
301 absence (Lawson et al. 2014; Phillips and Elith 2013). If we only treated a model as

302 predicting a species present if the probability of presence was over 50%, models for rare  
303 species would only predict absences (and vice versa for common species), so using a fixed  
304 probability threshold would lead to all models of rare (common) species having the same  
305 predictive performance as a model that just predicts the species always being absent (present).

306 Therefore, instead of using a fixed probability threshold to convert the probabilities into  
307 presence-absence predictions, we used a prevalence-based approach. For each species, we set  
308 a threshold equal to the true occurrence (prevalence) rate of the species across a given  
309 landscape (e.g., Liu et al. 2005). We used this threshold to generate a predicted presence-  
310 absence matrix for each site and each species in each landscape for a given model. This was  
311 achieved by determining whether the expected abundance by the model for that site was  
312 greater (present) or lower (absent) than the threshold value. We then compared the  
313 performance of each model to the oracle model using a range of metrics, the equations for  
314 which are provided in Table 3. Using the predicted presence-absence matrices, we calculated  
315 the True Skill Statistic (TSS, Peirce 1884; Table 3) for each model, species and landscape  
316 replicate. The TSS, which ranges from -1 to +1, measures the difference between the  
317 sensitivity and specificity of the model. A score of +1 indicates a perfect agreement between  
318 the model's predictions and the true presence-absence, while a score of 0 or lower signifies  
319 performance no better than random (Allouche, Tsoar, and Kadmon 2006). We calculated the  
320 ratio of the TSS of the model over the TSS of the oracle and computed the mean for each  
321 model, species and landscape. Then, we grouped species into bins based on occurrence rates  
322 across different landscapes. A TSS ratio of  $\geq 1$  indicates that the model performed as well or  
323 better than the oracle, while a TSS ratio of  $\leq 0$  or less means that the model predicted presence  
324 as badly or worse than random chance.

325 To compare whether including latent predictors increased model performance relative to just  
326 using environmental variables, we also calculated the delta TSS, defined as the TSS of

327 environmental model minus the TSS of corresponding latent model (i.e., models containing  
328 the same environmental variables where the only difference in specification was the inclusion  
329 of latent variables as predictors). A positive delta TSS indicates the environmental model to  
330 have the best performance, whereas a negative value suggests that the model including of  
331 latent variables performs best.

### 332 *Metrics for evaluating abundance predictions*

333 When evaluating how each model predicted species abundance, we limited comparisons to  
334 sites where the species was present (i.e., abundance of 1 or higher). To evaluate how well each  
335 model predicted species abundance we calculated the following prediction metrics for each  
336 model, species and landscape replicate: Mean Absolute Percentage Error (MAPE), Root Mean  
337 Squared Percentage Error (RMSPE), Relative Mean Squared Error (RMSE), Symmetric Mean  
338 Absolute Percentage Error (SMAPE), and Root Mean Ratio Percentage Error (RMRPE) (see  
339 Table 3 for definitions of these metrics). We calculated the ratio of each metric to the  
340 corresponding metric calculated for the oracle model (i.e., best possible scenario) and  
341 calculated the average ratio for each model, species and landscape (referred to as the ratio  
342 metric in the results). We also calculated the delta metric, defined as the metric calculated for  
343 a model containing only environmental variables minus the metric calculated for a model with  
344 the same environmental variables as well as latent variables. As above, a negative delta metric  
345 indicated that the latent model performed better than the same model lacking latent variables.

346 To illustrate how different metric performances varied with species abundance across  
347 simulations, we grouped species in different landscapes into percentile bins, based on the  
348 average (true) abundance of the species in its own landscape, and then calculated average  
349 ratio metrics and delta metrics for each percentile bin across landscapes and replicates.

350

## 351 **Results**

### 352 **Number of latent variables needed to capture environmental variation**

353 We first focus on determining the optimal number of latent dimensions to select when using  
354 Gaussian copulas. To assess the goodness of fit of the models, we examined both the RDA  
355 adjusted  $R^2$ , which represents the proportion of variance explained by the model, and the  
356 Bayesian Information Criterion (BIC), which is typically used to determine the optimal  
357 number of latent variables to retain. The RDA enabled us to estimate how effectively the  
358 latents characterize the original environmental variables (gradients) based on community  
359 composition, while the BIC helped us determine whether this criterion indeed allows for  
360 selection of an appropriate number of latents to represent community composition.

361 The adjusted  $R^2$  consistently increased with the number of latent dimensions until it equaled  
362 the actual number of environmental variables used to simulate the data, at which point it  
363 plateaued (Figure 3, Appendix S1: Figure S2). This indicates that additional latent variables  
364 did not improve the model's ability to predict the environmental state of a given location. The  
365 maximum fraction of variance explained was not significantly affected by the number of true  
366 environmental variables used to generate (simulate) species abundances; capturing variation  
367 from one environmental gradient was as feasible as capturing it from three or four  
368 environmental gradients (i.e., variables). Note, again, that the interpretation here as gradients  
369 is possible because environmental variables were generated independently. The adjusted  $R^2$   
370 was not sensitive to the number of sites in the landscape used to estimate the latent variables,  
371 but it was sensitive to the number of species used: models based on 10 species could only  
372 explain about 30% of the variation in environmental variables, regardless of the number of  
373 latent variables used, whereas models based on 30 species could explain ~60% of variation in  
374 the environmental matrix (Appendix S1: Figure S2).



375 In contrast, the Bayesian Information Criterion (BIC) consistently increased with the number  
376 of latent dimensions, without showing any signs of reaching a plateau (Appendix S1: Figure  
377 S1). While models with lower BIC are generally expected to have better predictive ability for  
378 unobserved data - suggesting that the best model would always retain one latent variable  
379 regardless of the environmental dimension - this expectation did not align with our  
380 observations for the adjusted  $R^2$ . This discrepancy indicates that BIC (as calculated by  
381 ecoCopula) is not a good metric of the predictive performance of the latent model, at least  
382 when applied to gradients driving abundances while their latents were extracted from  
383 presence-absence data. Therefore, we did not report BIC of the estimated latent models for the  
384 remainder of our simulations.

## 385 **Models' performance**

### 386 *Presence-absence predictions*

387 We now focus on the models' performance in predicting presence-absence, including the ratio  
388 TSS (representing how well each model performed compared to the oracle model) and delta  
389 TSS (represented how well models without latent variables performed relative to models  
390 including latent variables). The ratio of the TSS had a mean of 0.7 and ranged from -1.6 to 1.7  
391 (recall that any value below 0 indicates that the model did not perform better than random,  
392 while any value above 1 represents better performance compared to the oracle). Initially  
393 examining the TSS across species occurrence percentiles, there were no obvious patterns  
394 (Figure 4). In this case, the number of occurrences of a target species did not influence  
395 model's performance. When comparing models, models containing two environmental  
396 variables performed better on average than those with only one, regardless of whether latents  
397 are included or not.

398 When comparing models with and without latent variables, any delta TSS value above 0  
399 indicates that the environmental model performs better, while any negative value indicates a  
400 better performance by the latent model. Models containing latent variables generally  
401 performed better on average across all (target) species, especially for those with high  
402 occurrence and in models containing only one environmental predictor (Figure 4). The  
403 differences are less pronounced when comparing models that contain two environmental  
404 variables (i.e., where only one environmental predictor is missing from the model). Reducing  
405 the number of sites used to fit the model did not affect the performance of the TSS, sensitivity,  
406 or specificity (Appendix S1: Figure S3).

407 When comparing the TSS as performance of the oracle (i.e., a model using the true  
408 coefficients of the environmental variables to generate the species' conditional expectations),  
409 benchmark (i.e., a model containing all three environmental variables), and latent models (i.e.,  
410 a model containing only the latent variables), we can notice that they are very correlated  
411 across species occurrence percentiles (Figure 5). The benchmark and oracle models have  
412 extremely similar performances. Regarding sensitivity, the benchmark and oracle models are  
413 also highly correlated, while the latent model demonstrates good correlation for species with  
414 low occurrence. For specificity, the benchmark and oracle models are correlated for high  
415 occurrence species, while the benchmark and latent models are correlated for low occurrence  
416 species.

#### 417 *Abundance predictions*

418 To assess the goodness of fit for abundance-based models (i.e., target species include  
419 abundance information while latents are based on presence-absence of the other species), we  
420 calculated six metrics to assess the extent to which the models mispredict species abundances.  
421 Again, we used the ratio of each metric over the same metric calculated for the oracle model  
422 (i.e., representing the best possible predictive scenario), along with the delta metric to

423 compare models that differ in composition due to the inclusion or exclusion of latent  
424 variables.

425 To assess across all species the impact on model performance of removing any given  
426 environmental predictor, we had to consider the varying strengths in the relationship between  
427 each species abundance and each environmental variable to compare the predictive ability of  
428 latents. As a reminder, in models containing one environmental variable as predictor, we  
429 labeled the predictors as “high”, “intermediate”, and “low”, corresponding to the decreasing  
430 coefficients of the environmental variables. For models incorporating two environmental  
431 variables, we designated the model with the two highest coefficients as “high”, the model  
432 with the highest and lowest coefficient as “intermediate”, and the model with the two lowest  
433 coefficients as “low”. Regardless of the metric considered, we observe the following patterns:  
434 prediction error increases as species abundance increases, and models containing two  
435 environmental variables outperform models containing only one environmental variable  
436 (Figure 6, Appendix S1: Figure S4). When comparing models with or without latent variables,  
437 highly abundant species were best predicted by models containing latent variables (Figure 6,  
438 Appendix S1: Figure S4). For species with low and medium abundances, the inclusion or  
439 exclusion of latent did not impact the performance of the models; they exhibited very similar  
440 values of error.

441 When comparing the metrics in relation to the performance of the oracle (i.e., a model using  
442 the true coefficients of the environmental variables to generate the species’ conditional  
443 expectations), benchmark (i.e., a model containing all three environmental variables) and  
444 latent models (i.e., a model containing only the latent variables), we observe identical trends  
445 across all metrics. The performance of the three models was very similar for low abundance  
446 species; however, the latent model diverged when the abundance percentile was higher than  
447 70%, with an increase in predictive error (Appendix S1: Figure S5). The metrics were not

448 sensitive to the number of sites in the landscape used to fit the models (Appendix S1: Figure  
449 S6).

450

## 451 **Discussion**

### 452 **Number of latent variables needed to capture environmental variation**

453 Our first goal was to establish rules for determining the number of latent variables used in  
454 modeling single species abundances. To achieve this, we examined the behavior of two  
455 metrics, the BIC and the adjusted  $R^2$ , within a simulated landscape. Our results indicate that  
456 the BIC was not a useful metric for deciding the appropriate number of latent variables when  
457 employing Gaussians copulas. Instead of plateauing once the latent variables captured as  
458 much of the environment as possible, it continued to increase, implying that the best number  
459 of latent variables was consistently one even in cases where multiple independent  
460 environmental gradients were set to simulate species distributions. It is plausible that current  
461 calculation method for BIC is incorrect or does not employ an appropriate penalty measure  
462 (number of parameters and sample size). Note that there is a general lack of consensus about  
463 the best criteria for assessing latent models (Weller, Bowen, and Faubert 2020). On one hand,  
464 the BIC is generally regarded as a reliable metric for latent models (Nylund, Asparouhov, and  
465 Muthén 2007); however, it is also criticized for being overly conservative (Mindrila 2023) as  
466 it was the case here. Note, however, that the underperformance of BIC to decide the number  
467 of latents to use in species abundance models may be due to the fact that, in our simulations,  
468 species' responses to environmental gradients were in the form of abundances, whereas latent  
469 predictors were extracted from presence-absence data. Consequently, the more liberal AIC  
470 might be a preferable option for the Gaussian copulas used in our study. Note that regardless  
471 of whether we use AIC or BIC to assess the number of latents to retain, this assessment is

472 intrinsic and solely based on the community data used to estimate the latent variables, which  
473 are then used as predictors in abundance distribution models of single species. As we will  
474 discuss, an extrinsic selection, in which latents that improve abundance predictive accuracy  
475 are chosen, may prove to be a better strategy when using latent models based on co-  
476 occurrence data to predict abundance of single (target) species.

