

# eDNAPlus: A unifying modelling framework for DNA-based biodiversity monitoring

Alex Diana<sup>1</sup>, Eleni Matechou<sup>1</sup>, Jim Griffin<sup>2</sup>, Douglas W. Yu<sup>3,4</sup>, Mingjie Luo<sup>5</sup>, Marie Tosa<sup>5</sup>, Alex Bush<sup>6</sup>, Richard Griffiths<sup>7</sup>

<sup>1</sup> School of Mathematics, Statistics and Actuarial Science, University of Kent, UK,

<sup>2</sup> Department of Statistical Science, University College London, UK,

<sup>3</sup> School of Biological Sciences, University of East Anglia, UK,

<sup>4</sup> Center for Excellence in Animal Evolution and Genetics & State Key Laboratory of Genetic Resources and Evolution &

Yunnan Key Laboratory of Biodiversity and Ecological Security of Gaoligong Mountain & Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming, China,

<sup>5</sup> Kunming College of Life Sciences, University of Chinese Academy of Sciences, China,

<sup>6</sup> Department of Fisheries, Wildlife, & Conservation Sciences, Oregon State University, Oregon USA,

<sup>7</sup> Lancaster Environment Centre, University of Lancaster, UK,

<sup>8</sup> Durrell Institute of Conservation and Ecology, University of Kent, UK

## Abstract

DNA-based biodiversity surveys, which involve collecting physical samples from survey sites and assaying them in the laboratory to detect species via their diagnostic DNA sequences, are increasingly being adopted for biodiversity monitoring and decision-making. The most commonly employed method, metabarcoding, combines PCR with high-throughput DNA sequencing to amplify and read ‘DNA barcode’ sequences, generating count data indicating the number of times each DNA barcode was read. However, DNA-based data are noisy and error-prone, with several sources of variation, and cannot alone estimate the species-specific amount of DNA present at a surveyed site (*DNA biomass*). In this paper, we present a unifying modelling framework for DNA-based survey data that allows estimation of *changes in DNA biomass within species, across sites* and their links to environmental covariates, whilst for the first time simultaneously accounting for key sources of variation, error and noise in the data-generating process, and for between-species and between-sites correlation. Bayesian inference is performed using MCMC with Laplace approximations. We describe a re-parameterisation scheme for crossed-effects models designed to improve mixing, and an adaptive approach for updating latent variables, which reduces computation time. Theoretical and simulation results are used to guide study design, including the level of replication at different survey stages and the use of quality control methods. Finally, we demonstrate our new framework on a dataset of Malaise-trap samples, quantifying the effects of elevation and distance-to-road on each species, and produce maps identifying areas of high biodiversity and species DNA biomass.

*Keywords: crossed-effects model, environmental DNA, joint species distribution modelling, observation error, occupancy modelling*

# 1 Introduction

Ecology is undergoing a technology revolution that is making it possible to rapidly generate species inventories via automated and high-throughput DNA sequencers and via electronic sensors, such as drones, satellites, camera traps, and acoustic recorders. These techniques can, if coupled with appropriate algorithms and databases, simultaneously identify large numbers of target species, including those that are cryptic, difficult-to-access, tiny, and low-abundance (Bush et al., 2017; Besson et al., 2022; Piper et al., 2019; Ley, 2022). So far, the most efficient method for generating species-resolution inventories is DNA-based surveys, which rely on reading DNA barcodes: short, standardized sections of the genome that can be compared to a reference library to enable taxonomic identifications without the need to examine organism morphologies (Ratnasingham and Hebert, 2007).

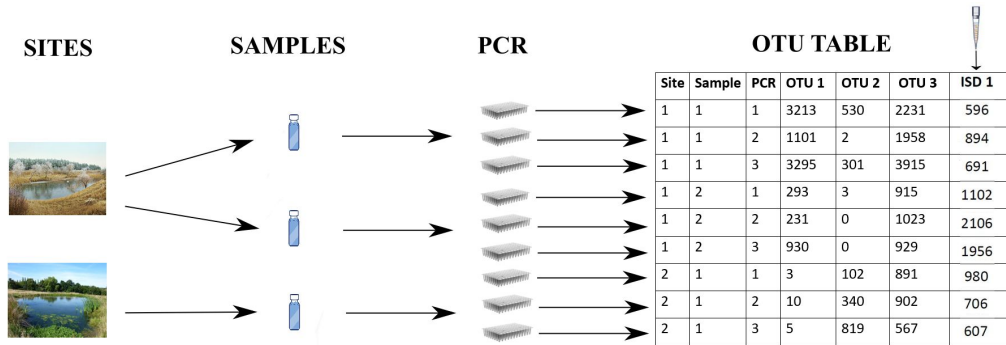
DNA barcoding refers to the identification of single species (Hebert et al., 2003), and DNA *metabarcoding* refers to the detection of large numbers of species from environmental DNA (eDNA), which is the collective name for DNA isolated from environmental samples (Taberlet et al., 2018). These environmental samples include water (Thomsen and Willerslev, 2015), soil (Frøslev et al., 2019), air (Clare et al., 2022), and bulk tissue (i.e. mass-trapped organisms) (Ji et al., 2013). For instance, Thomsen and Sigsgaard (2019) demonstrated that traces of eDNA on flower petals could be analysed to describe the diversity of arthropods that visit wildflowers, including pollinators, parasitoids, predators, and herbivores. Ji et al. (2022) used the trace amounts of residual vertebrate blood left in 30,468 blood-sucking leeches to map vertebrate wildlife across a 677 km<sup>2</sup> nature reserve in China. Finally, Abrego et al. (2021) sequenced 542 mixed-species, bulk-tissue samples of arctic arthropods captured over 14 years and showed that species richness in the study site had declined by 50% during a time period in which local mean temperature had increased by 2C.

The potential of DNA-based surveys for monitoring and managing biodiversity comes with a number of statistical challenges. Firstly, species-specific absolute abundances can-

66 not be estimated using DNA data alone. Secondly, DNA-based surveys yield data that  
67 are overdispersed (including zero-inflation) relative to a Poisson distribution due to several  
68 types of error and noise (see Section 1.1), some of which are species-specific. The framework  
69 presented in this paper addresses these challenges by developing a novel model and corre-  
70 sponding efficient inferential tools. Using our framework, we model *within-species change*  
71 *in DNA biomass across sites* (described in Section 1.1), which under certain conditions can  
72 be considered as a proxy for change in abundance, hence addressing the first challenge. To  
73 address the second challenge, we propose a hierarchical crossed-effects model that expresses  
74 key sources of variation, error and noise in the data collection and analysis pipeline, whilst  
75 accounting for correlation across species and across sites, and for covariate effects on DNA  
76 biomass. We also model frequently employed controls at the PCR stage and evaluate their  
77 effect on inference.

## 78 **1.1 DNA-based surveys and associated challenges**

79 Each individual of a species sheds tissue and waste products, and thus its DNA, into  
80 the environment. We will refer to this as *DNA biomass*. As we explain in Section 2,  
81 the estimates of species DNA biomass obtained from DNA-based surveys alone are only  
82 meaningful in comparison between sites, and for that reason, in this paper we focus on  
83 modelling *changes in DNA biomass within species, across sites*, referred to as changes in  
84 DNA biomass throughout. We achieve this by assuming that the processes are standardised  
85 across sites, samples, and PCR replicates and that any differences in the efficiencies of the  
86 processes are explained by covariates that can be included in the model. We highlight  
87 that, theoretically, the overall amount of DNA biomass for each species is proportional to  
88 the species' abundance at that site, but the rate at which each species sheds DNA into  
89 the environment is unknown and not estimable using eDNA data alone. Additionally, the  
90 relationship between DNA biomass and abundance can vary between species and sites due  
91 to environmental conditions, such as DNA degradation rates, and we return to this point  
92 in Section 6. Under the assumption that this relationship does not vary with sites then we



**Figure 1:** Representation of the DNA biomass collection stage (Stage 1, Sites to Samples) and the DNA biomass analysis stage (Stage 2, Samples to PCR to OTU table). Each of the selected sites to be surveyed hosts a community of species, and hence a certain amount of DNA biomass for each species. One or more physical samples are collected from each surveyed site, and a ‘spike-in’ or ‘internal standard’ ISD, can be added to each sample (last column). Each sample is PCR’d one or more times and then sequenced. This process gives rise to the OTU table.

93 can interpret changes in species DNA biomass as corresponding changes in abundance.

94 DNA-based surveys comprise two stages (Figure 1): the sample collection stage (Stage  
 95 1), taking place in the field, and the sample analysis stage (Stage 2), taking place in the  
 96 lab.

97 In Stage 1, physical samples are collected from each surveyed site. However, the amount  
 98 of DNA biomass of each species collected in each sample is the result of a noisy and error-  
 99 prone process (see Table 1). Specifically, the sampling method inevitably favours some  
 100 species over others, and as a result, DNA biomass collection rates, conditional on the  
 101 available DNA biomasses, are species-specific (*Stage 1 species effect*). The amount of DNA  
 102 biomass collected for each species also varies between samples collected at the same site  
 103 (*Stage 1 noise*). Finally, there are non-negligible probabilities that (a) no DNA biomass  
 104 is collected for a species even if there was DNA biomass of that species at the site (false  
 105 negative error) and (b) the DNA biomass in the sample is not the result of species presence,  
 106 but instead reflects contamination or deposition from elsewhere (false positive error) (Stage  
 107 1 false negative and false positive errors are jointly referred to as *Stage 1 error*).

108 In Stage 2, the physical samples are assayed in the lab. The most frequently used method  
 109 for reading DNA barcodes from eDNA samples is ‘amplicon sequencing’ (see Lindahl et al.,

110 2013, for an excellent review). In short, from each sample, all DNA is extracted and purified.  
111 After extraction, a small aliquot of DNA from each sample is subjected to Polymerase Chain  
112 Reaction (PCR), which selectively amplifies (makes many copies of) just the DNA-barcode  
113 sequences. It is common practice in Stage 2 for a sample to be PCR-assayed multiple times,  
114 known as technical replicates to distinguish them from sample replicates in Stage 1. The  
115 PCR outputs ('amplicons') from all the samples and their technical replicates are pooled  
116 and read on a high-throughput DNA sequencer. This procedure ultimately leads to a list  
117 of many millions of individual DNA sequences (known as reads), which are processed in a  
118 bioinformatic pipeline that removes low-quality reads, groups the remainder into clusters of  
119 similar reads that are species hypotheses known as OTUs (Operational Taxonomic Units),  
120 and apports each OTU's reads back to its original samples and PCRs. The resulting  
121 *OTU table* dataset indicates the number of reads for each OTU in each PCR in each sample  
122 in each site (Figure 1), with columns representing the species and rows representing the  
123 PCR runs. For simplicity, we hereafter use the terms OTUs and species interchangeably.

124 A real-world complication in DNA-based laboratory pipelines is that samples are typ-  
125 ically 'normalised' one or more times. For instance, after the samples are enzymatically  
126 digested to break down cells and release their DNA into their 'lysis-buffer' solutions, each  
127 sample constitutes a larger volume of liquid than can be used for DNA extraction. The  
128 samples are thus normalised by taking a fixed volume from each sample for processing.  
129 Another normalisation step happens after PCR, because different PCR replicates can gen-  
130 erate different amounts of product. In this case, the PCR products are normalised by  
131 taking a certain amount of liquid from each PCR output, either inversely proportional to  
132 their concentration, or fixed across PCRs. In the first (lysis buffer) normalisation step,  
133 the numerator (amount of lysis buffer taken for extraction) is fixed, while the denominator  
134 (total volume of lysis buffer) varies. In the second (PCR product) normalisation step, the  
135 numerator (amount of PCR liquid taken for sequencing) varies, while the denominator (to-  
136 tal volume of PCR liquid) is fixed. It is standard procedure to record these normalisation

137 fractions, and in Section 2, we show how this information is incorporated into the model.

138 Generally, we should expect a positive relationship between the DNA biomass of a  
139 species in a sample and the count of reads obtained for that species in that sample (Luo  
140 et al., 2022), but this relationship is imperfect, due to noise and error (see Table 1). First,  
141 even given best practice, there are small but non-negligible probabilities (a) that a species’  
142 DNA in a sample fails to be amplified or sequenced, leading to false-negative error and (b)  
143 that a species’ DNA cross-contaminates other samples and is amplified, leading to false-  
144 positive error (Stage 2 false negative and false positive errors are jointly referred to as *Stage*  
145 *2 error*). We say that a PCR yields non-negligible reads for a species when the PCR product  
146 of that species is successfully read by the DNA sequencer (i.e. the PCR is successful), and  
147 otherwise, a PCR yields zero or non-zero but negligible reads, in which case we say that  
148 the PCR is not successful for that species. We note that a PCR can be successful, that  
149 is, yield non-negligible reads, not only when the biomass is present in the sample but also  
150 when it is not, in the latter case because of contamination. Additionally, PCR amplification  
151 also inevitably favours some species over others, due to PCR primer mismatch, resulting  
152 in species-specific amplification rates (*Stage 2 species effect*, equal within columns of the  
153 OTU table), and PCR and sequencing stochasticity results in different total numbers of  
154 reads across all species, even for the same sample (*Stage 2 pipeline effect*, equal within rows  
155 of the OTU table). Finally, due to the inherent stochasticity of the PCR and sequencing  
156 process, there is added noise in the resulting reads *in each cell* of the OTU table (*Stage 2*  
157 *noise*).

158 In Stage 2, in addition to recording the normalisation fractions, different approaches  
159 are employed to understand and monitor some of the noise and error. One such approach  
160 is the so-called internal standard or *spike-in*, during which a known amount of DNA of a  
161 synthetic sequence or of a species that is known to be absent from all surveyed sites, is  
162 added to each sample. In addition, negative controls, which are samples that are known to  
163 not include DNA of any species, can be introduced in Stage 1 and Stage 2 (Ficetola et al.,

**Table 1:** Description of noise, error, and species/pipeline effects in the two stages of DNA-based surveys.

<b>Stage 1 - DNA biomass collection</b>	
<i>Species effect</i>	Every sample contains a certain amount of DNA biomass of each species, with the amount proportional to the DNA biomass available at the site. However, the proportionality constant is unknown and species-specific, since the DNA of different species can be collected at different rates.
<i>Noise</i>	The amount of DNA biomass collected for each species varies stochastically between samples collected at the same site and time.
<i>Error</i>	It is possible for the DNA of a target species that is present at a site not to be sampled (false negative error), or traces of DNA from one sample to contaminate another sample (false positive error).
<b>Stage 2 - DNA biomass analysis</b>	
<i>Species effect</i>	As a result of differences in gene copy number, DNA extraction efficiency, and PCR amplification efficiency, the correspondence between the source sample DNA biomass and the number of amplicon reads is species-specific (each column of the OTU table).
<i>Pipeline effect</i>	PCR stochasticity and the passing of small aliquots of liquid along the laboratory pipeline affects the total number of reads per technical replicate for all species (each row of the OTU table).
<i>Noise</i>	In addition to the species and pipeline effect, there is added noise in the number of reads per OTU and PCR (each cell of the OTU table).
<i>Error</i>	It is possible for the DNA of a target species that is present in the sample not to be amplified in the lab (false negative error), or traces of DNA of one sample to contaminate and be detected in other samples (false positive error), due to the high species-detection power of amplicon sequencing.