477 Note that the goal of the RDA analysis, based on the  $R^2$  metric, was to assess whether the  
478 latent structures used here could serve as a good proxy for the true environmental variables  
479 used to simulate species distributions. Given that the adjusted  $R^2$  plateaued when the number  
480 of latent variables equalled the true number of environmental dimensions, it instills  
481 confidence that these latents serve as robust proxies. However, it is important to note that this  
482 analysis cannot generally be performed, as in true empirical cases we do not know whether  
483 the measured predictors are important. Further, this plateau of latent predictive ability when  
484 the number of latent predictors equals the number of environmental predictors is likely due to  
485 the fact that our abundance simulations only used linear environment-abundance  
486 relationships; it is likely that if abundance-environment relationships were nonlinear (e.g. uni-  
487 or multi-modal), a larger number of latent variables would be needed to capture the same  
488 number of environmental dimensions.

489 Additionally, although the RDA analysis demonstrated that the correct number of latents can  
490 represent the true number of environmental gradients structuring co-occurring species, it is  
491 important to note that the original simulations generated abundance values that were then  
492 transformed into presence-absence for generating latents. Although using presence-absence  
493 data allows our models to be applicable across many systems - given that researchers often  
494 only have abundance data for a few target species and presence-absence data for multiple  
495 other co-occurring species - there is certainly loss of environmental signal by doing so. This  
496 explains why the adjusted  $R^2$  is generally not very high.

497 **Model performance**

498 Our second and third objectives were aimed at contrasting model performance that contained  
499 varying levels of information (i.e., number of predictors) about the true underlying drivers  
500 versus latent predictors and assessing how predictive performance varied as a function of  
501 sample size. We first compared model performance based on the presence-absence  
502 predictions, with the goal of assessing accuracy and comparing it to current models used by  
503 management which in most cases, do not contain all relevant environmental drivers. Although  
504 our study was primarily designed to predict abundance, the ability to derive accurate  
505 presence-absence predictions would enable researchers to apply an even more general  
506 framework for species distribution modeling based on latent predictors.

507 *Presence-absence predictions*

508 As to be expected, adding relevant environmental variables to the models improves  
509 predictions. Since the species' abundance - and consequently presence-absence - is linearly  
510 related to these variables, any environmental information enables the model to capture more  
511 variation and thus predict abundance more accurately. Including all environmental variables  
512 leads to a perfect prediction. Although our goal was to develop and assess the performance of  
513 a general framework for predicting species distributions of target species based on latents of  
514 co-occurring species, different issues could be considered in future studies. For instance, the  
515 perfect prediction including all predictors was an outcome to be expected given that we did  
516 not include measurement error for environmental predictors or species abundances (i.e., white  
517 noise) in our simulations (see McInerny and Purves 2011 for potential approaches for  
518 attenuating the potential effects of environmental measurement error species distributional  
519 models). It would be interesting to perform a sensitivity analysis after including measurement  
520 errors either in the way environment (e.g., spatial variation within sites, temporal lags in

521 species responses to environments) or abundance (e.g., estimates based on mark-recapture)  
522 are measured.

523 The inclusion of species co-occurrence patterns through latent variables also leads to an  
524 improvement in predictions, indicating that the latent variables can capture unobserved  
525 environmental variation and serve as a proxy for missing (but relevant) environmental drivers.

526 Indeed, models that incorporate two environmental variables and latent variables tended to  
527 perform better than models containing only two environmental variables. This is particularly  
528 important since empirical datasets are unlikely to capture all relevant environmental drivers.

529 Although presence-absence datasets are common, a model capable of predicting the presence  
530 and absence of an invasive species or a rare species based on the rest of the community  
531 composition could be useful for conservation efforts, especially with methods such as eDNA  
532 surveys that can collect information on presence from relatively few samples (Rees et al.  
533 2014).

534 The lack of influence of number of sites sampled on model performance may initially seem  
535 surprising. However, the training set of sites used to fit the models was sampled  
536 independently of the values of the environmental variables and without measurement error.

537 This means that regardless of number of sites used to fit the model, the relationship between  
538 abundance and environment would have been accurately captured. It would be interesting to  
539 assess how changing the relationship from linear to quadratic would influence the results; as  
540 there would be increased complexity in the link, we'd expect to have a greater impact of  
541 number of sites sampled on the predictions.

#### 542 *Abundance predictions*

543 The species' average abundance was generally low in our simulations. However, since we  
544 were interested in relative abundance error rather than true abundance error, we made a

545 deliberate decision not to adjust the parameters of our simulations, maintaining a low average  
546 abundance. The shape of the abundance density curve was, to us, the most salient  
547 characteristic we aimed to replicate. Keeping the average abundance low also allowed us to  
548 maintain the occurrence of species within an ecologically meaningful range (i.e., between  
549 10% and 90% of occurrence across the landscape).

550 As expected, adding environmental variables improved the abundance predictions. Since no  
551 measurement error was included in the simulations for either environmental variables or  
552 species abundances, the inclusion of any environmental variable is likely to improve  
553 predictive accuracy. However, it is interesting to note that adding community composition  
554 only improved predictions for the high abundance species. One possible explanation for this is  
555 that the way we generated species abundances resulted in low-abundance species also being  
556 only weakly predictable from environmental variation (and thus only weakly predictable from  
557 community composition). In our simulations, a species would have low average abundance if  
558 it either had a small intercept ( $b_0$ ) and values of the environmental slopes ( $b_1$  to  $b_E$  values)  
559 close to zero (so it would be roughly equally distributed across the landscape), or if it had a  
560 very small intercept value ( $b_0$ ) and one large environmental slope value, so it was well-  
561 predicted by a single environmental variable. As such, the low predictive power of latent  
562 variables for rare species observed in our results may not generalize to species in natural  
563 settings. In fact, one might expect that species with intermediate abundances are likely to be  
564 best predicted due to the positive relationship typically observed between occupancy (number  
565 of sites occupied) and abundance (Gaston 1996; but see Wright 1991). Species with low  
566 abundances may not occupy all suitable habitats, while those with high abundances could be  
567 generalists, occupying an excess of environments. Additionally, many other non-  
568 environmental factors (e.g., biogeography, dispersal limitation, species interactions, species  
569 introductions) may play an important role in shaping patterns of species distributions and



570 biodiversity in local communities and regionally (Boulangeat, Gravel, and Thuiller 2012;  
571 Lewis et al. 2017; Guisan and Thuiller 2005). We suggest that future research could extent  
572 these simulations to incorporate nonlinear environmental gradients driving species  
573 abundances.

574 Unlike presence-absence predictions, where no pattern related to species incidence could be  
575 identified, we observe a clear trend for the abundance predictions. The more abundant a  
576 species is, the higher the model's predictive error. Since we measure the relative error in  
577 prediction and not the absolute error, this is not an artefact related to the total abundance of  
578 the species but rather it is related to the fact that the high abundance sites are poorly predicted.  
579 However, it may be due to the fact that we simulated species abundance from a Poisson  
580 distribution, where the variance in outcome increases linearly with the mean abundance,  
581 which would lead to higher variability in abundance even between sites with identical  
582 environmental variables. This does not make this result an artifact of our simulations,  
583 however; positive mean-variance relationships are typical in ecological populations (He and  
584 Gaston 2003), so we expect that it should be more difficult in general to predict abundances of  
585 common species compared to rare ones. It is important to highlight the fact that using a  
586 different statistical family to model species' abundance might allow for a better fit of the  
587 model with empirical data and further improve the predictions (see review by Waldock et al.  
588 2022). Note, however, that the main component of our framework - the use of latents based on  
589 species co-occurrence patterns to predict species abundances - can be directly applied to any  
590 modeling procedure, whether it is based on maximum likelihood, Bayesian or machine  
591 learning models.

592 One intriguing result was observing the convergence of the models' performance for low-  
593 abundance species. Indeed, for species in the 0 to 50 percentiles of abundance, regardless of  
594 the metric used, a model containing only community composition can perform as well as one

595 containing all environmental variables. This result may demonstrate the true potential for our  
596 framework as a management tool. However, again, this may be due to the Poisson expectation  
597 of our simulations as explained earlier. This performance does not apply to high abundance  
598 species, where there is a significant divergence in the models' performance, likely caused by a  
599 few sites with very high abundances. Applications to empirical datasets may require  
600 downweighing the importance of sites containing high abundances to avoid skewing the  
601 model's predictive accuracy. The use of more robust models that may account for different  
602 types of overdispersion (e.g., very low and high abundances) can be considered within the  
603 context of our framework (e.g., Poisson-log normal model, Harrison 2014).

604 Additionally, increasing the number of sites sampled did not influence predictive  
605 performance, a result we anticipated since we sampled uniformly across the landscapes and  
606 captured the entire range of variation when fitting the model. However, such uniform  
607 sampling across landscapes is unlikely to be realistic when using empirical data, particularly  
608 in complex and patchy landscapes in which environmental features are clumped and spatially  
609 autocorrelated. This issue extends beyond our study. Various approaches have been proposed  
610 to mitigate the impact of complex landscapes on the predictive performance of species  
611 distribution models based on environmental features. Different sampling methods (Fortin,  
612 Drapeau, and Legendre 1989; Christianson and Kaufman 2016), model validation techniques  
613 (Wenger and Olden 2012), and modeling frameworks (e.g., Dormann 2007 for a review,  
614 Guélat and Kéry 2018) are among these proposed solutions and could, in principle, be  
615 incorporated into our modeling framework given its flexibility.

616 We did not include any species interactions in our model simulations: as such, our results  
617 demonstrate that latent community composition variables can capture similar patterns of  
618 environmental interactions even in the absence of species interacting with one another.

619 Although latent variable models can represent species interactions (e.g., competition, trophic

620 interactions) via networks (e.g., Ovaskainen et al. 2016), adjustments to the latent extraction  
621 may be necessary in order to incorporate more complex processes underlying pattern of  
622 species co-occurrences. It is likely that including direct species interactions (e.g., competition  
623 or predation) would increase the power of latent parameters for predicting species abundances  
624 as long as strong species interactions were relatively rare, or species interaction networks are  
625 relatively sparse; strong species interactions and dense species interaction networks can result  
626 in complex feedbacks, such that the net effect of presence or absence of a given species on a  
627 focal species may be indeterminant (Tunney, Carpenter, and Vander Zanden 2017).

628 Finally, it is important to consider that we used all species in any given simulated landscape to  
629 generate latents. However, it is likely that certain reduced number of species combinations  
630 would better serve as inputs for latent generation. For instance, consider a scenario involving  
631 two species and two independent environmental predictors. If one species is highly associated  
632 with one environmental predictor but randomly associated with the other; and the second  
633 species shows the reverse pattern, then the two species will not effectively predict each other.  
634 One possible solution is to cluster species based on their environmental affinities prior to  
635 latent generation (see Hui et al. 2013 for a discussion). As such, latents could be tailored to  
636 only consider species that increase the model performance of the target species.

637 Our proposed framework offers considerable promise for several compelling reasons. First, it  
638 is highly flexible in terms of parameter estimation, as it can accommodate any regression style  
639 approach. This allows to predict both presence-absence and abundance, and it demonstrates  
640 very good performance in predicting low-abundance species. Moreover, one can also use  
641 other latent modeling procedures and not necessarily Gaussian copulas. The framework could  
642 also be used to predict biomass rather than abundance by replacing the family of the GLM  
643 used, depending on the variable of highest interest for management. Overall, our proposed

644 framework is incredibly versatile, allowing for significant flexibility and adaptability to  
645 accommodate the available data.

## 646 **Acknowledgments**

647 The data presented in this article were sampled and curated by the Ontario Ministry of Natural  
648 Resources and Forestry; we thank Dr. Cindy Chu for graciously giving us access to the  
649 Ontario Broadscale Monitoring Program data. AS was funded by the NSERC-CREATE  
650 BIOS2 program the Genome Canada GEN-FISH Network. AS and PPN were funded by the  
651 Canada Research Chair in Spatial Ecology and Biodiversity.

## 652 **Author Contributions**

653 Conceptualisation: AS, EP, PPN. Coding: AS, EP, PPN. Methodology: AS, EP, PPN.  
654 Visualisation: AS. Writing: AS, EP, PPN. All authors read and approved the final manuscript.

## 655 **Conflict of Interest Statement**

656 The authors declare no conflicts of interest.

## 657 **References**

- 658 Allouche, Omri, Asaf Tsoar, and Ronen Kadmon. 2006. “Assessing the Accuracy of Species  
659 Distribution Models: Prevalence, Kappa and the True Skill Statistic (TSS).” *Journal of*  
660 *Applied Ecology* 43 (6): 1223–32. <https://doi.org/10.1111/j.1365-2664.2006.01214.x>.
- 661 Anderson, Marti J., Perry de Valpine, Andrew Punnett, and Arden E. Miller. 2019. “A  
662 Pathway for Multivariate Analysis of Ecological Communities Using Copulas.” *Ecology*  
663 *and Evolution* 9 (6): 3276–94. <https://doi.org/10.1002/ece3.4948>.
- 664 Bengtsson, Jan, Stephen R. Baillie, and John Lawton. 1997. “Community Variability  
665 Increases with Time.” Edited by David Warton. *Oikos* 78 (2): 249–56.

666 <https://doi.org/10.2307/3546291>.

667 Blanchet, F. Guillaume, Kevin Cazelles, and Dominique Gravel. 2020. “Co-occurrence Is Not  
668 Evidence of Ecological Interactions.” Edited by Elizabeth Jeffers. *Ecology Letters* 23 (7):  
669 1050–63. <https://doi.org/10.1111/ele.13525>.

670 Boulangeat, Isabelle, Dominique Gravel, and Wilfried Thuiller. 2012. “Accounting for  
671 Dispersal and Biotic Interactions to Disentangle the Drivers of Species Distributions and  
672 Their Abundances.” *Ecology Letters* 15 (6): 584–93. [https://doi.org/10.1111/j.1461-  
673 0248.2012.01772.x](https://doi.org/10.1111/j.1461-0248.2012.01772.x).

674 Boyce, Mark S., Chris J. Johnson, Evelyn H. Merrill, Scott E. Nielsen, Erling J. Solberg, and  
675 Bram van Moorter. 2016. “REVIEW: Can Habitat Selection Predict Abundance?” Edited  
676 by Luca Börger. *Journal of Animal Ecology* 85 (1): 11–20. [https://doi.org/10.1111/1365-  
677 2656.12359](https://doi.org/10.1111/1365-2656.12359).

678 Braak, Cajo J.F. ter. 1985. “Correspondence Analysis Data: In Terms of a Unimodal  
679 Properties Response Model.” *Biometrics* 41 (4): 859–73.

680 Braak, Cajo J.F. ter, and I. Colin Prentice. 1988. “A Theory of Gradient Analysis.” In  
681 *Advances in Ecological Research*, 18:271–317. [https://doi.org/10.1016/S0065-  
682 2504\(08\)60183-X](https://doi.org/10.1016/S0065-2504(08)60183-X).

683 Bradley, Bethany A. 2016. “Predicting Abundance with Presence-Only Models.” *Landscape  
684 Ecology* 31 (1): 19–30. <https://doi.org/10.1007/s10980-015-0303-4>.

685 Brosse, S., Jean-François Guegan, Jean-Nöel Tourenq, and Sovan Lek. 1999. “The Use of  
686 Artificial Neural Networks to Assess Fish Abundance and Spatial Occupancy in the  
687 Littoral Zone of a Mesotrophic Lake.” *Ecological Modelling* 120 (2–3): 299–311.  
688 [https://doi.org/10.1016/S0304-3800\(99\)00110-6](https://doi.org/10.1016/S0304-3800(99)00110-6).

- 689 Carreira-Perpinán, MA. 1997. “A Review of Dimension Reduction Techniques.” *Department*  
690 *of Computer Science. University of Sheffield. Tech. Rep. CS-96-09*, no. 9: 1–69.  
691 <http://www.pca.narod.ru/DimensionReductionBrifReview.pdf>.
- 692 Christianson, Danielle S., and Cari G. Kaufman. 2016. “Effects of Sample Design and  
693 Landscape Features on a Measure of Environmental Heterogeneity.” Edited by Robert  
694 Freckleton. *Methods in Ecology and Evolution* 7 (7): 770–82.  
695 <https://doi.org/10.1111/2041-210X.12539>.
- 696 Cunningham, Pádraig. 2008. “Dimension Reduction.” In *Machine Learning Techniques for*  
697 *Multimedia*, 91–112. Berlin, Heidelberg: Springer Berlin Heidelberg.  
698 [https://doi.org/10.1007/978-3-540-75171-7\\_4](https://doi.org/10.1007/978-3-540-75171-7_4).
- 699 Degnbol, P., and A. Jarre. 2004. “Review of Indicators in Fisheries Management – a  
700 Development Perspective.” *African Journal of Marine Science* 26 (1): 303–26.  
701 <https://doi.org/10.2989/18142320409504063>.
- 702 Dormann, Carsten F. 2007. “Effects of Incorporating Spatial Autocorrelation into the  
703 Analysis of Species Distribution Data.” *Global Ecology and Biogeography* 16 (2): 129–  
704 38. <https://doi.org/10.1111/j.1466-8238.2006.00279.x>.
- 705 Dornelas, Maria, Anne E Magurran, Stephen T Buckland, Anne Chao, Robin L Chazdon,  
706 Robert K Colwell, Tom Curtis, et al. 2013. “Quantifying Temporal Change in  
707 Biodiversity: Challenges and Opportunities.” *Proceedings of the Royal Society B:*  
708 *Biological Sciences* 280 (1750): 20121931. <https://doi.org/10.1098/rspb.2012.1931>.
- 709 Flecker, Alexander S., and William J. Matthews. 1999. “Patterns in Freshwater Fish  
710 Ecology.” *Copeia* 1999 (1): 229–30. <https://doi.org/10.2307/1447409>.
- 711 Fortin, Marie-jose, Pierre Drapeau, and Pierre Legendre. 1989. “Spatial Autocorrelation and

712 Sampling Design in Plant Ecology.” *Vegetatio* 83 (1–2): 209–22.  
713 <https://doi.org/10.1007/BF00031693>.

714 Gaston, Kevin J. 1996. “The Multiple Forms of the Interspecific Abundance-Distribution  
715 Relationship.” *Oikos* 76 (2): 211–20. <https://doi.org/10.2307/3546192>.

716 González-Salazar, Constantino, Christopher R. Stephens, and Pablo A. Marquet. 2013.  
717 “Comparing the Relative Contributions of Biotic and Abiotic Factors as Mediators of  
718 Species’ Distributions.” *Ecological Modelling* 248 (January): 57–70.  
719 <https://doi.org/10.1016/j.ecolmodel.2012.10.007>.

720 Guélat, Jérôme, and Marc Kéry. 2018. “Effects of Spatial Autocorrelation and Imperfect  
721 Detection on Species Distribution Models.” Edited by Nick Isaac. *Methods in Ecology  
722 and Evolution* 9 (6): 1614–25. <https://doi.org/10.1111/2041-210X.12983>.

723 Guisan, Antoine, and Wilfried Thuiller. 2005. “Predicting Species Distribution: Offering  
724 More than Simple Habitat Models.” *Ecology Letters* 8 (9): 993–1009.  
725 <https://doi.org/10.1111/j.1461-0248.2005.00792.x>.

726 Harrison, Xavier A. 2014. “Using Observation-Level Random Effects to Model  
727 Overdispersion in Count Data in Ecology and Evolution.” *PeerJ* 2 (1): e616.  
728 <https://doi.org/10.7717/peerj.616>.

729 He, Fangliang, and Kevin J. Gaston. 2003. “Occupancy, Spatial Variance, and the Abundance  
730 of Species.” *The American Naturalist* 162 (3): 366–75. <https://doi.org/10.1086/377190>.

731 Hui, Francis K C, David I Warton, Scott D Foster, and Piers K Dunstan. 2013. “To Mix or  
732 Not to Mix: Comparing the Predictive Performance of Mixture Models vs. Separate  
733 Species Distribution Models.” *Ecology* 94 (9): 1913–19. <https://doi.org/10.1890/12-1322.1>.

735 Jackson, Donald A, and Harold H Harvey. 1997. “Qualitative and Quantitative Sampling of  
736 Lake Fish Communities.” *Canadian Journal of Fisheries and Aquatic Sciences* 54 (12):  
737 2807–13. <https://doi.org/10.1139/f97-182>.

738 Kéry, Marc, and J. Andrew Royle. 2015. *Applied Hierarchical Modeling in Ecology: Analysis*  
739 *of Distribution, Abundance and Species Richness in R and BUGS*. Vol. 1. Elsevier.  
740 <https://doi.org/10.1016/C2015-0-04070-9>.

741 Lawson, Callum R., Jenny A. Hodgson, Robert J. Wilson, and Shane A. Richards. 2014.  
742 “Prevalence, Thresholds and the Performance of Presence–Absence Models.” Edited by  
743 Robert Freckleton. *Methods in Ecology and Evolution* 5 (1): 54–64.  
744 <https://doi.org/10.1111/2041-210X.12123>.

745 Lek, Sovan, Alain Belaud, Philippe Baran, Ioannis Dimopoulos, and Marc Delacoste. 1996.  
746 “Role of Some Environmental Variables in Trout Abundance Models Using Neural  
747 Networks.” *Aquatic Living Resources* 9 (1): 23–29. <https://doi.org/10.1051/alr:1996004>.

748 Lewis, Jesse S., Matthew L. Farnsworth, Chris L. Burdett, David M. Theobald, Miranda Gray,  
749 and Ryan S. Miller. 2017. “Biotic and Abiotic Factors Predicting the Global Distribution  
750 and Population Density of an Invasive Large Mammal.” *Scientific Reports* 7 (1): 44152.  
751 <https://doi.org/10.1038/srep44152>.

752 Liu, Canran, Pam M Berry, Terence P Dawson, and Richard G Pearson. 2005. “Selecting  
753 Thresholds of Occurrence in the Prediction of Species Distributions.” *Ecography* 28 (3):  
754 385–93. <http://www.jstor.org/stable/3683850>.

755 McNerny, Greg J., and Drew W. Purves. 2011. “Fine-scale Environmental Variation in  
756 Species Distribution Modelling: Regression Dilution, Latent Variables and Neighbourly  
757 Advice.” *Methods in Ecology and Evolution* 2 (3): 248–57.  
758 <https://doi.org/10.1111/j.2041-210X.2010.00077.x>.



759 Mindrila, Diana. 2023. “Bayesian Latent Class Analysis: Sample Size, Model Size, and  
760 Classification Precision.” *Mathematics* 11 (12): 2753.  
761 <https://doi.org/10.3390/math11122753>.

762 Myers, Ransom A. 1998. “When Do Environment-Recruitment Correlations Work?” *Reviews*  
763 *in Fish Biology and Fisheries* 8 (3): 285–305. <https://doi.org/10.1023/A:1008828730759>.