## 165 1.2 Existing approaches

166 A common approach for modelling metabarcoding data is to convert them to detection/non-  
167 detection data by thresholding the number of reads in the OTU table, with user-specified  
168 criteria. This allows the use of a generalized linear model (GLM) framework (Saine et al.,  
169 2020), which has also been extended to account for species correlation, for example using  
170 joint species distribution models (JSDMs) (Ovaskainen and Abrego, 2020). However, this  
171 approach does not account for the two stages or the noise and error inherent in DNA-based  
172 surveys (Table 1).

173 To that end, several different but related approaches have been proposed. A common  
174 approach applies occupancy models that account for false negative observation error to

175 the binary detection/no detection data (Ficetola et al., 2015). More recently, multi-scale  
176 extensions of these occupancy models have been proposed to account for false negative  
177 error in both stages (Mordecai et al., 2011; Schmidt et al., 2013) and for false positive  
178 error (Guillera-Arroita et al., 2017; Griffin et al., 2020) for a single species. However, the  
179 occupancy model framework disregards the information in the reads and relies on arbitrary  
180 thresholds about what constitutes a detection. Alternatively, the reads have also been  
181 modelled within a GLM framework (Takahara et al., 2012; Carraro et al., 2018) but without  
182 considering the errors in each stage. A joint model of species occupancy and corresponding  
183 reads was developed by Fukaya et al. (2022) but without considering the direct link between  
184 species DNA biomass at the site and species reads, or the correlation between species.

185 Finally, we note that an area of research similar to DNA-based biodiversity surveys  
186 is microbiome biology, which is the genetic material of all microbial life in an abiotic  
187 substrate (e.g. soil) or in a living host (e.g. the human microbiome). When modelling  
188 microbiome data, analysis has usually focused on understanding changes in the relative  
189 composition of each taxon across different samples. As a result, modelling approaches in  
190 this field have revolved around the Dirichlet-Multinomial, which allows inference of the  
191 changes, across samples, of the proportions of the species DNA biomasses (Fordyce et al.,  
192 2011; Coblenz et al., 2017; McLaren et al., 2019; Clausen and Willis, 2022), although  
193 within-species changes in DNA biomass are argued to be informative (Tkacz et al., 2018).  
194 A more detailed comparison between the model we introduce in this paper and models for  
195 microbiome data is given in Section 2.1.

### 196 **1.3 Structure of the paper**

197 In this paper, we present a unifying hierarchical modelling framework for OTU reads  
198 that considers key sources of variation, noise, and error at both stages of DNA-based  
199 biodiversity surveys (Table 1), while also modelling correlation between species and between  
200 sites. The model allows us to infer changes in DNA biomass and to link these changes to  
201 site-specific covariates.



202 We use state-of-the-art MCMC (Markov chain Monte Carlo) methods that build on  
203 recent work for hierarchical and crossed-effects models (Zanella and Roberts, 2021) as well  
204 as adaptive MCMC techniques (Andrieu and Thoms, 2008). In particular, we develop a  
205 novel sampling technique to improve mixing in the special case of a multivariate crossed-  
206 effect model with PCR-specific random effects, and we use adaptive updates of latent  
207 variables to focus sampling effort. This allows us to fit our model (with many latent  
208 variables across the different stages of DNA surveys) to data from large numbers of sites,  
209 samples per site, PCRs per sample, and species.

210 The new model, its properties, and links to existing models are presented in Section  
211 2. Details on our approach to inference are given in Section 3. Issues of study design are  
212 explored and corresponding simulations are presented in Section 4. A case study of a large  
213 Malaise-trap metabarcoding dataset is presented in Section 5, and the paper closes with a  
214 discussion in Section 6.

## 215 2 Model

216 We assume that  $M_i$  physical samples are collected from site  $i$ ,  $i = 1, \dots, n$ , and  $K_{im}$   
217 PCR replicates are performed on the  $m$ -th sample from site  $i$ . We denote by  $y_{imk}^s$  the  
218 number of DNA reads of the  $s$ -th species,  $s = 1, \dots, S$  in the  $k$ -th PCR replicate of the  
219  $m$ -th sample collected at the  $i$ -th site. We have  $n_z$  site covariates and  $X_i^z$  represents their  
220 value at site  $i$  and  $n_w$  sample covariates, represented as  $X_{im}^w$  for the  $m$  sample at the  $i$ -th  
221 site. In what follows,  $i$  indexes sites,  $m$  samples,  $k$  PCR replicates, and  $s$  species.

222 Our proposed model (see Figure 2) is hierarchical, with three levels. The first level  
223 models the amount of DNA biomass of each species at the surveyed sites, which is a  
224 function of environmental and landscape covariates as well as between-species and between-  
225 sites correlation (**DNA biomass availability**). The second level models the amount of  
226 DNA biomass collected for each species in each physical sample from each site (**DNA**  
227 **biomass collection**). Lastly, the third level models the number of reads obtained for  
228 each species in each PCR from each physical sample (**DNA biomass analysis**). Data are

229 observed only at the third level, as the result of Stage 2 of the survey, with levels one and  
 230 two corresponding to latent states.

**DNA biomass availability**  $L = \{l_i^s\} \sim \text{MN}(B_0 + X_z B, \Sigma, T), \quad T^{-1} \sim \text{GH}$

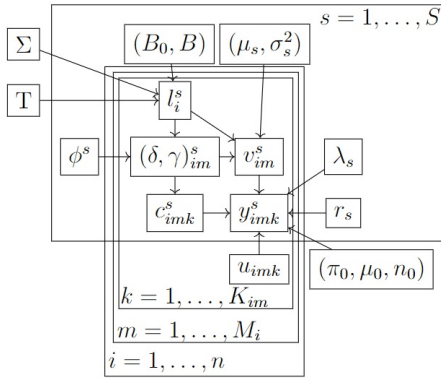
**DNA biomass collection**

$$\begin{aligned} \text{logit}(\theta_{im}^s) &= \phi_0^s + \phi_1^s l_i^s + X_{im}^w \phi^s \\ \mathbb{P}(\delta_{im}^s = 1) &= \theta_{im}^s, \\ \mathbb{P}(\gamma_{im}^s = 1 \mid \delta_{im}^s = 0) &= \zeta^s, \end{aligned} \quad v_{im}^s \sim \begin{cases} \text{N}(\eta_s + l_i^s + X_{im}^w \beta_s^W, \sigma_s^2) & \text{if } \delta_{im}^s = 1 \\ \text{N}(\mu_s, \nu_s^2) & \text{if } \delta_{im}^s = 0, \gamma_{im}^s = 1 \end{cases}$$

**DNA biomass analysis**

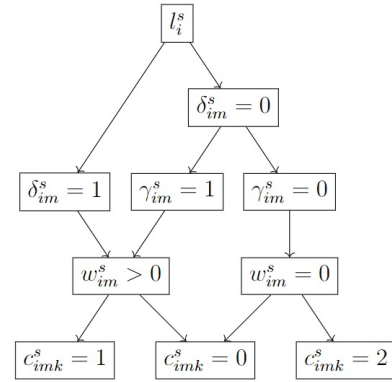
$\delta_{im}^s$	$\gamma_{im}^s$	$\mathbb{P}(c_{imk}^s = x \mid \delta_{im}^s, \gamma_{im}^s)$		
		$x = 0$	$x = 1$	$x = 2$
1	—	$1 - p_s$	$p_s$	0
0	1	$1 - p_s$	$p_s$	0
0	0	$1 - q_s$	0	$q_s$

$$y_{imk}^s \sim \begin{cases} \pi \delta_0 + (1 - \pi)(1 + \text{NB}(\mu_0, n_0)) & \text{if } c_{imk}^s = 0 \\ \begin{cases} \text{NB}(\exp(m_{imk}^s), r_s) \\ m_{imk}^s = \lambda_s + v_{im}^s + u_{imk} + o_{imk} \\ u_{imk} \sim \text{N}(0, \sigma_u^2) \end{cases} & \text{if } c_{imk}^s = 1 \\ \text{Pois}(\tilde{\mu}) & \text{if } c_{imk}^s = 2 \end{cases}$$



(b)

(a)



(c)

**Figure 2:** (a): Model summary, (b): Directed acyclic graph representing the relationships between the variables in the model. (c) Graphical representation of the latent indicator variables in the model.

231 **DNA biomass availability** We denote the logarithm of the amount of DNA biomass of  
 232 species  $s$  in site  $i$  available for collection by  $l_i^s$  and denote the  $n \times S$  matrix  $L$  by  $\{L\}_{is} = l_i^s$ .  
 233 We model DNA biomass correlation between species and spatial correlation between sites  
 234 by assuming that  $L$  follows a matrix normal distribution,  $L \sim \text{MN}(B_0 + X^z B, \Sigma, T)$  (Dawid,  
 235 1981), where  $B_0$  is an  $n \times S$  matrix with columns  $1_n \beta_0^s$ , with  $\beta_0^s$  a species-specific intercept,  
 236  $X^z$  is a design matrix whose rows are  $X_i^z$ ,  $B$  is an  $n_z \times S$  matrix of regression coefficients,  
 237  $\Sigma$  is an  $n \times n$  matrix modelling the correlation across sites, and  $T$  is an  $S \times S$  matrix  
 238 modelling the correlation across species. We note that, within this framework, the amount

239 of DNA biomass of a species at the surveyed site cannot be exactly 0, but can be negligible  
 240 for modelling purposes as we describe below. We employ a graphical horseshoe (GH) prior  
 241 (Li et al., 2019) for the inverse species covariance matrix  $Q = T^{-1}$ , which is defined by  
 242 specifying the following *a priori* independent distributions on each element

$$Q_{ss} \propto \text{Exp}\left(\frac{\lambda}{2}\right), s = 1, \dots, p, \quad Q_{ts} = Q_{st} \sim \text{N}(0, \lambda_{st}^2 \tau^2), \quad \lambda_{st} \sim C^+(0, 1), \quad s < t \leq S$$

243 subject to the constraint  $T \in \Omega_S$ , where  $\Omega_S$  is the space of the positive definite  $S \times S$   
 244 matrices,  $C^+$  represents the half-Cauchy distribution (Gelman, 2006), and  $\tau \sim C^+(0, 1)$ .  
 245 Unlike Li et al. (2019) who specified a flat prior  $Q_{ss} \propto 1$ , we follow Wang (2012) and  
 246 define a proper prior  $Q_{ss} \sim \text{Exp}(\frac{\lambda_{GH}}{2})$ , ensuring that  $T$ , which is latent, has a proper  
 247 posterior. We model the spatial correlation matrix  $\Sigma$  using an exponential kernel function,  
 248 so that  $\Sigma_{i_1 i_2} = \sigma^2 \exp\left\{-\frac{(x_{i_1} - x_{i_2})^2}{l^2}\right\}$ , where  $x_{i_1}$  and  $x_{i_2}$  are the locations of site  $i_1$  and  $i_2$ ,  
 249 respectively. We note that we have accounted for species correlations in the DNA biomass  
 250 availability stage, but any residual correlations of this type could also be the result of  
 251 species correlations in the collection or analysis stages, discussed below. It is not possible,  
 252 with metabarcoding data alone, to identify the source of these inferred correlations, and  
 253 therefore, species correlations should be interpreted with caution.

254 **DNA biomass collection** We denote by  $w_{im}^s$  the amount of DNA biomass of species  $s$   
 255 collected in sample  $m$  from site  $i$  and  $v_{im}^s := \log(w_{im}^s)$ . To account for *Stage 1 false negative*  
 256 *error* at this stage, we introduce the latent variable  $\delta_{im}^s$  that is equal to 1 if DNA biomass  
 257 for species  $i$  has been collected in the  $m$ -th physical sample from site  $i$ , and 0 otherwise.  
 258 We assume that  $\delta_{im}^s = 1$  with probability  $\theta_{im}^s$ , which is a function of covariates  $X_{im}^w$ , and  
 259 of  $l_i^s$ , since higher amounts of DNA biomass are expected to lead to a higher probability  
 260 of collecting DNA biomass in the sample, leading to  $\text{logit}(\theta_{im}^s) = \phi_0^s + \phi_1^s l_i^s + X_{im}^w \phi^s$ . We  
 261 note that as  $l_i^s$  tends to  $-\infty$ ,  $\theta_{im}^s$  tends to 0, and therefore the species becomes practically  
 262 impossible to detect. If the amount of DNA biomass collected is greater than 0 ( $\delta_{im}^s = 1$ ),  
 263 we model  $v_{im}^s \sim \text{N}(\eta_s + l_i^s + X_{im}^w \beta_s^w, \sigma_s^2)$ , where  $\eta_s$  models *Stage 1 species effects* on the  
 264 DNA biomass collection rate and  $\sigma_s^2$  models the species-specific *Stage 1 noise* in the DNA

265 biomass collection rate. To account for *Stage 1 false positive error*, we introduce latent  
 266 variable  $\gamma_{im}^s$ , which is equal to 1 with probability  $\zeta^s$  if the collected DNA biomass is the  
 267 result of contamination and 0 otherwise. We assume that  $\gamma_{im}^s$  can be 1 only if  $\delta_{im}^s = 0$   
 268 and that  $v_{im}^s \sim N(\mu_s, \nu_s^2)$  if  $\gamma_{im}^s = 1$ . In this way, we assume that a sample which already  
 269 contains DNA biomass of a species cannot be further contaminated by the DNA of the  
 270 same species from another sample or site. We make this assumption as there is not enough  
 271 information in the data to partition the collected DNA biomass between that which was  
 272 truly collected from the site and that which was contamination from elsewhere.

273 **DNA biomass analysis** As mentioned above, by non-negligible reads we mean that some  
 274 of the PCR product is successfully read by the DNA sequencer. We introduce latent variable  
 275  $c_{imk}^s$  to model the success of PCR  $k$ , sample  $m$ , and site  $i$  for species  $s$ , i.e. *Stage 2 error*.  
 276 Firstly, if sample  $m$  from site  $i$  contains DNA biomass of species  $s$  ( $w_{im}^s > 0$ ), PCR run  $k$   
 277 can be successful, i.e. non-negligible reads (true positive),  $c_{imk}^s = 1$ , or not successful, i.e.  
 278 negligible reads (false negative),  $c_{imk}^s = 0$ , and we assume that  $c_{imk}^s = 1$  with probability  
 279  $p_s$ . We note that we have assumed here that  $p_s$  only varies by species and not across sites  
 280 or replicates in either stage. However,  $p_s$  could depend (negatively) on the total amount  
 281 of DNA biomass in the sample, particularly in cases of low DNA concentration for that  
 282 species or could vary across primers or between labs. We return to these issues in Section  
 283 6. Secondly, if sample  $m$  from site  $i$  does not contain DNA biomass of species  $s$  ( $w_{im}^s = 0$ ),  
 284 PCR run  $k$  can be successful if it yields non-negligible reads due to lab contamination  
 285 (false positive),  $c_{imk}^s = 2$ , or not successful (again,  $c_{imk}^s = 0$ , true negative) and assume  
 286 that  $c_{imk}^s = 2$  with probability  $q_s$ .

287 We model the reads conditional on  $c_{imk}^s$  as follows. Conditional on  $c_{imk}^s = 1$ ,  $y_{imk}^s \sim$   
 288  $NB(\exp(\lambda_s + v_{im}^s + u_{imk} + o_{imk}), r_s)$ , where  $\lambda_s$  models the *Stage 2 species effect* on the  
 289 amplification rate,  $u_{imk}$  is the *Stage 2 pipeline effect*, with  $u_{imk} \sim N(0, \sigma_u^2)$ ,  $o_{imk}$  is an  
 290 offset modeling the normalisation steps described in Section 1.1, and  $r_s$  is a species-specific  
 291 variance of the *Stage 2 noise*. If more than one normalisation step is employed, then

292 they can all be incorporated into the same offset as a sum. Conditional on  $c_{imk}^s = 0$ ,  
 293  $y_{imk}^s \sim \pi\delta_0 + (1 - \pi)(1 + \text{NB}(\mu_0, n_0))$ , that is, there are zero reads with probability  $\pi$ , and  
 294 non-zero but negligible reads otherwise. Finally, conditional on  $c_{imk}^s = 2$ ,  $y_{imk}^s \sim \text{Pois}(\tilde{\mu}_s)$ .  
 295 The negative binomial is parameterised in terms of the mean and the number of failures.  
 296 A visual representation of the PCR process when  $c_{imk}^s = 1$  is shown in Figure 1 of the  
 297 Supplementary material.

298 Stage 2 negative control samples (which are known to not contain DNA of any species)  
 299 can be easily accounted for in our model by having additional samples for which  $\tilde{\delta}_i^s =$   
 300  $\tilde{\gamma}_i^s = 0$ . Accounting for spike-ins corresponds to having  $S^*$  additional species for which  
 301  $(v_{im}^{S+1}, \dots, v_{im}^{S+S^*})$  is known. Since the pipeline effect is shared across all species (including  
 302 spike-ins), the known values of  $v_{im}^s$  for the spike-ins help to better estimate  $u_{imk}$ . We further  
 303 investigate this effect in Section 4.

304 The model is summarised in Figure 2 (a), the directed acyclic graph of the model is  
 305 shown in Figure 2 (b), while a graphical representation of the latent variables introduced  
 306 across both stages is shown in Figure 2 (c). The model allows both zero-inflation and  
 307 overdispersion (even after accounting for zero-inflation) of the reads. In the case of true  
 308 positives (when  $c_{imk}^s = 1$ ), we allow overdispersion through the negative binomial distribu-  
 309 tion and the introduction of the offset. The use of negative binomial is a standard choice  
 310 for overdispersed data, particularly in Bayesian modelling. Ver Hoef and Boveng (2007)  
 311 discuss the merits of negative binomial and quasi-Poisson regression modelling in ecological  
 312 data. Datta and Dunson (2016) discuss how a scale mixture of negative-binomial regression  
 313 models can be used for so-called quasi-sparse counts, which are often small, not zero.

314 The model presented in Figure 2 is not identifiable in its general form unless certain  
 315 constraints are applied, as we discuss below. For example, choosing for simplicity  $\Sigma$  and  $T$   
 316 to be diagonal, if we define  $\tilde{v}_{im}^s := v_{im}^s - \eta_s - l_i^s$  and  $\tilde{l}_i^s := l_i^s - \beta_0^s$ , the model for  $\theta_{im}^s$  and  
 317  $y_{imk}^s$  conditional on  $c_{imk}^s = 1$  and all offsets  $o_{imk}$  set to 0 can be expressed as

$$\begin{cases} \tilde{l}_i^s \sim \text{N}(X_i \beta_s^z, \tau_s^2) \\ \tilde{v}_{im}^s \sim \text{N}(X_{im} \beta_s^w, \sigma_s^2) \\ \theta_{im}^s = \text{logit}(\phi_0^s + \phi_1^s \beta_0^s + \phi_1^s \tilde{l}_i^s + \phi^s X_{im}^s) \\ y_{imk}^s \sim \text{NB}\left(\exp(\beta_0^s + \tilde{l}_i^s + \eta_s + \tilde{v}_{im}^s + \lambda_s + u_{imk}), r_s\right) \end{cases} \quad (1)$$

It is evident that the model is invariant to transformations of the form

$$(\beta_0^s)^* = \beta_0^s + c + d, \quad (\lambda_s)^* = \lambda_s - c, \quad (\eta_s)^* = \eta_s - d, \quad (\phi_0^s)^* = \phi_0^s - \phi_1^s(c + d).$$

318 The reason for this unidentifiability is that data are observed only in the third level of  
 319 the model, and hence the following sets of species-specific parameters are confounded: the  
 320 baseline amount of DNA biomass across all sites ( $\beta_0^s$ ) with the baseline collection rate ( $\eta_s$ )  
 321 and the baseline amplification rate ( $\lambda_s$ ), and the former again with the baseline detection  
 322 rate  $\phi_0^s$ . However, by assuming that all these baseline rates are constant across sites,  
 323 samples, and PCRs, we are able to infer species-specific *changes* in DNA biomass across  
 324 sites and therefore covariate effects.

325 For inferential purposes, we reparameterise the model and set the new baseline (log)  
 326 amount of DNA biomass,  $(\beta_0^s)^*$ , equal to  $\beta_0^s + \eta_s$ , which means that we can only estimate  
 327 the sum of the baseline amount of available DNA biomass and the corresponding baseline  
 328 collection rate for the same species. Similarly, we set the new baseline (logit) collection  
 329 probability  $(\phi_0^s)^*$ , equal to  $\phi_0^s - \phi_1^s \eta_s$ , since the baseline collection probability is also con-  
 330 founded with the baseline collection rate (equivalent to setting  $\phi_0^s \equiv 0$  and  $\eta_s \equiv 0$  in  
 331 Equation (1)).

332 As a result, we cannot infer the amount of available DNA biomass separately from the  
 333 collection rate, and hence the estimates of log DNA biomass obtained, as mentioned above,  
 334 are only meaningful for comparison *within* each species. For the same reason, comparisons  
 335 of absolute amount of DNA biomass *across* species are not meaningful. We also note that  
 336 depending on the survey design in terms of the number of samples collected per site and  
 337 the number of PCR replicates per sample, additional sets of parameters can be confounded  
 338 and not estimable. Specifically, the following pairs of parameters are confounded:

- 339 •  $S = 1$ : pipeline effect  $u_{imk}$  and PCR variance  $r_s$ ,

- 340 •  $K = 1$ : PCR variance  $r_s$  and sample noise  $\tilde{v}_{im}^s$ ,
- 341 •  $M = 1$ : sample noise  $\tilde{v}_{im}^s$  and site noise  $\tilde{l}_i^s$ .

342 These are pathological cases that arise when there is no replication at the site/sample/PCR  
 343 levels. Replication is vital for being able to account for and to estimate the effects of the  
 344 different sources of noise and error (Buxton et al., 2021), an issue to which we return in  
 345 Section 4.1. Finally, we note that if the offsets  $o_{imk}$  introduced in the model due to the  
 346 several normalisations occurring in the pipeline are not recorded, the link between the  
 347 amount of DNA biomass in the environment and the reads is broken. However, a potential  
 348 way to restore this link is the introduction of spike-ins, which contribute to the estimation  
 349 of the “overall” pipeline effects  $\tilde{u}_{imk} = u_{imk} + o_{imk}$ .

## 350 2.1 Special cases

351 Two models available in the literature (Section 1.2) arise as special cases of our model.  
 352 First, the Dirichlet-Multinomial model (DMM) (Fordyce et al., 2011) is expressed through  
 353 the following hierarchy (omitting the indexes  $m$  and  $k$  to simplify notation):

$$\begin{cases} (y_i^1, \dots, y_i^S) \sim \text{Multi}(N_i, \pi_i^1, \dots, \pi_i^S) \\ (\pi_i^1, \dots, \pi_i^S) \sim \text{Dirichlet}(w\alpha^1, \dots, w\alpha^S) \end{cases} \quad (2)$$

354 where  $N_i = \sum_{s=1}^S y_i^s$ . The DMM can be seen as a special case of the model described in  
 355 Section 2, for the Stage 2 process, conditional on  $\delta_i^s = 1$ . Specifically,  $y_i^s \sim \text{NB}(\exp(\lambda_s +$   
 356  $v_i^s + u_i), r_s)$ , and therefore, assuming  $\lambda_s = u_i = 0$ , if  $r_s \rightarrow \infty$ , the distribution for  $y_i^s$  con-  
 357 verges to a  $\text{Pois}(\exp(v_i^s))$ . Conditional on  $N_i$ , the model is a  $\text{Multi}(N_i, \pi_i^1, \dots, \pi_i^S)$ , where  
 358  $(\pi_i^1, \dots, \pi_i^S) = \left( \frac{\exp(v_i^1)}{\sum_s \exp(v_i^s)}, \dots, \frac{\exp(v_i^S)}{\sum_s \exp(v_i^s)} \right)$ . Next, assuming  $\exp(v_i^s) \sim \text{Gamma}(w\alpha_s, \theta)$ , we  
 359 obtain  $(\pi_i^1, \dots, \pi_i^S) \sim \text{Dirichlet}(w\alpha_1, \dots, w\alpha_S)$ . Finally, as the DMM does not take errors  
 360 into account, the equivalence with our model can be obtained by setting  $p_s \equiv 1$ .

361 McLaren et al. (2019) propose to account for the Stage 2 species effect in the DMM  
 362 framework by modelling the probabilities  $(\pi_i^1, \dots, \pi_i^S)$  as  $\left( \frac{e^1 \tilde{\pi}_i^1}{\sum_s e^s \tilde{\pi}_i^s}, \dots, \frac{e^S \tilde{\pi}_i^S}{\sum_s e^s \tilde{\pi}_i^s} \right)$ , where  $e_s$   
 363 models the species-specific efficiencies, which in our model is achieved by using a species-  
 364 specific  $\lambda_s$ . The DMM can be extended hierarchically if nested treatments are considered

365 (Coblentz et al., 2017) by defining a nested prior  $(\alpha^1, \dots, \alpha^S) \sim \text{Dirichlet}(\alpha_0^1, \dots, \alpha_0^S)$  for  
 366 each level. In our model, this is achieved by a hierarchy of normal priors. This highlights  
 367 a key difference between the DMM approach and the approach we introduce in this paper,  
 368 since we model the propagation of the *absolute* amount of DNA biomass across the different  
 369 stages, while the DMM models the propagation of the *relative* amount of DNA biomass.

370 Secondly, the occupancy model of Griffin et al. (2020), in the simple case of no covariates,

$$\begin{cases} z_i \sim \text{Be}(\psi) \\ w_{im} \sim \text{Be}(z_i \xi_1 + (1 - z_i) \xi_0) \\ y_{imk} \sim \text{Be}(w_{im} p + (1 - w_{im}) q) \end{cases} \quad (3)$$

371 designed for (single-species) qPCR, can be seen as a special case of our model when the  
 372 information in the counts is not considered. Specifically, letting  $l_i$  be binary, with  $l_i \in$   
 373  $\{-\infty, 0\}$ , and defining  $z_i = \exp(l_i)$ , we obtain  $\theta_{im}|(l_i = -\infty) = 0$  and  $\theta_{im}|(l_i = 0) =$   
 374  $\text{logit}(\phi_0)$ . Hence, the model for  $\delta$  and  $c$  becomes

$$\begin{cases} \delta_{im} \sim \text{Be}(z_i(\text{logit}(\phi_0) + (1 - \text{logit}(\phi_0))\zeta) + (1 - z_i)\zeta) \\ c_{imk} \sim \text{Be}(\delta_{im} p + (1 - \delta_{im}) q) \end{cases},$$

375 which is identical to the Griffin et al. (2020) model after defining  $\xi_1 = \text{logit}(\phi_0) + (1 -$   
 376  $\text{logit}(\phi_0))\zeta$  and  $\xi_0 = \zeta$ .

### 377 **3 Inference**

378 Samples can be drawn from the posterior distribution of the parameters using a Gibbs  
 379 sampler. Posterior sampling is greatly helped by representing the negative binomial dis-  
 380 tribution as a Gamma-Poisson mixture, which allows many parameters to be updated in  
 381 closed form from their full conditional distribution.

382 For the parameters  $\sigma_s$ ,  $\mu_s$ ,  $B$  and  $B_0$ , the full conditional distribution is available in  
 383 closed form, and therefore posterior sampling is straightforward. We use simple random  
 384 walk Metropolis-Hastings steps for parameters  $\pi$ ,  $\mu_0$ ,  $n_0$ , and  $r_s$  and Metropolis-Hastings  
 385 steps with a Laplace approximation proposal for the parameters  $l_i^s$ ,  $\lambda_s$ ,  $v_{im}^s$ ,  $u_{imk}$  and  $r_s$ .  
 386 However, on its own, this naive Gibbs sampler will mix slowly since we have a complex  
 387 hierarchical model with crossed-effects and many latent variables. We address this by



388 updating parameters in blocks using re-parameterisation and an adaptive updating scheme  
 389 for the discrete latent variables.