764 Nylund, Karen L., Tihomir Asparouhov, and Bengt O. Muthén. 2007. “Deciding on the  
765 Number of Classes in Latent Class Analysis and Growth Mixture Modeling: A Monte  
766 Carlo Simulation Study.” *Structural Equation Modeling: A Multidisciplinary Journal* 14  
767 (4): 535–69. <https://doi.org/10.1080/10705510701575396>.

768 Öglü, Burak, Upendra Bhele, Ain Järvalt, Lea Tuvikene, Henn Timm, Siim Seller, Juta  
769 Haberman, et al. 2019. “Is Fish Biomass Controlled by Abiotic or Biotic Factors?  
770 Results of Long-Term Monitoring in a Large Eutrophic Lake.” *Journal of Great Lakes*  
771 *Research* 46 (4): 881–90. <https://doi.org/10.1016/j.jglr.2019.08.004>.

772 Oksanen, Jari, Gavin L Simpson, F Guillaume Blanchet, Roeland Kindt, Pierre Legendre,  
773 Peter R Minchin, R B O’Hara, et al. 2022. “Vegan: Community Ecology Package.”  
774 <https://cran.r-project.org/package=vegan>.

775 Olkeba, Beekam Kebede, Pieter Boets, Seid Tiku Mereta, Mesfin Yeshigeta, Geremew  
776 Muleta Akessa, Argaw Ambelu, and Peter L. M. Goethals. 2020. “Environmental and  
777 Biotic Factors Affecting Freshwater Snail Intermediate Hosts in the Ethiopian Rift  
778 Valley Region.” *Parasites & Vectors* 13 (1): 292. [https://doi.org/10.1186/s13071-020-](https://doi.org/10.1186/s13071-020-04163-6)  
779 [04163-6](https://doi.org/10.1186/s13071-020-04163-6).

780 Ovaskainen, Otso, Nerea Abrego, Panu Halme, and David Dunson. 2016. “Using Latent  
781 Variable Models to Identify Large Networks of Species-to-species Associations at  
782 Different Spatial Scales.” Edited by David Warton. *Methods in Ecology and Evolution* 7

783 (5): 549–55. <https://doi.org/10.1111/2041-210X.12501>.

784 Peirce, C S. 1884. “The Numerical Measure of the Success of Predictions.” *Science (New*  
785 *York, N.Y.)* 4 (93): 453–54. <https://doi.org/10.1126/science.ns-4.93.453-a>.

786 Peres-Neto, Pedro R., Pierre Legendre, Stéphane Dray, and Daniel Borcard. 2006. “Variation  
787 Partitioning of Species Data Matrices: Estimation and Comparison of Fractions.”  
788 *Ecology* 87 (10): 2614–25. [https://doi.org/10.1890/0012-9658\(2006\)87](https://doi.org/10.1890/0012-9658(2006)87).

789 Persson, Lennart. 2008. “Community Ecology of Freshwater Fishes.” In *Handbook of Fish*  
790 *Biology and Fisheries, Volume 1*, 321–40. Oxford, UK: Blackwell Publishing Ltd.  
791 <https://doi.org/10.1002/9780470693803.ch15>.

792 Phillips, Steven J., and Jane Elith. 2013. “On Estimating Probability of Presence from Use–  
793 Availability or Presence–Background Data.” *Ecology* 94 (6): 1409–19.  
794 <https://doi.org/10.1890/12-1520.1>.

795 Popovic, Gordana C., David I. Warton, Fiona J. Thomson, Francis K.C. Hui, and Angela T.  
796 Moles. 2019. “Untangling Direct Species Associations from Indirect Mediator Species  
797 Effects with Graphical Models.” *Methods in Ecology and Evolution* 10 (9): 1571–83.  
798 <https://doi.org/10.1111/2041-210X.13247>.

799 Popovic, Gordana C, Francis K C Hui, and David I Warton. 2022. “Fast Model-based  
800 Ordination with Copulas.” *Methods in Ecology and Evolution* 13 (1): 194–202.  
801 <https://doi.org/10.1111/2041-210X.13733>.

802 Rees, Helen C., Ben C. Maddison, David J. Middleditch, James R.M. Patmore, and Kevin C.  
803 Gough. 2014. “The Detection of Aquatic Animal Species Using Environmental DNA – a  
804 Review of EDNA as a Survey Tool in Ecology.” Edited by Erika Crispo. *Journal of*  
805 *Applied Ecology* 51 (5): 1450–59. <https://doi.org/10.1111/1365-2664.12306>.

806 Sobrino, Ignacio, Lucia Rueda, Maria Pilar Tugores, Candelaria Burgos, Miguel Cojan, and  
807 Graham J. Pierce. 2020. “Abundance Prediction and Influence of Environmental  
808 Parameters in the Abundance of Octopus (*Octopus Vulgaris* Cuvier, 1797) in the Gulf of  
809 Cadiz.” *Fisheries Research* 221 (August 2019): 105382.  
810 <https://doi.org/10.1016/j.fishres.2019.105382>.

811 Tunney, Tyler D., Stephen R. Carpenter, and M. Jake Vander Zanden. 2017. “The  
812 Consistency of a Species’ Response to Press Perturbations with High Food Web  
813 Uncertainty.” *Ecology* 98 (7): 1859–68. <https://doi.org/10.1002/ecy.1853>.

814 VanDerWal, Jeremy, Luke P. Shoo, Christopher N. Johnson, and Stephen E. Williams. 2009.  
815 “Abundance and the Environmental Niche: Environmental Suitability Estimated from  
816 Niche Models Predicts the Upper Limit of Local Abundance.” *The American Naturalist*  
817 174 (2): 282–91. <https://doi.org/10.1086/600087>.

818 Waldock, Conor, Rick D. Stuart-Smith, Camille Albouy, William W. L. Cheung, Graham J.  
819 Edgar, David Mouillot, Jerry Tjiputra, and Loïc Pellissier. 2022. “A Quantitative Review  
820 of Abundance-based Species Distribution Models.” *Ecography* 2022 (1): 1–18.  
821 <https://doi.org/10.1111/ecog.05694>.

822 Walker, Steven C., and Donald A. Jackson. 2011. “Random-Effects Ordination: Describing  
823 and Predicting Multivariate Correlations and Co-Occurrences.” *Ecological Monographs*  
824 81 (4): 635–63. <https://doi.org/10.1890/11-0886.1>.

825 Wang, Yi, Ulrike Naumann, Dirk Eddelbuettel, John Wilshire, and David Warton. 2022.  
826 “Mvabund: Statistical Methods for Analysing Multivariate Abundance Data.”  
827 <https://cran.r-project.org/package=mvabund>.

828 Weller, Bridget E., Natasha K. Bowen, and Sarah J. Faubert. 2020. “Latent Class Analysis: A  
829 Guide to Best Practice.” *Journal of Black Psychology* 46 (4): 287–311.

830 <https://doi.org/10.1177/0095798420930932>.

831 Wenger, Seth J., and Julian D. Olden. 2012. “Assessing Transferability of Ecological Models:  
832 An Underappreciated Aspect of Statistical Validation.” *Methods in Ecology and*  
833 *Evolution* 3 (2): 260–67. <https://doi.org/10.1111/j.2041-210X.2011.00170.x>.

834 Wright, David Hamilton. 1991. “Correlations between Incidence and Abundance Are  
835 Expected by Chance.” *Journal of Biogeography* 18 (4): 463–66.  
836 <https://doi.org/10.2307/2845487>.

837 Zou, Hui, and Hao Helen Zhang. 2009. “On the Adaptive Elastic-Net with a Diverging  
838 Number of Parameters.” *The Annals of Statistics* 37 (4): 1733–51.  
839 <https://doi.org/10.1214/08-AOS625>.

840

841 **Tables**

842 **Table 1.** Variable symbols and indexes, and their associated values and distributions used in  
 843 the simulation study. Bold letters indicate that the variable is a vector or a matrix.

<b>Variable name</b>	<b>Variable</b>	<b>Values</b>
<b>A</b>	Abundance	0 to $\infty$
<b>S, s</b>	Number of species, species index	{10, 20, 30}
<b>U, u</b>	Number of landscapes, landscape index	30
<b>J, j</b>	Number of sites, site index	
<b>E</b>	Number of environmental variables	3
$b_{0,s,u}$	Intercept for species s and landscape u	Uniform(-2.4, 1.2)
$b_{1,s,u}$ to $b_{E,s,u}$	Slopes for species s, landscape u and environmental variables 1 to E	Uniform(-0.8, 0.8)
$X_{1,u,j}$ to $X_{E,j,u}$	Environmental variables 1 to E for site j of landscape u	Normal(0,1)
<b>L</b>	Number of latent variables	3
<b>X</b>	Environmental variable	
<b>Z</b>	Latent variable	

844

845 **Table 2.** All models considered in this study based on combinations of environmental  
846 variables and community composition (latents). The best model is expected to be the true  
847 model considering all environmental variables.  $\mathbf{A}$  refers to the abundance matrix,  $\mathbf{X}_1$  to  $\mathbf{X}_3$  to  
848 the environmental variables, and  $\mathbf{Z}_1$  to  $\mathbf{Z}_3$  to the community composition (latent variables).

Variables included	Model specification	Regression formula
Environmental variables	3 environmental variables	$\mathbf{A} \sim \mathbf{X}_1 + \mathbf{X}_2 + \mathbf{X}_3$
	2 environmental variables	$\mathbf{A} \sim \mathbf{X}_1 + \mathbf{X}_2$
		$\mathbf{A} \sim \mathbf{X}_1 + \mathbf{X}_3$
		$\mathbf{A} \sim \mathbf{X}_2 + \mathbf{X}_3$
	1 environmental variable	$\mathbf{A} \sim \mathbf{X}_1$
		$\mathbf{A} \sim \mathbf{X}_2$
$\mathbf{A} \sim \mathbf{X}_3$		
Environmental variables and community composition	2 environmental variables and community composition	$\mathbf{A} \sim \mathbf{X}_1 + \mathbf{X}_2 + \mathbf{Z}_1 : \mathbf{Z}_3$
		$\mathbf{A} \sim \mathbf{X}_1 + \mathbf{X}_3 + \mathbf{Z}_1 : \mathbf{Z}_3$
		$\mathbf{A} \sim \mathbf{X}_2 + \mathbf{X}_3 + \mathbf{Z}_1 : \mathbf{Z}_3$
	1 environmental variable and community composition	$\mathbf{A} \sim \mathbf{X}_1 + \mathbf{Z}_1 : \mathbf{Z}_3$
		$\mathbf{A} \sim \mathbf{X}_2 + \mathbf{Z}_1 : \mathbf{Z}_3$
		$\mathbf{A} \sim \mathbf{X}_2 + \mathbf{Z}_1 : \mathbf{Z}_3$

849

850 **Table 3.** Metrics used for assessing model predictive performance based on presence-absence  
851 and abundance of target species.  $J$  represents the number of sites,  $A_s$  the true abundance of the  
852 (target) species,  $P_s$  the predicted abundance, TP the true positives, FP the false positives, TN  
853 the true negatives, and FN the false negatives. Bold letters indicate that the variable is a vector  
854 or a matrix. The True Skill Statistic (TSS), sensitivity, and specificity are calculated for all  
855 sites of the landscape. Having evaluated the presence-absence predictions of the models and  
856 to avoid artificially inflating the error rate of the abundance metrics, the Mean Absolute  
857 Percentage Error (MAPE), Root Mean Squared Percentage Error (RMSPE), Relative Mean  
858 Squared Error (RMSE), Symmetric Mean Absolute Percentage Error (SMAPE), and Root  
859 Mean Ratio Percentage Error (RMRPE) are calculated for sites where the species is truly  
860 present (i.e., abundance of 1 or more).