390 To illustrate our approach to blocking and re-parameterisation, we consider the error-  
 391 free version of our model

$$\begin{cases} l_i^s \sim N(0, \tau_s^2) \\ v_{im}^s \sim N(l_i^s, \sigma_s^2) \\ u_{imk} \sim N(0, \sigma_u^2) \\ y_{imk}^s \sim \text{NB}(\exp(\lambda_s + v_{im}^s + u_{imk}), r_s) \end{cases} \quad (4)$$

392 A naive Gibbs sampler updating each parameter from its full conditional leads to pro-  
 393 hibitively slow mixing, due to the form of the likelihood where  $\lambda_s$ ,  $v_{im}^s$  and  $u_{imk}$  appear as  
 394 a sum. To address the slow mixing in the nested effects,  $\lambda_s$  and  $v_{im}^s$ , the use of a centred  
 395 parameterisation (Papaspiliopoulos et al., 2007) has been suggested, which corresponds to  
 396 defining  $\bar{v}_{im}^s := \lambda_s + v_{im}^s$  and  $\bar{l}_i^s := \lambda_s + l_i^s$ . However, issues of slow mixing still exist between  
 397  $\bar{v}_{im}^s$  and  $u_{imk}$  and, as noted by Zanella and Roberts (2021), re-parameterisation does not  
 398 improve mixing in the case of crossed-effects models. In a classic crossed-effect model of the  
 399 form  $y_{jkl} \sim N(\lambda + v_j + u_k, \sigma^2)$ , Papaspiliopoulos et al. (2020) propose a collapsed Gibbs sam-  
 400 pler by first jointly sampling  $\lambda$  with  $v_j$  and then  $\lambda$  jointly with  $u_k$ . However, this approach  
 401 does not scale well in our setup, since it would involve sampling all the  $\lambda_s$  and  $u_{imk}$  jointly,  
 402 which have dimensions  $S$  and the total number of PCR technical replicates  $\sum_{i,m} K_{im}$  re-  
 403 spectively. Zanella and Roberts (2021) propose the use of identifiability constraints on  
 404 the model, which in Equation (4) correspond to assuming  $\sum_s v_{im}^s = \sum_k u_{imk} = 0$ . Since  
 405 sampling conditionally on constraints can be challenging, we propose a simpler strategy to  
 406 improve mixing that is more suited to our framework. We consider re-parameterising to  
 407 the *factor averages*  $\hat{v}_{im} = \frac{1}{S} \sum_{s=1}^S \bar{v}_{im}^s$  and  $\hat{u}_{im} = \frac{1}{K} \sum_{k=1}^K u_{imk}$  and the *factor increments*  
 408  $\tilde{v}_{im}^s = \bar{v}_{im}^s - \hat{v}_{im}$  and  $\tilde{u}_{imk} = u_{imk} - \hat{u}_{im}$  and performing an update by first sampling jointly  
 409 the factor means conditional on the increments, that is, from  $(\hat{v}_{im}, \hat{u}_{im} | \tilde{v}_{im}^s, \tilde{u}_{imk})$  and next  
 410 using the standard updates  $(u_{imk} | v_{im}^1, \dots, v_{im}^S)$  and  $(v_{im}^j | u_{im1}, \dots, u_{imK})$ . In our simula-  
 411 tions, we have found that jointly updating the factor means considerably improves mixing.

412 The sampling scheme for the complete model is presented in the Supplementary material.

413 The indicator variables  $(\delta_{im}^s, \gamma_{im}^s, c_{imk}^s)$  can be updated directly from their full condi-  
414 tional distributions but, since there are  $nMS(\bar{K} + 2)$  (where  $\bar{K}$  is the average number of  
415 PCR replicates) of these variables and often one value of  $(\delta_{imk}, \gamma_{imk}, c_{imk})$  has probability  
416 very close to 1, evaluating every full conditional distribution in every iteration can be very  
417 time-consuming and computationally wasteful. Therefore, we use a cheap approximation as  
418 a proposal in a Metropolis-Hastings step. Specifically, every  $B$  iterations, we update the ap-  
419 proximation  $\hat{p}((\delta_{im}^s, \gamma_{im}^s, c_{imk}^s) = (\epsilon_1, \epsilon_2, \epsilon_3)) = \frac{1}{T} \sum_{t=1}^T \mathbf{I}((\delta_{im}^s)^{(t)}, (\gamma_{im}^s)^{(t)}, (c_{imk}^s)^{(t)}) = (\epsilon_1, \epsilon_2, \epsilon_3))$ ,  
420 where  $(\delta_{im}^s)^{(t)}, (\gamma_{im}^s)^{(t)}, (c_{imk}^s)^{(t)}$  is the value of  $(\delta_{im}^s, \gamma_{im}^s, c_{imk}^s)$  at the  $t$ -th iteration,  $\mathbf{I}(A)$  is  
421 the indicator function, which takes the value 1 if  $A$  is true and 0 otherwise, and  $T$  is the  
422 number of current iterations. Using this update scheme, we only need to evaluate the  
423 likelihood if the state is proposed to change. If the probability of one state is close to  
424 one, the adaptive scheme often proposes the current state, which can be accepted without  
425 computation. The adaptive scheme does not affect convergence of the MCMC algorithm  
426 since the approximation clearly has diminishing adaptation, and the state space of the  
427 indicator variables is discrete (see *e.g.* Roberts and Rosenthal, 2009, for more discussion of  
428 conditions for convergence of adaptive MCMC schemes).

## 429 4 Study design

430 In this section, we use a simplified version of the model to investigate the properties  
431 of our modelling approach under different study designs in terms of the number of sites,  
432 samples per site, and PCRs per sample, as well as the number of spike-ins. In each section,  
433 we consider the estimates of the differences in log DNA biomass, when log DNA biomass  
434 is not a function of site-specific covariates (no covariate case), and the estimates of the  
435 covariate coefficients when log DNA biomass is a function of a single continuous covariate  
436 (regression case). In Section 4.1 we present theoretical results using a continuous version  
437 of our model that does not account for error in either stage. In Section 4.2 we fit our model  
438 as presented in Section 2 under different scenarios for study design by varying the number

439 of sites, number of samples per site, and number of PCRs per sample. Finally, in Section  
 440 4.3, we explore the effect of spike-ins for different levels of noise in each stage of the process  
 441 and different study designs.

## 442 4.1 Theoretical results for a simplified version of the model

443 We consider a normal approximation of the model presented in Section 2, which assumes  
 444 no species or site correlations, that both stages are error-free by setting  $\theta_{im}^s = p_s = 1$ , and  
 445 that the variances of the distributions of the noise at each stage are the same across species.  
 446 As mentioned in Section 2, the use of spike-ins corresponds to the presence of species in  
 447 the sample for which  $(v_{im}^{S+1}, \dots, v_{im}^{S+S^*})$  is known. We assume, without loss of generality,  
 448 that  $v_{im}^{S+j} = 0$  for  $j = 1, \dots, S^*$ . We have the following proposition.

**Proposition 4.1.** *Consider the model  $\lambda_s \sim N(0, \sigma_\lambda^2)$  for  $s = 1, \dots, S + S^*$  and, for  
 $i = 1, \dots, n$ ,  $k = 1, \dots, K$ ,  $m = 1, \dots, M$ ,*

$$u_{imk} \sim N(0, \sigma_u^2), \quad v_{im}^s \begin{cases} \sim N(l_i^s, \sigma^2), & s = 1, \dots, S \\ = 0, & s = S + 1, \dots, S + S^* \end{cases},$$

$$y_{imk}^s \sim N(u_{imk} + \lambda_s + v_{im}^s, \sigma_y^2), \quad s = 1, \dots, S + S^*$$

449 where  $\sigma^2$ ,  $\sigma_u^2$  and  $\sigma_y^2$  are known.

(a) *If we assume  $p(l_i^s) \propto 1$  and  $\sigma_\lambda^2 \in (0, \infty)$  is known, then*

$$\text{Var}(l_1^s - l_2^s | y) = \frac{1}{M} \left( \sigma^2 + \frac{\sigma_y^2}{K} \left( 1 + \frac{\frac{\sigma_y^2}{\sigma_u^2}}{\frac{\sigma_y^2}{\sigma_u^2} S^* + 1} \right) \right). \quad (5)$$

450 (b) *If we observe a single covariate  $X_i \stackrel{i.i.d.}{\sim} N(0, 1)$  for the  $i$ -th site and assume  $l_i^s \sim$   
 $N(X_i \beta_s, \tau^2)$  with  $\sigma_\lambda^2 = \infty$  (i.e.  $p(\lambda_s) \propto 1$ ) and  $p(\beta_s) \propto 1$ , then*

$$\text{Var}(\beta_s | y) = \frac{1}{n-1} \left( \tau^2 + \frac{1}{M} \left( \sigma^2 + \frac{\sigma_y^2}{K} \right) \right) \times (1 + C) \quad (6)$$

451 where  $C = \frac{\sigma_u^2}{\sigma_y^2 + (M\tau^2 + \sigma^2)K(1 + S^* \frac{\sigma_y^2}{\sigma_u^2}) + \sigma_u^2(S + S^* - 1)}$ .

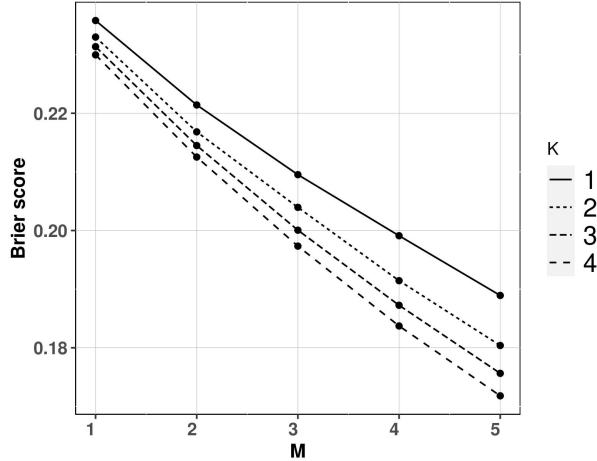
452 Here  $\sigma_y^2$  models the variance of the noise in Stage 2, as was the case for  $r_s$  in the original  
 453 model. Equations (5) and (6) show the contributions of the variances at each stage to the  
 454 posterior variance of the corresponding estimates (changes in biomass between sites, on the  
 455 log scale, and covariate coefficients, respectively) in this special case.

456 The results for this special case suggest that, for both  $\text{Var}(l_1^s - l_2^s|y)$  and  $\text{Var}(\beta|y)$ ,  
 457 increasing replication at a given stage decreases the contribution of the error variance at  
 458 that stage and all downstream stages. For example, increasing the number of samples  $M$   
 459 per site reduces the contribution of the noise variance  $\sigma^2$  at Stage 1 and at all downstream  
 460 stages, i.e.  $\sigma_u^2$  and  $\sigma_v^2$  in Stage 2. Whereas, increasing the number of PCR replicates,  $K$ ,  
 461 only reduces the contribution of the Stage 2 variances ( $\sigma_u^2$  and  $\sigma_v^2$ ). Additionally, the benefit  
 462 of the spike-in is greater as the ratio of variances  $\frac{\sigma_u^2}{\sigma_v^2}$  increases. Moreover, in the case of  
 463  $\text{Var}(\beta|y)$ , if  $\sigma^2$  is much greater than  $\sigma_y^2$ , the benefit of the spike-in is negligible, as the noise  
 464 induced by  $\sigma^2$  greatly outweighs the noise that can be mitigated via the use of spike-ins.

## 465 4.2 Simulated results for the full model; varying $n$ , $M$ , and $K$

466 We turn our attention to the full model in Fig. 2 and again consider two cases: no  
 467 covariates and a single covariate,  $X_i \sim N(0, 1)$ . In the no covariate case, we consider the  
 468 model's ability to estimate the correct sign of the difference of species log DNA biomasses  
 469 at two sites. We use the Brier score  $b(i_1, i_2, s) := (\bar{p}(l_{i_1}^s > l_{i_2}^s) - \delta_{i_1, i_2})^2$ , where  $\bar{p}(l_{i_1}^s > l_{i_2}^s)$  is  
 470 the posterior probability of  $l_{i_1}^s > l_{i_2}^s$  and  $\delta_{i_1, i_2}$  is 1 if the true value of  $l_{i_1}^s$  is greater than the  
 471 true value of  $l_{i_2}^s$  and 0 otherwise. We generate  $l_i^s \sim \begin{cases} N(1, \tau_s^2) & i \text{ odd} \\ N(0, \tau_s^2) & i \text{ even} \end{cases}$  which separates  
 472 the sites between those with "high" DNA biomass and those with "low" DNA biomass. We  
 473 use  $S = 40$  species,  $n = 300$  sites,  $M \in \{1, 2, 3, 4, 5\}$  samples per site and  $K \in \{1, 2, 3, 4\}$   
 474 PCR replicates. The values of the other parameters are reported in the Supplementary  
 475 Material. We have performed 50 replications for each combination of values of the design  
 476 parameters,  $M$  and  $K$ . We report the average  $b(i_1, i_2, j)$  spanning  $i_1$  across the sites with  
 477 low DNA biomass,  $i_2$  across the sites with high DNA biomass, and  $s$  across all species

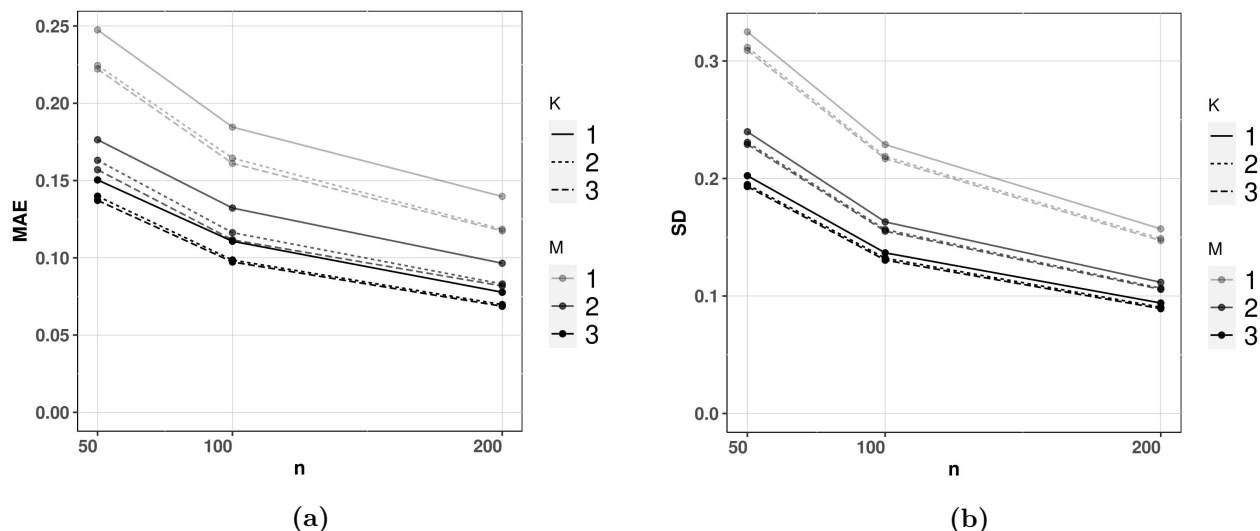
478 and across the replicates. As expected, the Brier score decreases, and hence the ability to  
 479 distinguish between sites with low and high DNA biomass increases, as  $M$  and  $K$  increase  
 480 (Figure 3). However, the benefit of increasing  $K$  decreases with  $M$ , which highlights the  
 481 greater importance of multiple sample replicates per site in Stage 1.