Metric	Equation
TSS	$TSS = \frac{TP}{TP + FN} + \frac{TN}{TN + FP} - 1$
Sensitivity	$Sensitivity = \frac{TP}{TP + FN}$
Specificity	$Specificity = \frac{TN}{TN + FP}$
MAPE	$MAPE = \frac{1}{J} \sum_s \frac{ A_s - P_s }{A_s} \times 100$
RMSPE	$RMSPE = \sqrt{\frac{1}{J} \sum_s \left( \frac{A_s - P_s}{A_s} \right)^2} \times 100$
RMSE	$RMSE = \sqrt{\frac{1}{J} \sum_s \frac{(A_s - P_s)^2}{A_s^2}} \times 100$
SMAPE	$SMAPE = \frac{1}{J} \sum_s \frac{ A_s - P_s }{ A_s  +  P_s } \times 100$
RMRPE	$RMRPE = \sqrt{\frac{1}{J} \sum_s \log \left( \frac{P_s}{A_s} \right)^2} \times 100$

861

862 **Figure captions**

863 **Figure 1.** The rationale underlying our model framework and simulation workflow to assess  
864 its performance. First, species abundances were simulated for all species (top left panel) as a  
865 function of multiple environmental factors. In this example, two environmental variables were  
866 used to simulate species abundances ( $X_1$  and  $X_2$ ; bottom left panel). Species abundances are  
867 then transformed into presence-absence data and used to derive latent variables (bottom left  
868 panel). Here, only one latent variable is presented for simplicity, allowing one to more easily  
869 its association with the abundances of the original simulated species. Variation in species  
870 abundances (target species) across sites is then modeled against latent and environmental  
871 variables or reduced combinations (e.g., removing an environmental variable and assess the  
872 conditions that affect latent performances), depending on specific simulation scenarios. The  
873 model can produce either abundance or presence-absence predictions for each site. The black  
874 rectangular outline highlights the target species (species 10) that the model aims at predicting.

875 **Figure 2.** The density of average species abundance across sites within each landscape. For  
876 each landscape, we calculated the average abundance of each species and plotted the density  
877 of abundances in each of the 30 landscapes (grey lines). We also plotted the density of  
878 abundances across all landscapes to represent the average landscape (black line). The red line  
879 is a reference line indicates the probability density function of a log-normal distribution with  
880 the same log-mean and log-standard deviation of the average abundance distribution across  
881 replicates.

882 **Figure 3.** Variation in adjusted  $R^2$  as a function of the number of latent variables used, as well  
883 as the true dimensions of the environment and the number of species in the landscape. Here  
884 we used 500 sites, and variations according to other number of sites are presented in  
885 Appendix S1: Figure S2. Colors represent the varying number of species in the landscape, and



886 each panel indicates the true dimension of the environment (i.e., number of environmental  
887 variables used to simulate the abundance of a given target species).

888 **Figure 4.** Ratio TSS and delta TSS for each model and bin of species occurrence percentiles.  
889 The ratio TSS was averaged across all landscapes and replicates per model and species, with  
890 species binned by percentile of occurrence (percentage of sites occupied) and divided by the  
891 TSS of the oracle model. A value of 1 for the ratio TSS indicates an identical performance  
892 between the model and the oracle model, while a value below 0 represents a performance  
893 similar to that of a random model. To improve contrast between colors, we confined the color  
894 scheme between 0 and 1. Any value below 0 indicates a prediction of presence-absence no  
895 better than a random model, and any value above 1 a better prediction than the oracle model.  
896 The environment panel represents models containing only environmental variables, while the  
897 latent panel is for models containing latent variables (mix of latent and environmental  
898 predictors); the models were then ordered from bottom to top as fewest to the greatest number  
899 of environmental variables included and sorted by coefficients relative to each environmental  
900 variable (see Methods for more information, note that the “mid” model refers to the  
901 “intermediate” model). The delta TSS was measured as the TSS of the model with  
902 environmental variables minus the TSS of the model with the same combination of  
903 environmental variables and latent variables. A negative value indicates that the model with  
904 latent predicts the presence-absence of the species better than the model containing only  
905 environmental variables.

906 **Figure 5.** Correlation between the metrics studied (TSS, sensitivity, and specificity)  
907 depending on the model across species occurrence percentiles. The vertical panels indicate the  
908 different metrics, with models represented in different colors. The oracle model refers to the  
909 model using the true environmental coefficients, while the other models were fitted using all  
910 environmental variables (benchmark) or latent variables (latent). The True Skill Statistic

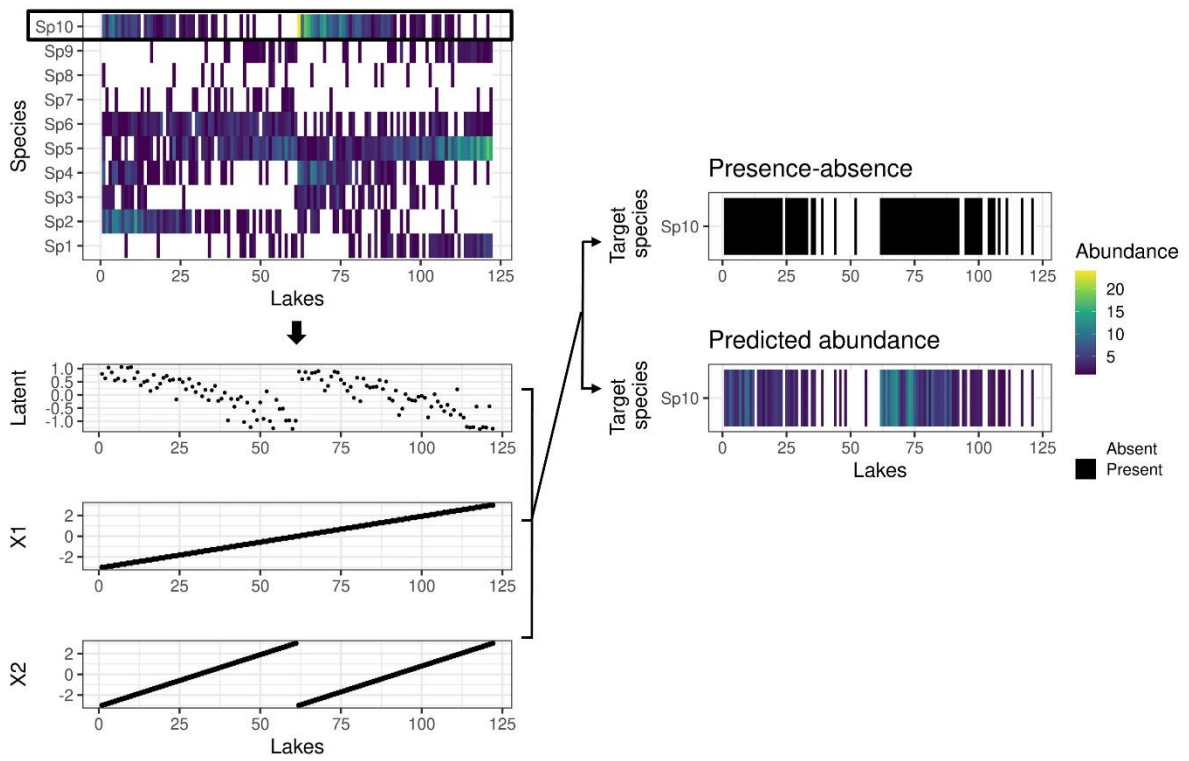
911 (TSS) measures the difference between sensitivity and specificity of the model and ranges  
912 from -1 to +1. A score of +1 indicates a perfect agreement between the predictions of the  
913 model and the true presence-absence, while a score of 0 or less represents a performance no  
914 better than random. Sensitivity represents the ability to correctly classify a species as  
915 “present”, while specificity represents the ability to correctly classify a species as “absent”.  
916 Their values can be interpreted as a percentage, with values of 1 indicating perfect  
917 classification of either presence or absence, and values of 0.5 no better than random. Here we  
918 used 500 sites, and variations according to other number of sites are presented in Appendix  
919 S1: Figure S3.

920 **Figure 6.** Ratio Mean Absolute Percentage (MAPE) and delta MAPE are presented for each  
921 model and bins of species abundance percentiles. The MAPE is averaged across all  
922 landscapes and replicates per model and species, with the species binned by percentile of  
923 abundance and divided by the MAPE of the oracle model to derive the ratio MAPE. The  
924 environment panel represents models containing only environmental variables, while the  
925 latent panel depicts models containing latent predictors. The models are then ordered from  
926 bottom to top, from the fewest to the greatest number of environmental variables included and  
927 sorted by coefficients relative to each environmental variable. See Methods for more  
928 information, note that the “mid” model refers to the “intermediate” model. Delta MAPE was  
929 measured as the MAPE of the model with environmental variables only minus the MAPE of  
930 the model with the same combination of environmental and latent predictors. A positive value  
931 indicates that the model with latent predicts the abundance of the species better than the  
932 model containing only environmental variables.

933

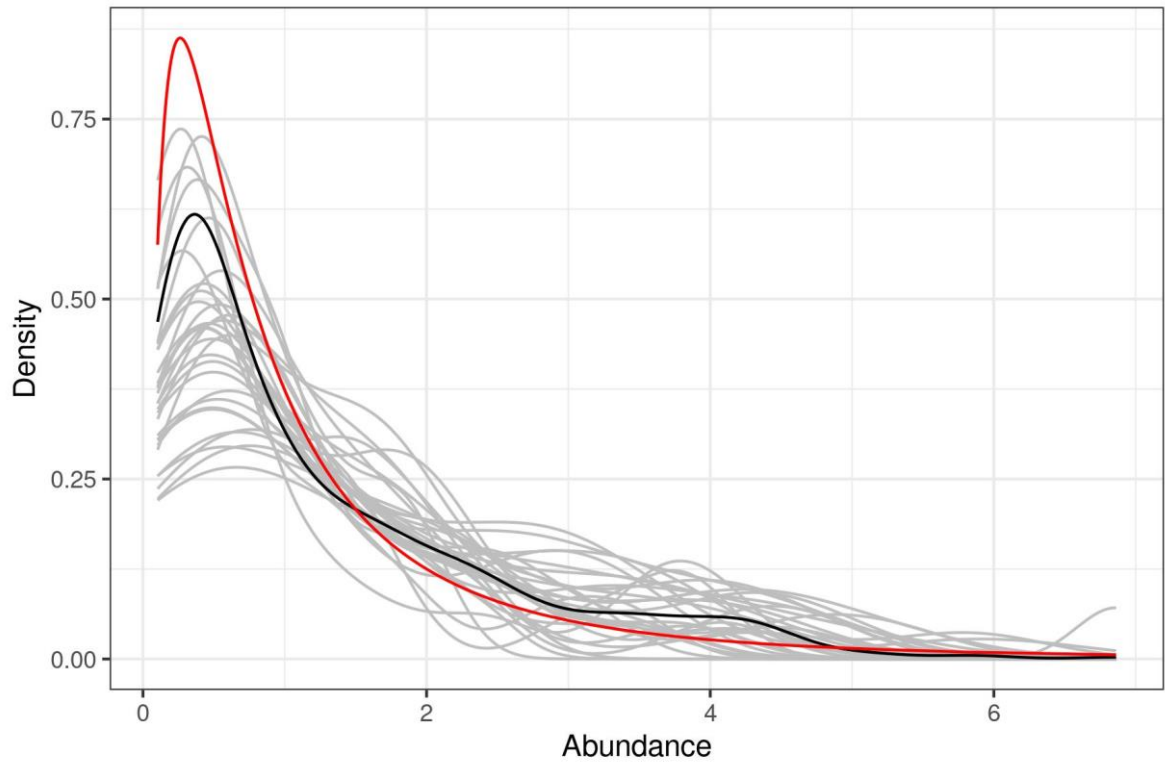
934

935 **Figures**



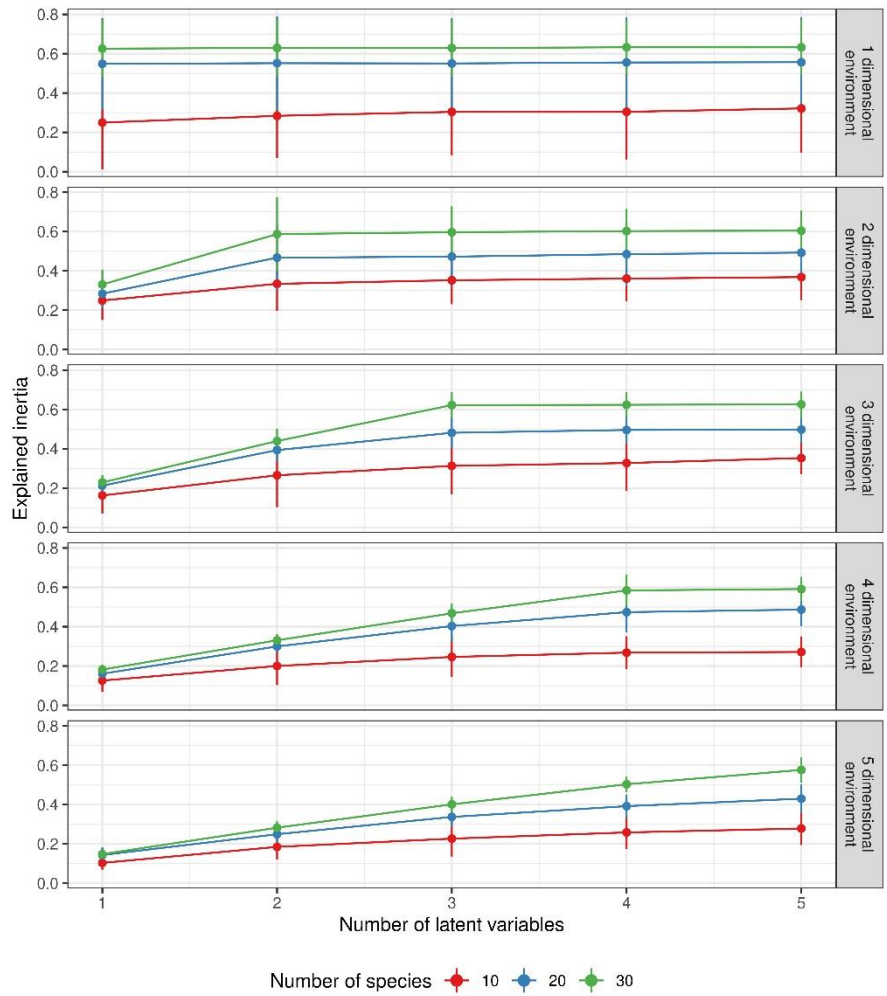
936

937 **Figure 1.**



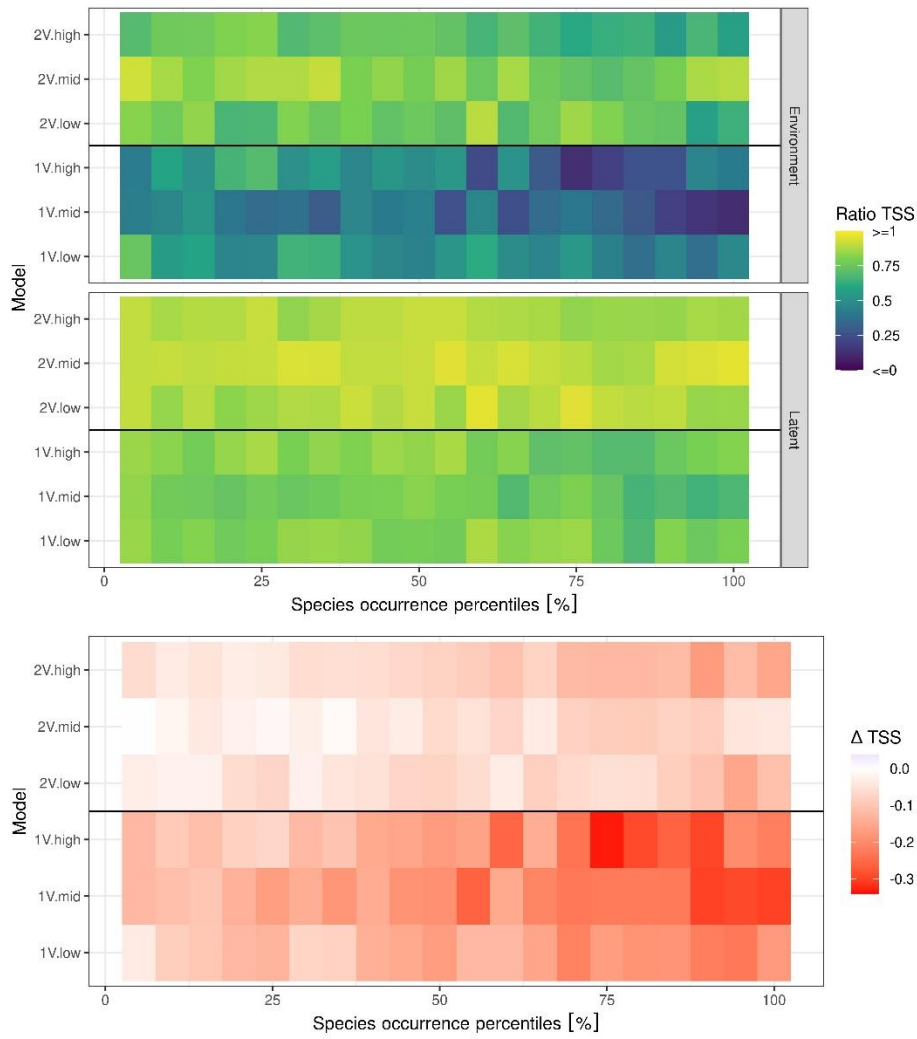
938

939 **Figure 2.**



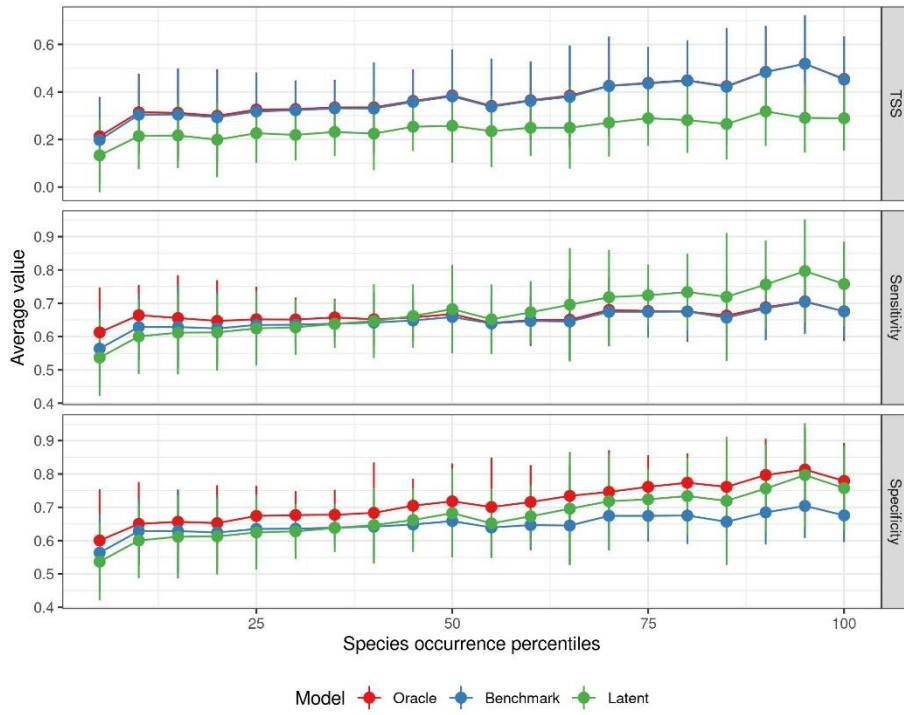
940

941 **Figure 3.**



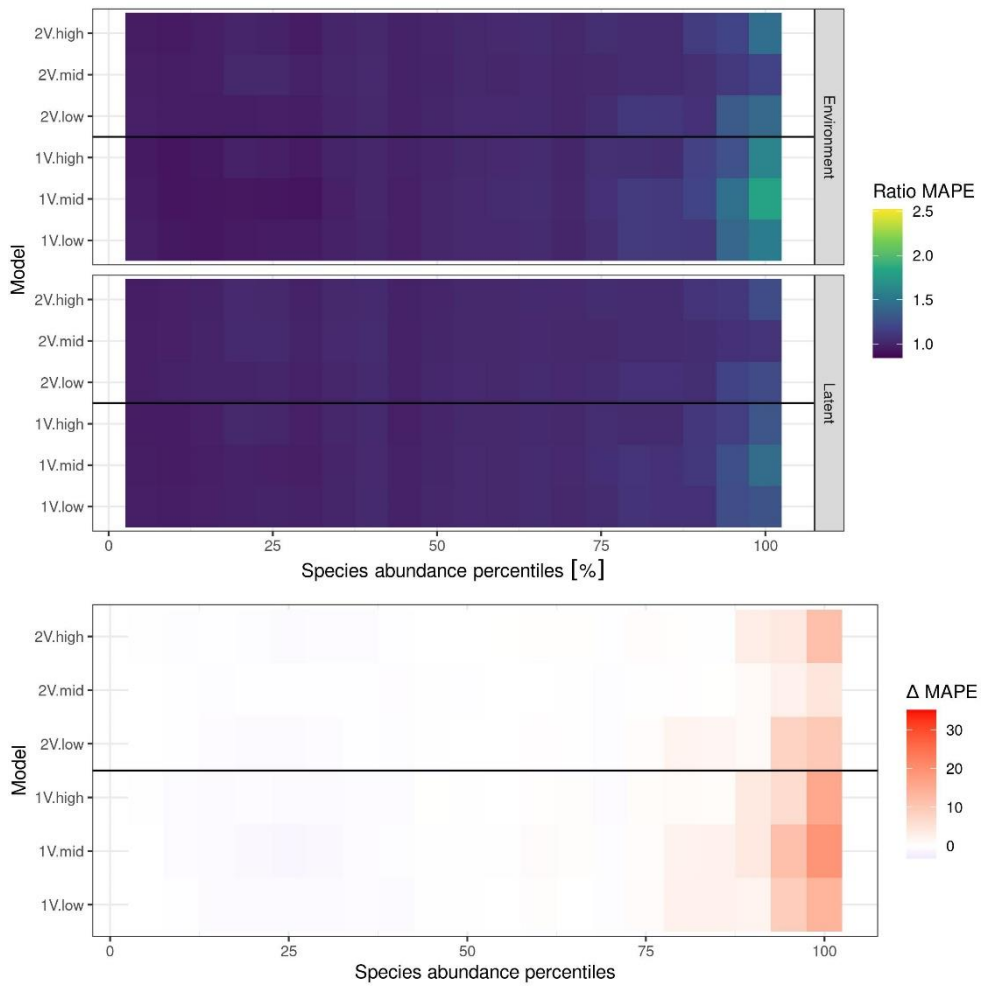
942

943 **Figure 4.**



944

945 **Figure 5.**



946

947 **Figure 6.**



## **Supporting Information for review and publication**

Journal name: Ecological Applications (ESA)