**Figure 3:** Brier score for distinguishing high and low DNA biomass sites, as a function of the number of samples ( $M$ ) and number of PCR replicates ( $K$ ). We have only considered  $M \leq 5$ , since greater  $M$  is unrealistic, and set  $n = 300$ .

482 In the regression case, we consider the absolute error and posterior standard deviation  
 483 of the covariate coefficient  $\beta_s$ . We use  $n \in \{50, 100, 200\}$  sites,  $M \in \{1, 2, 3\}$  samples per  
 484 site and  $K \in \{1, 2, 3\}$  PCR replicates per sample and  $S = 40$  species. The values of the  
 485 other parameters are reported in the Supplementary Material. We performed 50 replicates  
 486 for each combination of values of the design parameters and averaged results across all  
 487 replicates and species. Results are shown in Figure 4.

488 As expected, absolute error and posterior standard error both decrease with more sites  
 489  $n$ , samples per site  $M$ , and PCRs per sample  $K$ . Doubling the number of sites from 50  
 490 to 100 has a bigger effect than doubling them again from 100 to 200, suggesting that the  
 491 benefit of increasing the number of sampled sites decreases as the number of sites gets large.  
 492 Collecting two samples per site instead of one drastically decreases both absolute error and  
 493 posterior standard deviation, whereas the effect is less pronounced when the number of  
 494 samples is further increased to three compared to two, and the same can be said about the  
 495 number of PCRs.



**Figure 4:** Mean absolute error, (a), and posterior standard deviation, (b), averaged across all species and all simulations, of the covariate coefficient  $\beta^s$  for varying numbers of sites ( $n$ ), samples per site ( $M$ ), and numbers of PCR replicates per sample ( $K$ ).

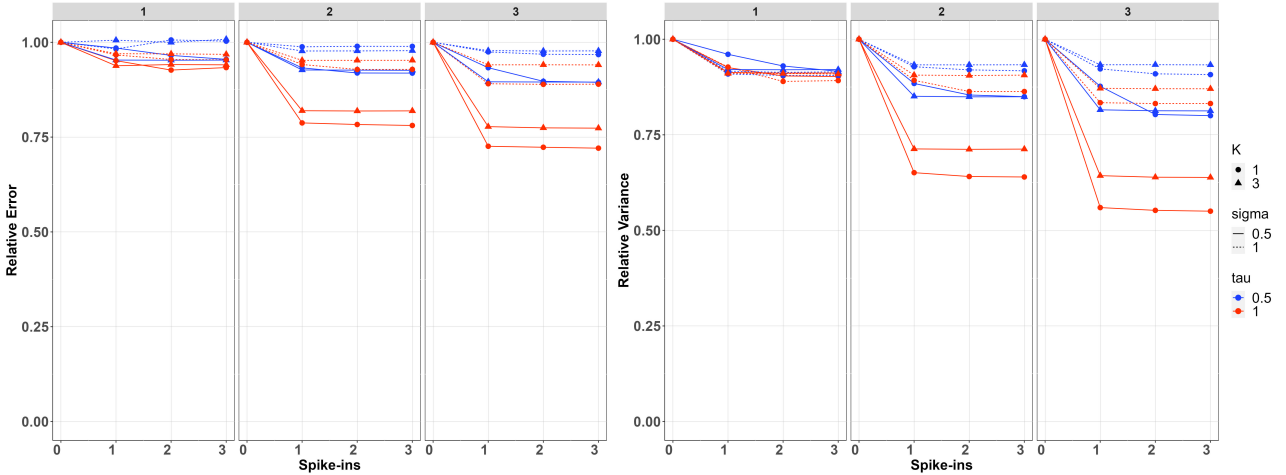
### 4.3 Spike-ins

In this section, we consider the improvement in inference when  $S^*$  spike-ins are employed in Stage 2. The effect of the spike-ins is maximised in the case of no false negative/positive errors, otherwise the benefit of the spike-ins is lower, and dependent upon the level of error. Therefore, in this section we consider data and corresponding model with no false positive/negative errors.

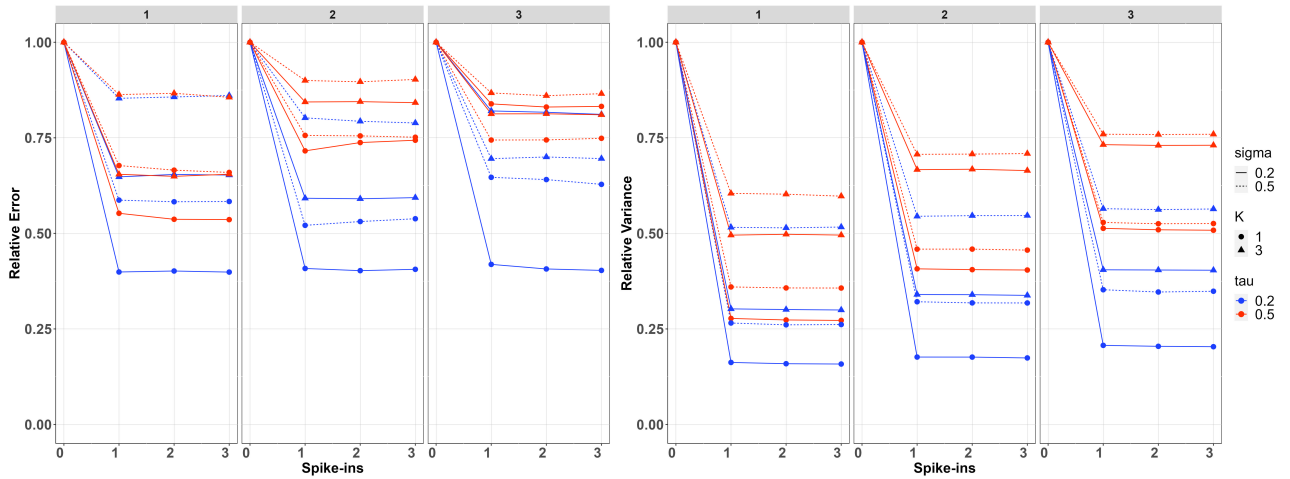
We simulated data on  $n = 300$  sites,  $M \in \{1, 2, 3\}$  samples per site, and  $K \in \{1, 3\}$  PCR replicates per sample on  $S = 10$  species. For each setting of  $M$  and  $K$ , we have fitted the model when  $S^* \in \{0, 1, 2, 3\}$  and report in each case the posterior relative error and posterior relative variance of the estimates, which are calculated by dividing the posterior error/variance by the corresponding error/variance when using  $S^* = 0$  (which is the case with the greatest error/variance).

Results of the simulation study are presented in Figure 5. In both cases, improvements diminish for  $S^* \geq 2$ , and in most cases  $S^* = 1$  already provides most of the improvement, suggesting that the benefit of more than one spike-ins is minimal. The no covariate case is shown in the first row of Figure 5. Spike-ins contribute more to reducing biomass-change estimation error and variance with  $M > 1$ , with  $M = 1$  resulting in virtually no

### No covariate



### Regression



**Figure 5:** Effect of spike-ins on inference. The three facets per figure represent simulations with  $M = 1/2/3$  samples per site. The between-samples standard deviation,  $\sigma$ , is represented by the line type, the between-sites standard deviation,  $\tau$ , is represented by the color, the number of PCR replicates,  $K$ , is represented by the symbols. The first column represents the posterior relative error of the estimates and the second column represents the posterior relative variance.

513 improvements for any setting considered in the simulation. When  $M > 1$ , improvement  
 514 is more pronounced when  $K = 1$  instead of  $K = 3$ , because in the latter case, thanks  
 515 to this replication at Stage 2, there is already increased information for estimating the  
 516 pipeline effect. This is particularly true when  $\tau$  is 1 instead of 0.5, because in this case,  
 517 the differences between sites are more pronounced. For both values of  $\tau$ , improvements are  
 518 bigger when the between-samples standard deviation ( $\sigma$ ) is smaller, since otherwise, Stage  
 519 1 noise dominates the process and understanding noise in Stage 2 decreases the overall

520 variance proportionally less.

521 The second row of Figure 5 shows the regression case. We have chosen smaller values  
522 for  $\sigma$  and  $\tau$  (.2 and .5), since the relative contribution of the spike-ins is negligible with  
523 larger values. Spike-ins contribute more to reducing error and variance when the between-  
524 samples standard deviation ( $\sigma$ ) and the between-sites standard deviation ( $\tau$ ) is smaller  
525 because, similar to before, the noise at early stages dominates the process, and therefore  
526 the relative contribution of the spike-ins is smaller. Also similar to the no covariate case,  
527 the contribution of the spike-ins is higher for  $K = 1$  PCR replicates compared to  $K = 3$ .  
528 However, unlike that case, the contribution does not appear to increase as the number of  
529 samples per site  $M$  increases.

## 530 5 Case study

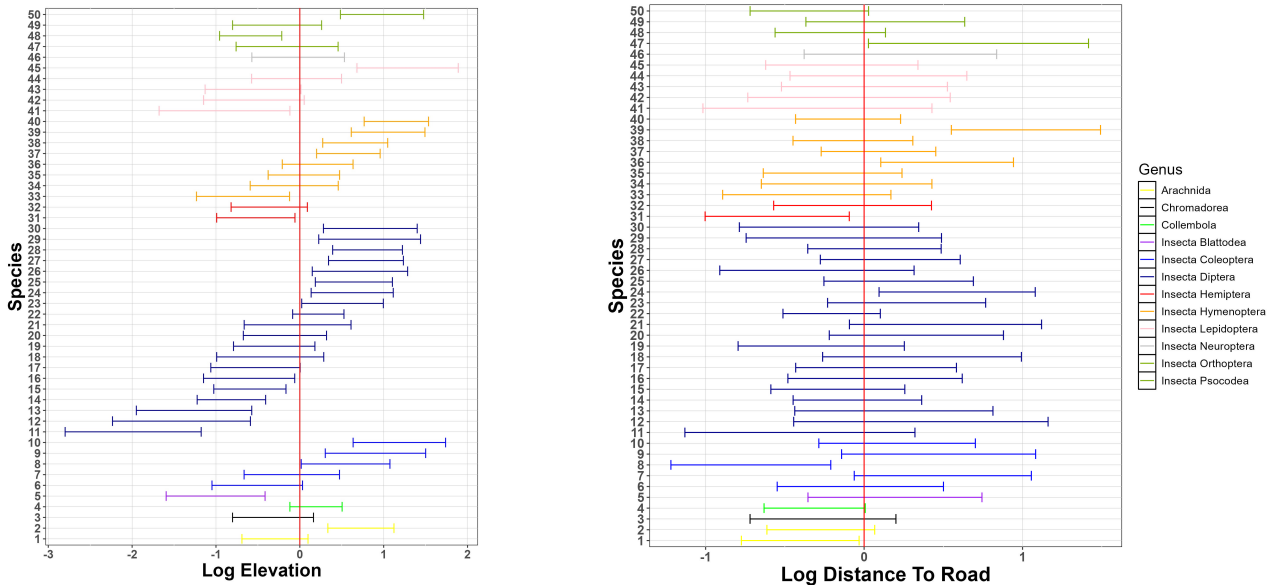
531 We apply our model to an unpublished amplicon sequencing dataset of arthropod inver-  
532 tebrates collected using 121 Malaise-trap samples from 89 sample sites in the H.J. Andrews  
533 Experimental Forest (HJA), Oregon, USA (225 km<sup>2</sup>) in July 2018 (site details are provided  
534 in Li et al. 2024). Each trap was left to collect for seven days, and samples were transferred  
535 to fresh 100% ethanol to store at room temperature until extraction. The management ob-  
536 jective that motivated the collection of this dataset is to interpolate continuous species  
537 distributions among the 89 sample points so that areas of higher and lower conservation  
538 value at the HJA can be identified.

539 For each sample, the collected invertebrate samples were combined with a lysis buffer,  
540 in an amount proportional to the starting sample mass, to digest the tissue, and a fixed  
541 aliquot was then taken from the overall mixture (and recorded) for DNA extraction and  
542 subsequent three PCRs. This normalization, as described in Section 2, was accounted for  
543 in the model by setting the offset  $o_{imk}$  equal to the log ratio between the aliquot and the  
544 overall amount of liquid mixture in each case. We included 50 species in the study by  
545 selecting the species that have the most non-zero counts across all PCR replicates. Log  
546 DNA biomass is modelled as a function of two environmental covariates: log elevation and



547 log distance-to-road.

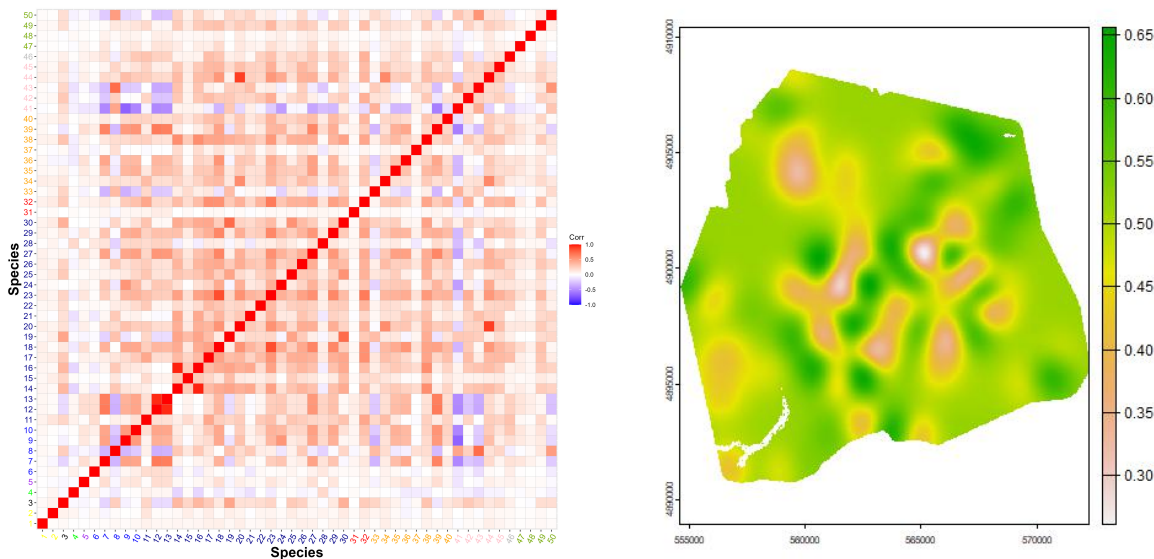
548 Figure 6 presents the 95% posterior credible intervals (PCIs) for the species-specific  
549 coefficients of log elevation and log distance-to-road in the model for log DNA biomass. The  
550 effects of the covariates on species DNA biomass are not consistent within each taxonomic  
551 order, which suggests low phylogenetic inertia at this rank for response to these landscape  
552 characteristics. Elevation is a stronger predictor for species DNA biomass than distance-to-  
553 road for this ecosystem. This makes ecological sense, since distance-to-road is only expected  
554 to exert an effect over about 100 meters, via canopy openness, whereas elevation exerts a  
555 pervasive effect via its effects on temperature, precipitation, and vegetation.



**Figure 6:** Case study: 95% PCI of the species-specific coefficients of log elevation (left) and distance to road (right) in the model for log DNA biomass. Species are grouped taxonomically.

556 Figure 7 (a) presents the posterior mean of the between-species residual correlations.  
557 We set  $\lambda_{GH} = 1$  in the GH prior and we emphasize that the GH prior assumes no prior  
558 structure imposed on the taxa. Species in the Diptera (flies, spp. 14-30) exhibit higher pos-  
559 itive correlations with each other, as well as with several species in the Hymenoptera (ants,  
560 bees, and wasps) and Lepidoptera (butterflies and months). We conservatively interpret  
561 these positive residual correlations as indicative of unmeasured environmental covariates,  
562 such as canopy openness, rather than of biotic interactions. We also note that two species in

563 the Lepidoptera, (spp. 41, 43), one in the Hymenoptera (sp. 33), and one in the Psocodea  
 564 (barklice, sp. 50) are among the few species showing strong negative residual correlation  
 565 with many of the other species, and again, we conservatively interpret these correlations  
 566 as indicative of unmeasured environmental covariates. There is a strongly positive, pair-  
 567 wise correlation between two tabanid fly species *Hybomitra liorhina* and *Hybomitra* sp.  
 568 (spp. 12, 13), which might indicate the oversplitting of one biological species into two  
 569 OTUs during the bioinformatic pipeline. Finally, there is also a strongly positive, pairwise  
 570 correlation between the moth species *Ceratodelia gueneata* (sp. 44) and the predatory fly  
 571 (Scathophagidae, *Microprosopa* sp., (sp. 20), which might indeed indicate a specialised  
 572 predator-prey relationship. All that said, we highlight that these inferred correlations have  
 573 been accounted for in the DNA availability stage of the model, but, as we discuss in Section  
 574 2, they can also be the result of the DNA biomass collection or analysis stages, so should  
 575 be interpreted with caution.



**Figure 7:** Case study. Left: Correlation plot of all species. Red represents positive correlations while blue represents negative correlations. Species are grouped taxonomically. Right: Posterior mean of biomass-weighted species richness across the study area. For each species, we rescale the log-biomass amount across all study sites into the range  $[0, 1]$  and next we compute the species richness as the sum of all the rescaled biomasses across all species.

576 In Figure 7 (b), we show the biodiversity map for the area, which is useful for identifying  
 577 areas of higher species richness and compositional distinctiveness, which together can be

578 used to identify areas of higher conservation value (i.e. higher ‘site irreplaceability’ *sensu*  
579 Baisero et al., 2022). The predicted mean log DNA biomasses on a continuous map over  
580 the HJA for all individual species are presented in the Supplementary Material. These can  
581 be used to identify species with a wide spatial range, such as the click beetle (*Megapenthes*  
582 *caprella*), or with a restricted range, such as the leafhopper (*Osbornellus borealis*).

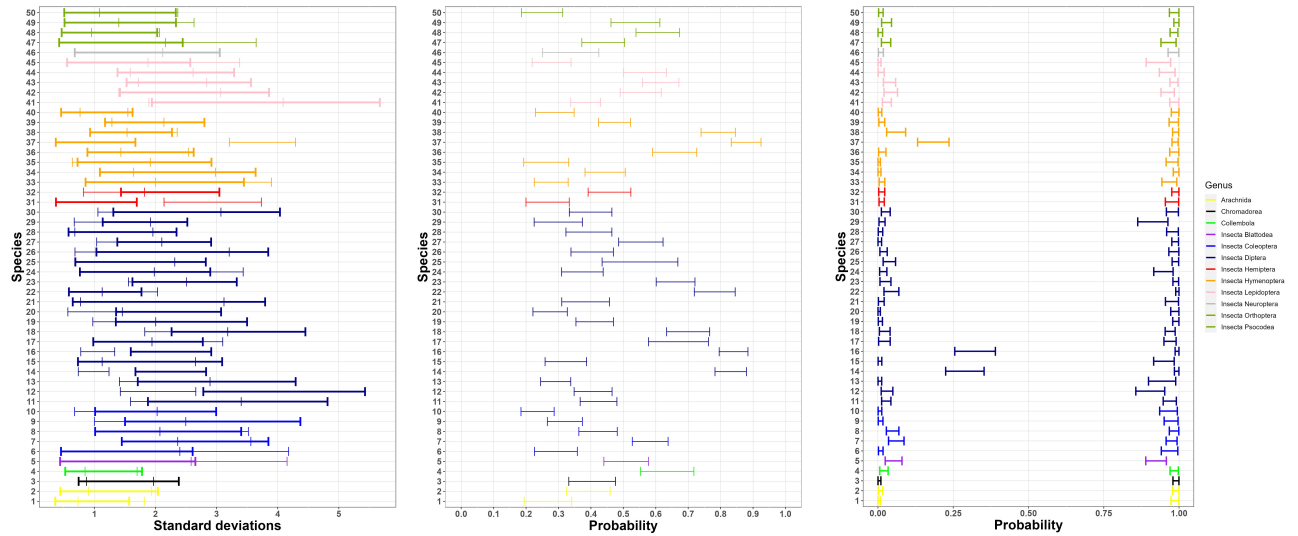
583 Finally, Figure 8 (a) suggests that generally, there is a similar amount of variation  
584 between sites and between samples for these species. As suggested by Figure 8 (b), the  
585 species that we have considered have similar collection probabilities across the several sites,  
586 possibly due to the fact that the most frequently detected species across PCRs have been  
587 selected. Figure 8 (c) demonstrates, as expected, that the Stage 2 true positive probability  
588 is close to 1 for all species. We highlight here that this probability is modelled as species-  
589 specific but assumed constant across all replicates. Similarly, the figure also suggests that  
590 the probability of a Stage 2 false negative error is very close to 0 for all but three species.  
591 One of these three (sp. 14) is in the fly family Tachinidae, which are parasitoids of other  
592 insects and thus might have been collected not only as adults but also occasionally as eggs  
593 attached to the adults of other (insect) host species, with the latter case being classified as  
594 false positives in Stage 2, given that an egg would contribute very low amounts of starting  
595 DNA biomass.

596 elfig:cov`results

## 597 **6 Discussion**

598

599 Over the last decade, DNA-based biodiversity studies, primarily using metabarcoding,  
600 have rapidly increased in popularity, and multivariate statistical models are now starting  
601 to be deployed to analyse metabarcoding data (e.g. Lin et al., 2021; Pichler and Hartig,  
602 2021; Abrego et al., 2021; Fukaya et al., 2022; Ji et al., 2022). Our paper provides the  
603 first unifying modelling framework that considers and quantifies key sources of variation,  
604 error and noise in metabarcoding surveys (Table 1). As a result, our modelling framework  
605 allows more reliable and more powerful biodiversity monitoring and inference on species



**Figure 8:** Case study: (left) 95% PCI of the species-specific between-samples standard deviation  $\sigma_s$  and between-sites standard deviation  $\sqrt{T_{ss}}$  (in bold). (center) 95% PCI of the species-specific average collection probability  $\theta_{im}^s$  across all sites. (right) 95% PCI of the species-specific Stage 2 false-positive probabilities  $q_s$  (on the left of the plot) and true-positive probabilities  $p_s$  (on the right of the plot). Species are grouped taxonomically.

606 responses to landscape characteristics than has been possible before. We have employed,  
 607 extended, and developed a number of inferential tools to deal with the complexity of the  
 608 proposed hierarchical model, which involves two latent stages and a large number of latent  
 609 variables. Finally, this is the first modelling approach that accounts for spike-ins and  
 610 negative controls (empty tubes), which are widely used quality-control methods in DNA-  
 611 based biodiversity surveys but rarely explicitly considered within a modelling framework.  
 612 We explored the benefits of spike-ins on inference and provided analytical and simulation  
 613 results of the effects of study design choices on parameter estimates. As is the case in all  
 614 models, we make certain assumptions about the data-generating process and if (any of)  
 615 these assumptions are violated, then inference can be biased. Below, we discuss the key  
 616 assumptions and corresponding model extensions, when appropriate.

617 Our new framework allows us to infer and map species DNA biomass change across  
 618 surveyed sites (Figure 7 (b)), and to link these to landscape characteristics (Figure 6).  
 619 The resulting maps can be used to identify areas of high conservation value, as well as  
 620 areas where particular species or groups of species are more or less prevalent, and to  
 621 detect species-specific shifts, expansions, and shrinkage. We are also able to study pairwise

622 correlations across large numbers of species (Figure 7 (a)), which is considerably more  
623 scalable using metabarcoding data than using standard observational data. We note that,  
624 as discussed in the corresponding sections of the model and the case study, we cannot  
625 unambiguously identify the sources of the estimated correlations using the available data  
626 alone, as factors other than the affinity between species, such as competition for primers,  
627 could affect the inferred species correlations. We have shown that using spike-ins can  
628 substantially increase inference accuracy for parameters of interest (Figure 5). Our results  
629 also demonstrate that the current practice of collecting a single sample from each surveyed  
630 site considerably reduces our ability to infer changes in species DNA biomass and that  
631 replication at both stages as well as the use of normalisation-ratio offsets or spike-ins is the  
632 optimal approach to designing metabarcoding studies (Figure 3).

633 In metabarcoding data, the baseline DNA biomass of each species is confounded with its  
634 species-specific collection and amplification rates. Hence, we cannot infer absolute values  
635 of species-specific DNA biomass across sites using metabarcoding data alone. However,  
636 by assuming that baseline species-specific collection and amplification rates are the same  
637 across sites, samples, and PCR replicates, we can infer species-specific DNA biomass *change*  
638 across sites, species-specific covariate effects, and pairwise species correlations. Finally, we  
639 model species amplification rates as independent random effects, but competition between  
640 species for primers, polymerases and nucleotides during PCR amplification might violate  
641 this independence assumption, and future experimental work, alongside model extensions,  
642 should explore this issue further.

643 We note that we have not allowed the probability of Stage 2 species detection,  $p_s$ , to vary  
644 between samples or PCR replicates, and hence we have assumed that it does not depend  
645 on the DNA biomass of other species in the sample/PCR replicate. However, because  
646 of the PCR product normalisation step, described in Section 1.1, in PCR replicates with  
647 relatively high resulting overall DNA biomass, relatively low-DNA-biomass species might  
648 be less likely to be drawn in high enough concentration to be detected, an issue that is

649 often referred to as PCR dropout. Empirically, it is known that such PCR competition can  
650 be mitigated by using a lower number of PCR cycles (Yang et al., 2021) and by sequencing  
651 each sample replicate more deeply. When extending the model of this paper, Stage 2 species  
652 detection can be modelled as a function of DNA biomass, so that  $\text{logit}(p_{imk}^s) = \beta_0^p + \beta_p(v_{im}^s +$   
653  $o_{imk})$ . Model extensions of this type are important but are expected to introduce further  
654 identifiability issues and computational challenges and hence require careful investigation.

655 Generally, modelling changes in (proxies of) abundance, such as changes in DNA  
656 biomass, is a more powerful monitoring tool than modelling changes in species presence  
657 across survey sites (Joseph et al., 2006). Metabarcoding studies yield count data without  
658 any consequence on associated cost, and hence overcome the time and cost implications  
659 associated with collecting count data for multiple species. Our model uses the raw count  
660 data, and does not rely on ad-hoc rules about what constitutes a practically zero count for  
661 converting them to binary data, which has been the standard practice thus far (Ovaskainen  
662 et al., 2017; Bush et al., 2020). To model changes in (log)biomass for each species across  
663 sites, we rely on the investigator being able to record any normalisation steps (or to include  
664 a spike-in), otherwise the relationship between change in read counts and change in the  
665 amount of biomass in the environment cannot be inferred, and instead the counts can only  
666 be used to infer composition, as is standard practice in metabarcoding studies. We have  
667 allowed for over-dispersion in the count data using a negative binomial distribution, but  
668 future work could consider alternative parametrisations, such as the discrete Weibull distri-  
669 bution. The model can also be extended to account for multiple primers or for differences  
670 between labs, if samples are processed by more than one lab, by introducing regression  
671 models for corresponding parameters.

672 Metabarcoding studies, particularly when applied to microbiomes and meiofauna (e.g.  
673 nematodes, micro-eukaryotes), can detect 1000s of species, which leads to large numbers of  
674 latent variables and coefficients in the model. There are several ways that the inferential  
675 tools presented here could be further extended to scale to these cases. Firstly, the posterior

676 distribution conditional on the  $u_{imk}$  is independent across species. If  $u_{imk}$  could be esti-  
677 mated at a first stage then inference across species could be easily parallelized. Secondly,  
678 variational Bayes methods could be applied to avoid the use of sampling methods. The  
679 choice of variational distribution will be important and can exploit the conditional normal-  
680 ity of much of the model. Alternatively, the model could be adapted by assuming that the  
681 coefficient matrices such as  $\beta^z = (\beta_1^z, \dots, \beta_S^z)$ , have a low-dimensional representation. We  
682 highlight that in its current format, the model assumes species-specific parameters, and  
683 hence there is potentially a large number of parameters to be estimated for each species.  
684 Therefore, if a species only has a few non-zero PCR reads from potentially only a few sites,  
685 estimating all of these species-specific parameters is difficult. Future work should explore  
686 sharing parameters between species, making inference for rarely-observed species possible.