Manuscript type: Article

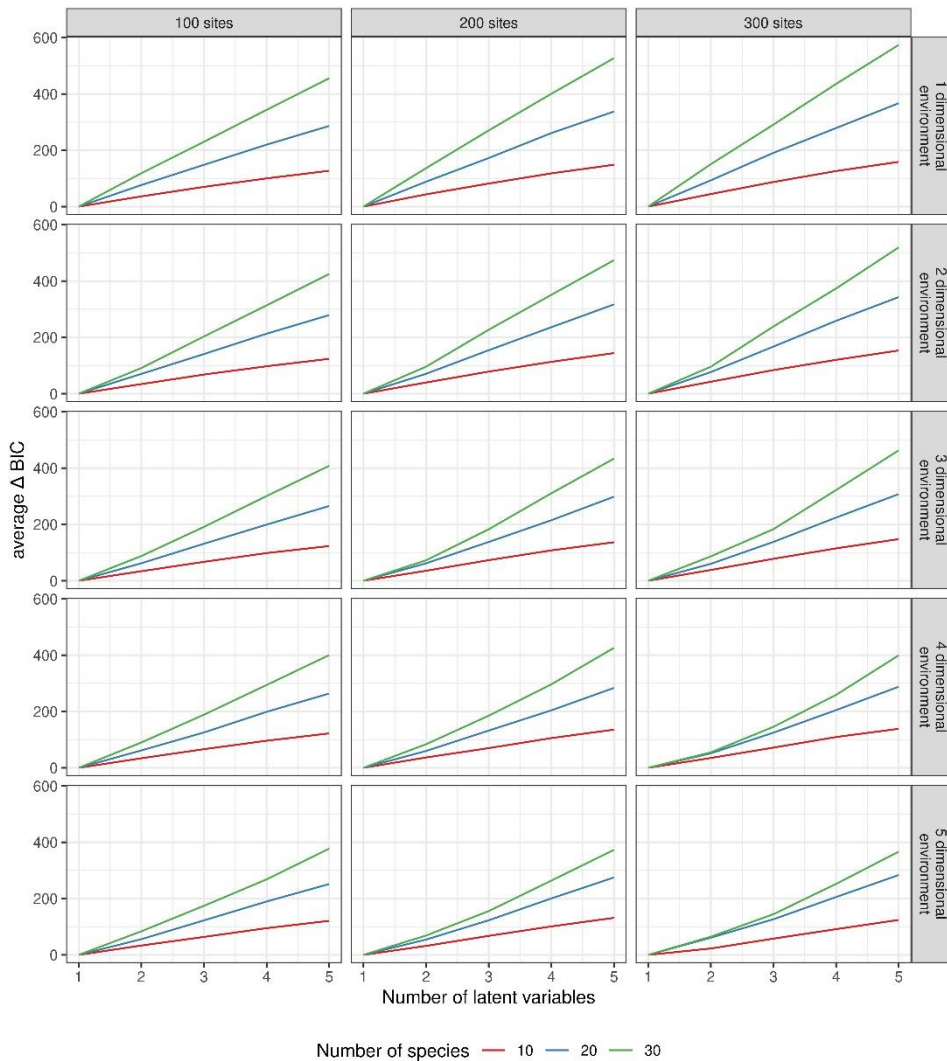
Manuscript title: Advancing single species abundance models: robust models for predicting abundance using co-occurrence from communities

Author names: Aliénor Stahl<sup>1</sup>, Eric J. Pedersen<sup>1</sup> and Pedro R. Peres-Neto<sup>1</sup>

Affiliations: 1: Concordia University, Department of Biology, Montreal, Canada

Corresponding author: Aliénor Stahl, [alienor.stahl@gmail.com](mailto:alienor.stahl@gmail.com)

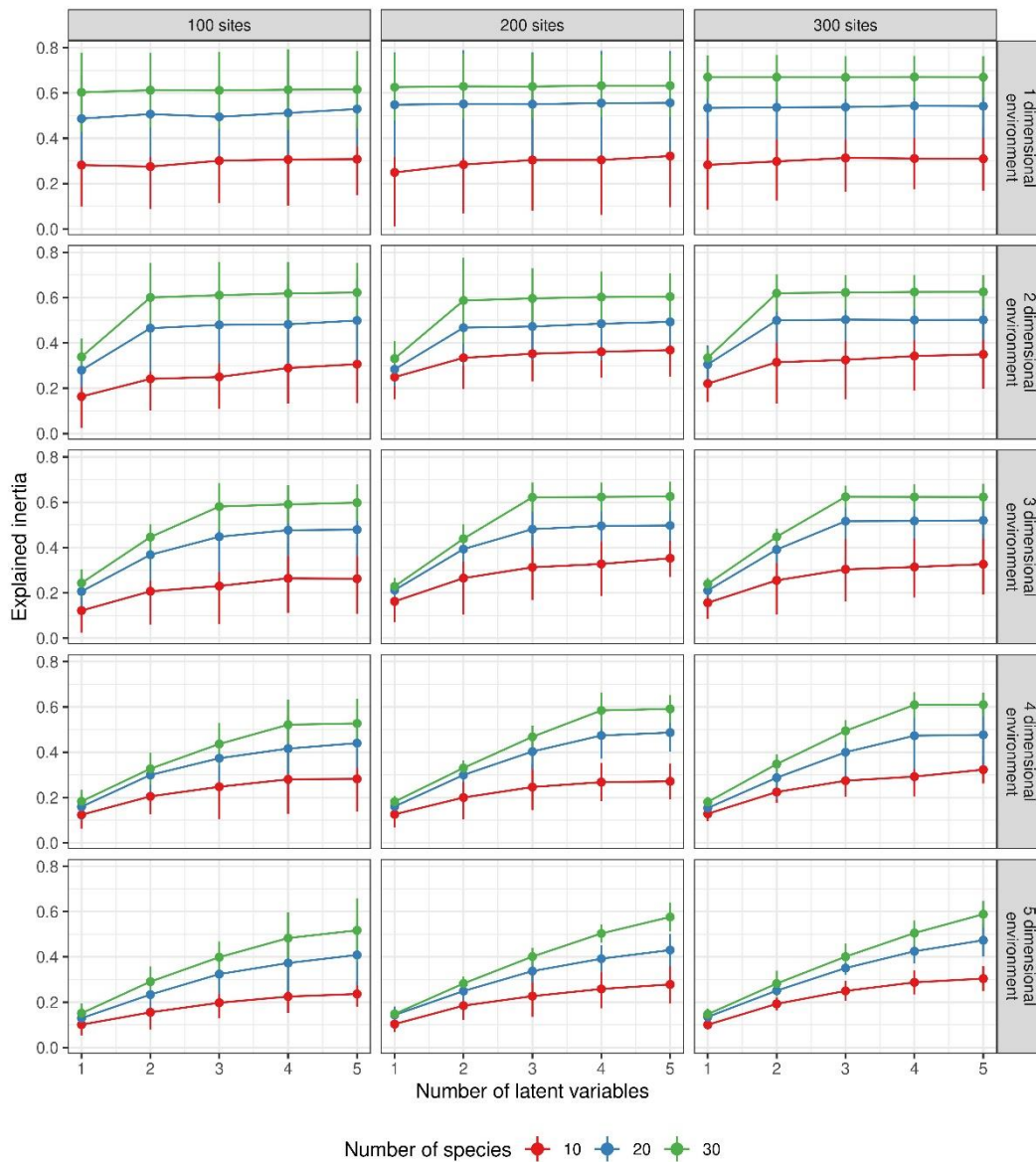
# 1 Appendix S1



2

3 **Figure S1.** Variation in delta BIC as a function of the number of latent variables used, as well  
4 as the true dimensions of the environment, the number of species in the landscape and the  
5 number of sites. Horizontal panels represent the number of sites, and each vertical panel  
6 indicates the true dimension of the environment (i.e., number of environmental variables used  
7 to simulate the abundance of a given target species). Colors represent the varying number of  
8 species in the landscape. The delta BIC is calculated as the BIC of the model minus the BIC  
9 of the best model for the ongoing simulation.

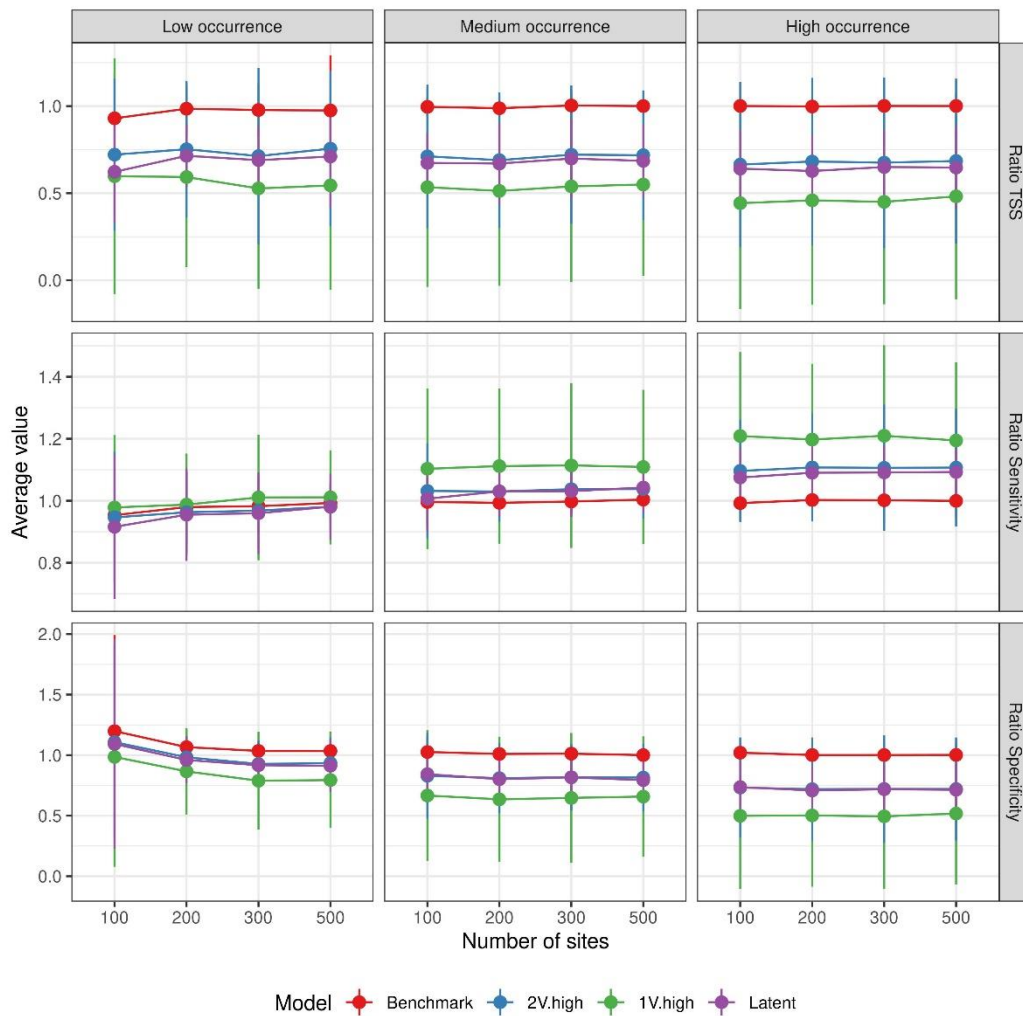
10



11

12 **Figure S2.** Variation in adjusted  $R^2$  as a function of the number of latent variables used, as  
 13 well as the true dimensions of the environment, the number of species in the landscape and  
 14 the number of sites. Horizontal panels represent the varying number of sites, and each vertical  
 15 panel indicates the true dimension of the environment (i.e., number of environmental  
 16 variables used to simulate the abundance of a given target species). Colors represent the  
 17 varying number of species in the landscape.

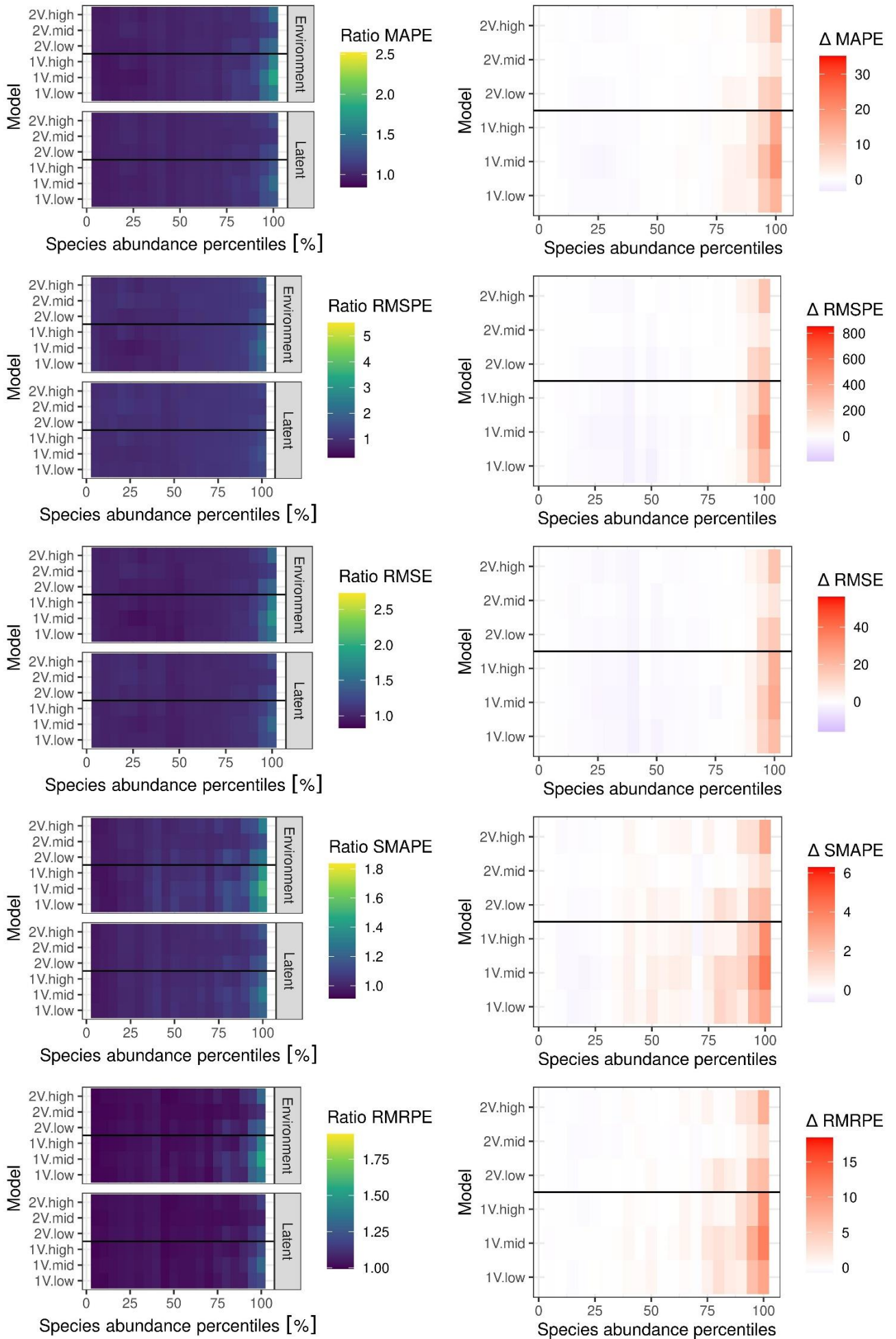
18



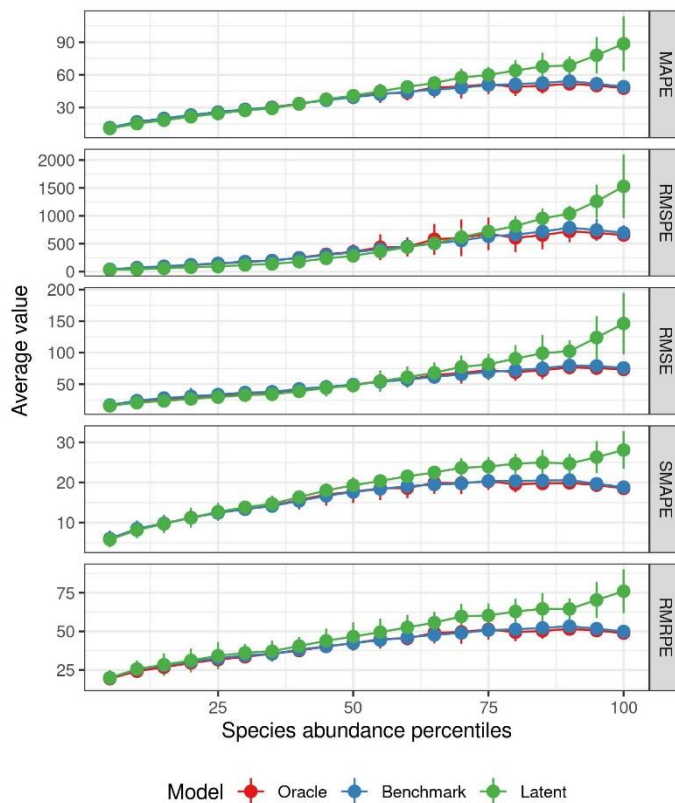
19

20 **Figure S3.** Average value of the studied metrics (Ratio TSS, ratio sensitivity, and ratio  
 21 specificity) depending on the number of sites used to fit the models, the model used, and the  
 22 occurrence of species. Horizontal panels represent the different occurrence: species with low,  
 23 medium and high occurrence corresponding respectively to bins of 15, 50, and 80 percentiles  
 24 of occurrence. Vertical panels indicate the metrics considered, with the models represented in  
 25 different colors. The ratio metric is calculated as the metric for the predictions of a model for  
 26 a species of the landscape divided by the same metric calculated for the oracle model. For the  
 27 ratio TSS, a score of 1 indicates a perfect agreement between the predictions of the considered  
 28 model and the oracle model, while a score of 0 or less represents a performance no better than  
 29 random. For the ratio sensitivity, it represents the ability to correctly classify a species as

30 “present”, while the ratio specificity represents the ability to correctly classify a species as  
31 “absent”. For both metrics, values above 1 indicate a better performance than the oracle  
32 model and values below 1 indicate a lesser performance. The benchmark model refers to the  
33 model containing all environmental variables, 2V.high the model with the two environmental  
34 variables with the highest coefficients, 1V.high the model with the environmental variable  
35 with the highest coefficient, and Latent the model containing the latent variables.

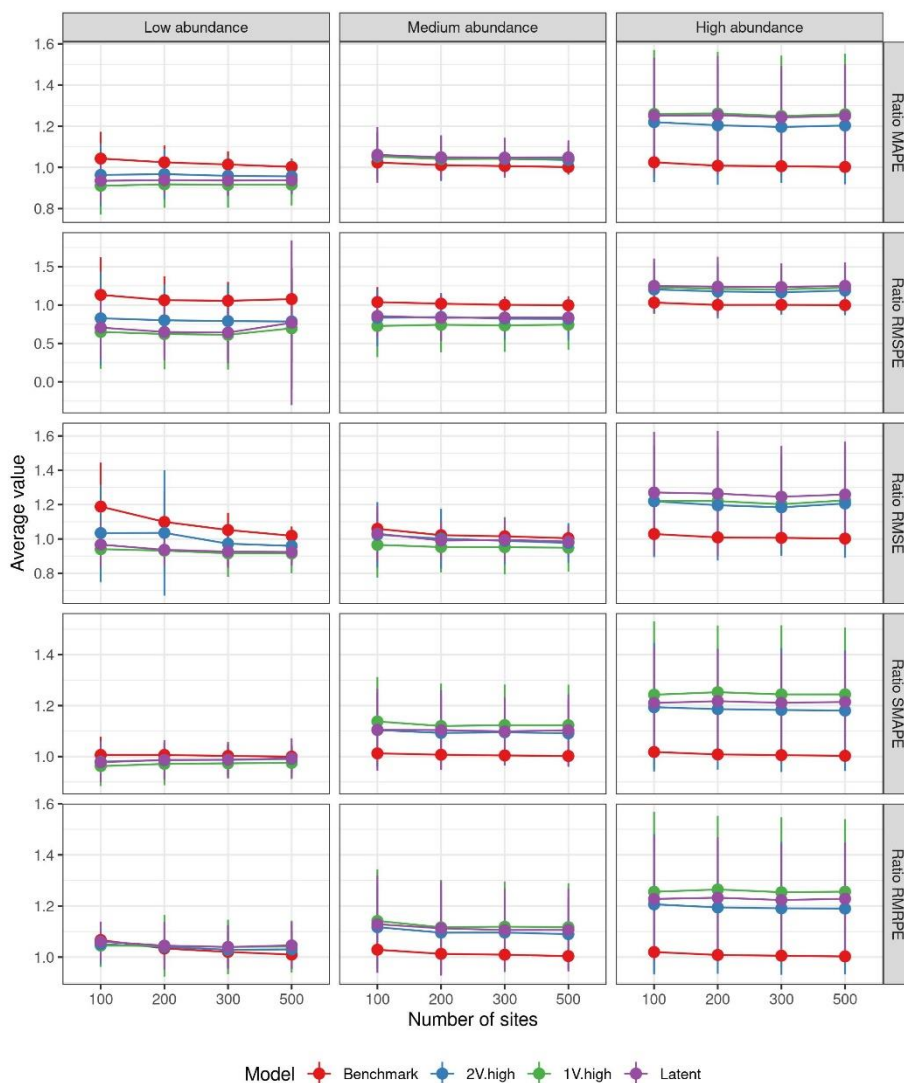


37 **Figure S4.** Abundance metrics and the comparison of performance between environmental  
 38 models and latent models measured as delta metrics. Each metric is averaged across all  
 39 landscapes and replicates per model and species, with the species binned by percentile of  
 40 abundance, and divided by the metric of the oracle model to give the ratio metric. The  
 41 environment panel represents models containing only environmental variables, while the  
 42 latent panel depicts models containing latent predictors. The models are then ordered from  
 43 bottom to top, from the fewest to the greatest number of environmental variables included and  
 44 sorted by coefficients relative to each environmental variable. See Methods for more  
 45 information, note that the “mid” model refers to the “intermediate” model. The delta metric  
 46 was measured as the metric of the model with environmental variables only minus the metric  
 47 of the model with the same combination of environmental and latent predictors. A positive  
 48 value indicates that the model with latent predicts the abundance of the species better than the  
 49 model containing only environmental variables.



50

51 **Figure S5.** Correlation between the metrics studied (MAPE, RMSPE, RMSE, SMAPE, and  
 52 RMRPE) depending on the model across species abundance percentiles. The vertical panels  
 53 indicate the different metrics, with models represented in different colors. Each metric is  
 54 averaged across all landscapes and replicates per model and species, with the species binned  
 55 by percentile of abundance. The oracle model refers to the model using the true environmental  
 56 coefficients while the other models were fitted using all environmental variables (benchmark)  
 57 or latent variables (latent).



58  
 59 **Figure S6.** Average value of the studied metrics (MAPE, RMSPE, RMSE, SMAPE, and  
 60 RMRPE) depending on the number of sites used to fit the models, the model used, and the



61 abundance of species. Horizontal panels represent the different abundances: species with low,  
62 medium and high occurrence corresponding respectively to bins of 15, 50, and 80 percentiles  
63 of occurrence. Vertical panels indicate the metrics considered, with the models represented in  
64 different colors. Each metric is averaged across all landscapes and replicates per model and  
65 species, with the species binned by percentile of abundance and divided by the metric of the  
66 oracle model to give the ratio metric. The benchmark model refers to the model containing all  
67 environmental variables, 2V.high the model with the two environmental variables with the  
68 highest coefficients, 1V.high the model with the environmental variable with the highest  
69 coefficient, and Latent the model containing the latent variables.

70

71