687 We are not modelling species presence/absence and instead we have focused on mod-  
688 elling biomass on a continuous scale. As a result, we cannot infer whether a species is  
689 absent from a particular study site, but instead only if its DNA biomass at a given site is  
690 practically zero. We have assumed that a sample which already contains DNA biomass of  
691 a species cannot be further contaminated by the DNA of the same species from another  
692 sample or site in Stage 1. This is a reasonable but also necessary assumption, because of  
693 model identifiability issues otherwise. It is possible that there exists contamination from  
694 other sites if their samples are all processed in the same laboratory, especially at the same  
695 time, or that there is contamination during the collection or transfer of samples. However,  
696 with only metabarcoding data to hand, it is not possible to identify the source of contami-  
697 nation, or to model the possibility that a sample that contains DNA of a species has been  
698 further contaminated by the DNA of the same species from another sample or site in Stage  
699 1. This is yet another reason to take measures that minimise contamination risk.

700 eDNA metabarcoding has revolutionised the cost-effectiveness, precision, and scale at  
701 which biodiversity assessment can be performed. Nevertheless, the multiple stages at which  
702 imperfect detection of DNA biomass can occur during the workflow are not insignificant. By

703 facilitating estimates of within-species changes in DNA biomass as a function of covariates,  
704 while accounting for workflow uncertainties, our modelling framework provides a substantial  
705 improvement in the design and analysis of eDNA metabarcoding data.

## 706 **Data Availability**

707 The sequence data, bioinformatic scripts, and the three sample by species tables and  
708 environmental covariates are archived on DataDryad at doi.org/10.5061/dryad.4f4qrfjbb.

## 709 **Acknowledgments**

710 The work was funded by NERC project NE/T010045/1 “Integrating new statistical  
711 frameworks into eDNA survey and analysis at the landscape scale” and benefited from the  
712 sCom Working Group at iDiv.de. DWY and MJL were supported by the Strategic Priority  
713 Research Program of Chinese Academy of Sciences, Grant No. XDA20050202, the Key  
714 Research Program of Frontier Sciences, CAS (QYZDY-SSW-SMC024), the State Key Lab-  
715 oratory of Genetic Resources and Evolution (GREKF19-01, GREKF20-01, GREKF21-01)  
716 at the Kunming Institute of Zoology, and the University of Chinese Academy of Sciences.

## 717 **Disclosure statement**

718 The authors declare that there are no conflicts of interest relevant to this article.

## 719 **References**

- 720 Abrego, N., Roslin, T., Huotari, T., et al. (2021). Accounting for species interactions is  
721 necessary for predicting how arctic arthropod communities respond to climate change.  
722 *Ecography*, 44(6):885–896.
- 723 Andrieu, C. and Thoms, J. (2008). A tutorial on adaptive MCMC. *Statistics and Comput-*  
724 *ing*, 18(4):343–373.
- 725 Baisero, D., Schuster, R., and Plumptre, A. J. (2022). Redefining and mapping global  
726 irreplaceability. *Conservation Biology*, 36(2):e13806.
- 727 Besson, M., Alison, J., Bjerger, K., et al. (2022). Towards the fully automated monitoring  
728 of ecological communities. *Ecological Letters*, to appear.
- 729 Bush, A., Monk, W. A., Compson, Z. G., et al. (2020). DNA metabarcoding reveals



730 metacommunity dynamics in a threatened boreal wetland wilderness. *Proceedings of the*  
731 *National Academy of Sciences*, 117(15):8539–8545.

732 Bush, A., Sollmann, R., Wilting, A., et al. (2017). Connecting Earth observation to high-  
733 throughput biodiversity data. *Nature Ecology & Evolution*, 1(7):0176.

734 Buxton, A., Matechou, E., Griffin, J., et al. (2021). Optimising sampling and analysis  
735 protocols in environmental DNA studies. *Scientific Reports*, 11(1):11637.

736 Carraro, L., Hartikainen, H., Jokela, J., et al. (2018). Estimating species distribution and  
737 abundance in river networks using environmental DNA. *Proceedings of the National*  
738 *Academy of Sciences*, 115(46):11724–11729.

739 Clare, E. L., Economou, C. K., Bennett, F. J., and others. (2022). Measuring biodiversity  
740 from DNA in the air. *Current Biology*, 32(3):693–700.e5.

741 Clausen, D. S. and Willis, A. D. (2022). Modeling complex measurement error in micro-  
742 biome experiments. *arXiv preprint arXiv:2204.12733*.

743 Coblenz, K. E., Rosenblatt, A. E., and Novak, M. (2017). The application of Bayesian  
744 hierarchical models to quantify individual diet specialization. *Ecology*, 98(6):1535–1547.

745 Datta, J. and Dunson, D. B. (2016). Bayesian inference on quasi-sparse count data.  
746 *Biometrika*, 103:971–983.

747 Dawid, A. P. (1981). Some matrix-variate distribution theory: notational considerations  
748 and a Bayesian application. *Biometrika*, 68(1):265–274.

749 Ficetola, G. F., Pansu, J., Bonin, A., et al. (2015). Replication levels, false presences  
750 and the estimation of the presence/absence from eDNA metabarcoding data. *Molecular*  
751 *Ecology Resources*, 15(3):543–556.

752 Fordyce, J. A., Gompert, Z., Forister, M. L., and Nice, C. C. (2011). A hierarchical Bayesian  
753 approach to ecological count data: a flexible tool for ecologists. *PloS One*, 6(11):e26785.

754 Frøslev, T. G., Kjølner, R., Bruun, H. H., et al. (2019). Man against machine: Do fungal  
755 fruitbodies and eDNA give similar biodiversity assessments across broad environmental  
756 gradients? *Biological Conservation*, 233:201–212.

757 Fukaya, K., Kondo, N. I., Matsuzaki, S.-i. S., and Kadoya, T. (2022). Multispecies site occu-  
758 pancy modelling and study design for spatially replicated environmental DNA metabar-  
759 coding. *Methods in Ecology and Evolution*, 13(1):183–193.

760 Gelman, A. (2006). Prior distributions for variance parameters in hierarchical models.  
761 *Bayesian Analysis*, 1(3):515–533.

762 Griffin, J. E., Matechou, E., Buxton, A. S., et al. (2020). Modelling environmental DNA

763 data; Bayesian variable selection accounting for false positive and false negative errors.  
764 *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 69(2):377–392.

765 Guillera-Arroita, G., Lahoz-Monfort, J., van Rooyen, A., Weeks, A., and Tingley, R. (2017).  
766 Dealing with false-positive and false-negative errors about species occurrence at multiple  
767 levels. *Methods in Ecology and Evolution*, 8(9):1081–1091.

768 Hebert, P. D., Cywinska, A., Ball, S. L., and DeWaard, J. R. (2003). Biological identifi-  
769 cations through DNA barcodes. *Proceedings of the Royal Society of London. Series B:*  
770 *Biological Sciences*, 270(1512):313–321.

771 Ji, Y., Ashton, L., Pedley, S. M., et al. (2013). Reliable, verifiable and efficient monitoring  
772 of biodiversity via metabarcoding. *Ecology Letters*, 16(10):1245–1257.

773 Ji, Y., Baker, C. C. M., Popescu, V. D., et al. (2022). Measuring protected-area effectiveness  
774 using vertebrate distributions from leech iDNA. *Nature Communications*, 13(1):1555.

775 Joseph, L. N., Field, S. A., Wilcox, C., and Possingham, H. P. (2006). Presence–absence ver-  
776 sus abundance data for monitoring threatened species. *Conservation Biology*, 20(6):1679–  
777 1687.

778 Ley, R. (2022). The human microbiome: there is much left to do. *Nature*,  
779 606(7914):435–435.

780 Li, Y., Craig, B. A., and Bhadra, A. (2019). The graphical horseshoe estimator for inverse  
781 covariance matrices. *Journal of Computational and Graphical Statistics*, 28(3):747–757.

782 Li, Y., Devenish, C., Tosa, et al. (2024). Combining environmental dna and remote sensing  
783 for efficient, fine-scale mapping of arthropod biodiversity. *Philosophical Transactions of*  
784 *the Royal Society B: Biological Sciences*, 379(1904):20230123.

785 Lin, M., Simons, A. L., et al. (2021). Landscape analyses using edna metabarcoding and  
786 earth observation predict community biodiversity in california. *Ecological Applications*,  
787 31(6):e02379.

788 Lindahl, B. D., Nilsson, R. H., Tedersoo, L., et al. (2013). Fungal community analysis  
789 by high-throughput sequencing of amplified markers—a user’s guide. *New Phytologist*,  
790 199(1):288–299.

791 Luo, M., Ji, Y., Warton, D., and Yu, D. W. (2022). Extracting abundance information  
792 from DNA-based data. *Molecular Ecology Resources*, to appear.

793 McLaren, M. R., Willis, A. D., and Callahan, B. J. (2019). Consistent and correctable bias  
794 in metagenomic sequencing experiments. *Elife*, 8:e46923.

795 Mordecai, R. S., Mattsson, B. J., Tzilkowski, C. J., and Cooper, R. J. (2011). Addressing

796 challenges when studying mobile or episodic species: hierarchical Bayes estimation of  
797 occupancy and use. *Journal of Applied Ecology*, 48(1):56–66.

798 Ovaskainen, O. and Abrego, N. (2020). *Joint Species Distribution Modelling: With Appli-*  
799 *cations in R*. Cambridge University Press.

800 Ovaskainen, O., Tikhonov, G., Dunson, D., et al. (2017). How are species interactions  
801 structured in species-rich communities? A new method for analysing time-series data.  
802 *Proceedings of the Royal Society B: Biological Sciences*, 284(1855):20170768.

803 Papaspiliopoulos, O., Roberts, G. O., and Sköld, M. (2007). A general framework for the  
804 parametrization of hierarchical models. *Statistical Science*, 22(1):59–73.

805 Papaspiliopoulos, O., Roberts, G. O., and Zanella, G. (2020). Scalable inference for crossed  
806 random effects models. *Biometrika*, 107(1):25–40.

807 Pichler, M. and Hartig, F. (2021). A new joint species distribution model for faster and  
808 more accurate inference of species associations from big community data. *Methods in*  
809 *Ecology and Evolution*.

810 Piper, A. M., Batovska, J., Cogan, N. O. I., et al. (2019). Prospects and challenges of  
811 implementing DNA metabarcoding for high-throughput insect surveillance. *GigaScience*,  
812 8(8):giz092.

813 Ratnasingham, S. and Hebert, P. D. (2007). Bold: The barcode of life data system  
814 (<http://www.barcodinglife.org>). *Molecular Ecology Notes*, 7(3):355–364.

815 Roberts, G. O. and Rosenthal, J. S. (2009). Examples of adaptive MCMC. *Journal of*  
816 *Computational and Graphical Statistics*, 18(2):349–367.

817 Saine, S., Ovaskainen, O., Somervuo, P., and Abrego, N. (2020). Data collected by fruit  
818 body-and DNA-based survey methods yield consistent species-to-species association net-  
819 works in wood-inhabiting fungal communities. *Oikos*, 129(12):1833–1843.

820 Schmidt, B. R., Kery, M., Ursenbacher, S., et al. (2013). Site occupancy models in the  
821 analysis of environmental DNA presence/absence surveys: a case study of an emerging  
822 amphibian pathogen. *Methods in Ecology and Evolution*, 4(7):646–653.

823 Taberlet, P., Bonin, A., Zinger, L., and Coissac, E. (2018). *Environmental DNA: for*  
824 *biodiversity research and monitoring*. Oxford University Press, Oxford, United Kingdom.

825 Takahara, T., Minamoto, T., Yamanaka, H., et al. (2012). Estimation of fish biomass using  
826 environmental DNA. *PloS one*, 7(4):e35868.

827 Thomsen, P. F. and Sigsgaard, E. E. (2019). Environmental DNA metabarcoding of wild  
828 flowers reveals diverse communities of terrestrial arthropods. *Ecology and Evolution*,

829 9(4):1665–1679.

830 Thomsen, P. F. and Willerslev, E. (2015). Environmental DNA – An emerging tool in con-  
831 servation for monitoring past and present biodiversity. *Biological Conservation*, 183:4–18.

832 Tkacz, A., Hortala, M., and Poole, P. S. (2018). Absolute quantitation of microbiota  
833 abundance in environmental samples. *Microbiome*, 6(1):110.

834 Ver Hoef, J. M. and Boveng, P. L. (2007). Quasi-Poisson vs. Negative Binomial Regression:  
835 How Should We Model Overdispersed Count Data? *Ecology*, 88:2766–2772.

836 Wang, H. (2012). Bayesian graphical lasso models and efficient posterior computation.  
837 *Bayesian Analysis*, 7(4):867–886.

838 Yang, C., Bohmann, K., Wang, X., and others. (2021). Biodiversity Soup II: A bulk-sample  
839 metabarcoding pipeline emphasizing error reduction. *Methods in Ecology and Evolution*,  
840 12(7):1252–1264.

841 Zanella, G. and Roberts, G. (2021). Multilevel linear models, Gibbs samplers and multigrid  
842 decompositions (with discussion). *Bayesian Analysis*, 16(4):1309–1391.

# Supplementary material of eDNAPlus: A unifying modelling framework for DNA-based biodiversity monitoring

Alex Diana<sup>1</sup>, Eleni Matechou<sup>1</sup>, Jim Griffin<sup>2</sup>, Douglas W. Yu<sup>3,4,5</sup>, Mingjie Luo<sup>4</sup>, Marie Tosa<sup>5</sup>, Alex Bush<sup>6</sup>, Richard Griffiths<sup>7</sup>

<sup>1</sup> School of Mathematics, Statistics and Actuarial Science, University of Kent, UK,

<sup>2</sup> Department of Statistical Science, University College London, UK,

<sup>3</sup> School of Biological Sciences, University of East Anglia, UK,

<sup>4</sup> Kunming College of Life Sciences, University of Chinese Academy of Sciences, China,

<sup>5</sup> Center for Excellence in Animal Evolution and Genetics &

State Key Laboratory of Genetic Resources and Evolution &

Yunnan Key Laboratory of Biodiversity and Ecological Security of Gaoligong Mountain &

Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming, China,

<sup>6</sup> Department of Fisheries, Wildlife, & Conservation Sciences, Oregon State University, Oregon USA,

<sup>7</sup> Lancaster Environment Centre, University of Lancaster, UK,

<sup>8</sup> Durrell Institute of Conservation and Ecology, University of Kent, UK

September 29, 2024

# 1 Inference

We can reparameterise the model by expressing the negative binomial using the Poisson-Gamma mixture and a centred parameterisation:

$$\left\{ \begin{array}{l}
 \bar{\beta}_0^s \sim N(\lambda_s, \sigma_\beta^2) \\
 \bar{L} = \{\bar{l}_i^s\} \sim \text{MN}(\bar{B}_0 + X_z B, \Sigma, T) \\
 T^{-1} \sim \text{GH} \\
 \text{logit}(\theta_{im}^s) = (\bar{l}_i^s - \lambda_s)\phi_{1s} + X_{im}^w \phi_s \\
 \mathbb{P}(\delta_{im}^s = 1) = \theta_{im}^s \\
 \mathbb{P}(\gamma_{im}^s = 1 \mid \delta_{im}^s = 0) = \zeta^s \\
 \bar{\mu}_s \sim N(\lambda_s, \sigma_\mu^2) \\
 \bar{v}_{im}^s \sim \begin{cases} N(\bar{l}_i^s + X_{im}^w \beta_s^W, \sigma_s^2) & \text{if } \delta_{im}^s = 1 \\ N(\bar{\mu}_s, \nu_s^2) & \text{if } \delta_{im}^s = 0, \gamma_{im}^s = 1 \end{cases} \\
 \mathbb{P}(c_{imk}^s = 1 \mid \delta_{im}^s = 1 \text{ or } \gamma_{im}^s = 1) = p_s \\
 \mathbb{P}(c_{imk}^s = 2 \mid \delta_{im}^s = 0, \gamma_{im}^s = 0) = q_s \\
 \eta_{imk}^s \sim \text{Gamma}\left(r_s, \frac{r_s}{\exp(\bar{v}_{im}^s + u_{imk} + o_{imk})}\right) & \text{if } c_{imk}^s = 1 \\
 y_{imk}^s \sim \begin{cases} \pi \delta_0 + (1 - \pi)(1 + \text{NB}(\mu_0, n_0)) & \text{if } c_{imk}^s = 0 \\ \text{Pois}(\eta_{imk}^s) & \text{if } c_{imk}^s = 1 \\ \text{Pois}(\tilde{\mu}) & \text{if } c_{imk}^s = 2 \end{cases} \\
 u_{imk} \sim N(0, \sigma_u^2)
 \end{array} \right.$$

where we have defined

$$\left\{ \begin{array}{l}
 \bar{\beta}_0^s = \beta_0^s + \lambda_s \\
 \bar{l}_i^s = l_i^s + \lambda_s \\
 \bar{v}_{im}^s = v_{im}^s + \lambda_s \\
 \bar{\mu}_s = \mu_s + \lambda_s
 \end{array} \right. .$$

If not stated, we use a Metropolis-Hastings update with a Laplace approximation as

proposal if a full conditional distribution is not tractable.

- Update  $\lambda_s$

The full conditional has the density

$$N(\bar{\beta}_0^s | \lambda_s, \sigma_\beta^2) N(\bar{\mu}_s | \lambda_s, \sigma_\mu^2) \left( \prod_{i,m} \text{Be}(\delta_{im}^s, \text{logit}(-\phi_{1s}\lambda_s + \phi_{1s}\bar{l}_i^s + X_{im}^w \phi_s)) \right)$$

and the parameter is updated using a Metropolis-Hastings random walk.

- Update  $\eta_{imk}^s$

$$\eta_{imk}^s \sim \text{Gamma} \left( r_s + y_{imk}^s, 1 + \frac{r_s}{\exp(\bar{v}_{im}^s + u_{imk} + o_{imk})} \right)$$

- Update  $r_s$

The full conditional distribution has density

$$p(r_s | \cdot) \propto N(r_s | \mu_r, \sigma_r^2) \prod_{i,m,k: c_{imk}^s=1} \text{NB}(y_{imk}^s | \exp(\bar{v}_{im}^s + u_{imk} + o_{imk}), r_s)$$

and the parameter is updated using a Metropolis-Hastings random walk.

- Update  $\bar{v}_{im}^s$  and  $u_{imk}$

We define the overall means  $\hat{v}_{im} = \frac{1}{n_{im}} \sum_{s: \delta_{im}^s=1 \text{ or } \gamma_{im}^s=1} \bar{v}_{im}^s$ , where  $n_{im} = \sum_{s: \delta_{im}^s=1 \text{ or } \gamma_{im}^s=1} 1$ , and  $\hat{u}_{im} = \frac{1}{K} \sum_k u_{imk}$ , and the increments  $\tilde{v}_{im}^s = \bar{v}_{im}^s - \hat{v}_{im}$  and  $\tilde{u}_{imk} = u_{imk} - \hat{u}_{im}$ .

We first sample from the joint full conditional of the overall means conditional on the increments  $(\hat{v}_{im}, \hat{u}_{im} | \tilde{v}_{im}^1, \dots, \tilde{v}_{im}^S, \tilde{u}_{im1}, \dots, \tilde{u}_{imK}, \dots)$ , which has density of the form

$$p(\hat{v}_{im}, \hat{u}_{im} | \cdot) \propto N \left( \hat{v}_{im} \left| \frac{1}{n_{im}} \sum_{s: \delta_{im}^s=1 \text{ or } \gamma_{im}^s=1} l_i^s, \frac{1}{n_{im}} \sum_{s: \delta_{im}^s=1 \text{ or } \gamma_{im}^s=1} \sigma_s^2 \right. \right) N \left( \hat{u}_{im} \left| 0, \frac{\sigma_u^2}{K_{im}} \right. \right)$$

$$\prod_{k,s: c_{imk}^s=1} \text{Gamma} \left( \eta_{imk}^s | r_s, \frac{r_s}{\exp(\hat{v}_{im} + \tilde{v}_{im}^s + \hat{u}_{im} + \tilde{u}_{imk} + o_{imk})} \right)$$

where  $n_{im} = \sum_{s: \delta_{im}^s=1} 1$ .

Next, we sample  $\bar{v}_{im}^s$  from its full conditional distribution which has density

$$p(\bar{v}_{im}^s|\cdot) \propto \text{N}(\bar{v}_{im}^s|\bar{l}_i^s + X_{im}^w \beta_s^W, \sigma_s^2)^{\delta_{im}^s} \text{N}(\bar{v}_{im}^s|\bar{\mu}_s, \nu_s^2)^{(1-\delta_{im}^s)\gamma_{im}^s} \\ \times \prod_{k: c_{imk}^s=1} \text{Gamma}\left(\eta_{imk}^s|r_s, \frac{r_s}{\exp(\bar{v}_{im}^s + u_{imk} + o_{imk})}\right),$$

and  $u_{imk}$  from its full conditional distribution which has density

$$p(u_{imk}|\cdot) \propto \text{N}(u_{imk}|0, \sigma_u^2) \prod_{s: c_{imk}^s=1} \text{Gamma}\left(\eta_{imk}^s|r_s, \frac{r_s}{\exp(\bar{v}_{im}^s + u_{imk} + o_{imk})}\right),$$

In all three cases, we use a Metropolis-Hastings update with a Laplace approximation as proposal.

- Update  $\bar{l}_i^s$

The parameter can be updated from its full conditional distribution which has density

$$p(\bar{l}_i^s|\cdot) \propto \text{N}(\bar{l}_i^s|\bar{\beta}_0^s + X_i \beta_s^z, \bar{l}_{-i}^s, T, \Sigma) \prod_{m: \delta_{im}^s=1} \text{N}(\bar{v}_{im}^s|\bar{l}_i^s + X_{im} \beta_s^W, \sigma_s^2) \\ \prod_m \text{Be}(\delta_{im}^s|\text{logit}(-\phi_{1s}\lambda + \phi_{1s}\bar{l}_i^s + X_{im}^W \phi_s))$$

using a Metropolis-Hastings update with a Laplace approximation as proposal. The conditional distribution  $\text{N}(\bar{l}_i^s|\bar{\beta}_0^s + X_i \beta_s^z, \bar{l}_{-i}^s, T, \Sigma)$  can be efficiently computed using the algorithm described in Section 1.1.

- Update  $(B_0, B)$

These parameters are updated from their joint full conditional distribution which is

$$p(B_0, B|l) \propto \text{N}(\text{vec}((B_0, B))|0, \sigma_\beta^2 I_{(p+1)s}) \text{MN}(B_0 + X_z B, \Sigma, T)$$

and can be written in closed form as

$$\text{vec}(B_0, B) \sim \text{N}(\Lambda_B^{-1} \mu_B, \Lambda_B^{-1})$$



where  $\Lambda_B = T^{-1} \otimes (\bar{X}^T \Sigma^{-1} \bar{X})$  and  $\mu_B = (I_S \otimes \bar{X}^T)(\Sigma^{-1} \text{vec}(L)T^{-1})$ , with  $\bar{X} = (1, X)$

- Update  $\beta_s^w$

$$\beta_s^w \sim \text{N}(\Lambda_\beta^{-1} \mu_\beta, \Lambda_\beta^{-1})$$

where  $\Lambda_\beta^{-1} = \frac{(X_\delta^W)^T X^W}{\sigma_s^2} + \sigma_\beta^2 I_{n_w}$  and  $\mu_\beta = \frac{(X_\delta^W)^T}{\sigma_s^2} (v_\delta^s - l_\delta^s)$ , and  $X_\delta^W$ ,  $v_\delta^s$ ,  $l_\delta^s$  are the subset of  $X^w$ ,  $v^s$  and  $l^s$ , respectively, such that  $\delta_{im}^s = 1$  for species  $s$ .

- Update  $\sigma_s^2$

$$\sigma_s^2 \sim \text{IG}(a_\sigma + n_\delta, b_\sigma + s_\delta)$$

where  $n_\delta = \sum_{i,m} 1_{\delta_{i,m}^s=1}$  and  $s_\delta = \sum_{i,m:\delta_{i,m}^s=1} (\bar{v}_{im}^s - \bar{l}_i^s - X_{im}^w \beta_s^w)^2$ .

- Update  $\bar{\mu}_s$

$$\bar{\mu}_s \sim \text{N}(m_\mu, \sigma_\mu^2)$$

where  $\sigma_\mu^2 = \left(\frac{1}{\sigma_\mu^2} + \frac{n_\gamma}{v_s^2}\right)^{-1}$ ,  $m_\mu = \left(\frac{\lambda_s}{\sigma_\mu^2} + \frac{\sum_{\gamma_{im}^s=1} \bar{v}_{im}^s}{v_s^2}\right) \sigma_\mu^2$  and  $n_\gamma = \sum_{i,m} 1_{\gamma_{im}^s=1}$

- Update  $\phi_{1s}$  and  $\phi_s$ .

Since these are coefficients of a logistic regression, we use the Pólya-gamma updating scheme (Polson et al., 2013) where  $\delta_{im}^s$  are responses and the regressors are  $\bar{l}_i^s - \lambda_s$  and  $X_{im}^w$ .

- Update  $c_{imk}^s$ ,  $\delta_{im}^s$  and  $\gamma_{im}^s$

We update  $c_{im1}^s, \dots, c_{imK}^s$ ,  $\delta_{im}^s$  and  $\gamma_{im}^s$  in a single block and, for ease of notation, we drop the indices  $i$ ,  $m$  and  $s$  when describing the update. During the burn-in phase, we update these parameters by sampling from their full conditional distribution. We note that since the variable  $\bar{v}$  is not present in the model if  $\delta = \gamma = 0$ , this update requires a reversible jump Markov Chain Monte Carlo (RJMCMC) move.

If  $\delta \neq 0$  or  $\gamma \neq 0$ , we propose  $c_1, \dots, c_K$ ,  $\delta$  and  $\gamma$  from their joint distribution evaluated using the currently sampled value of  $v^*$ . Otherwise, we propose  $v^*$  using the informed

proposal  $v^* \sim N(\mu^*, .5^2)$ , where  $\exp(\mu^*) = \frac{1}{K} \sum_{k=1}^K \frac{y_{imk}^s}{\exp(\lambda_s + u_{imk})}$  and propose  $c_1, \dots, c_K, \delta$  and  $\gamma$  from their joint distribution evaluated using this proposed value.

Using the notation  $\delta^*$ ,  $\gamma^*$  and  $c^*$  for the proposed value, the Metropolis-Hastings acceptance ratio are as follows. If  $\delta = \gamma = 0$ , the MH ratio takes the form

$$\begin{aligned} \min & \left\{ 1, \frac{p(y_1, \dots, y_K | \delta^*=0, \gamma^*=0, c_1^*, \dots, c_K^*) p(\delta^*=0, \gamma^*=0, c_1^*, \dots, c_K^*) q(\delta=0, \gamma=0, c_1, \dots, c_K)}{p(y_1, \dots, y_K | \delta=0, \gamma=0, c_1, \dots, c_K) p(\delta=0, \gamma=0, c_1, \dots, c_K) q(\delta^*=0, \gamma^*=0, c_1^*, \dots, c_K^*)} \right\} & \text{if } \delta^* = \gamma^* = 0 \\ \min & \left\{ 1, \frac{p(y_1, \dots, y_K | \delta^*=1, \gamma^*=1, c_1^*, \dots, c_K^*, v^*) p(\delta^*=1, \gamma^*=1, c_1^*, \dots, c_K^*) q(\delta=0, \gamma=0, c_1, \dots, c_K) p(v^*)}{p(y_1, \dots, y_K | \delta=0, \gamma=0, c_1, \dots, c_K) p(\delta=0, \gamma=0, c_1, \dots, c_K) q(\delta^*=1, \gamma^*=1, c_1^*, \dots, c_K^*, v^*)} \right\} & \text{if } \delta^* + \gamma^* = 1 \end{aligned}$$

or, if  $\delta + \gamma = 1$

$$\begin{aligned} \min & \left\{ 1, \frac{p(y_1, \dots, y_K | \delta^*=1, \gamma^*=1, c_1^*, \dots, c_K^*, v) p(\delta^*=1, \gamma^*=1, c_1^*, \dots, c_K^*) p(v) q(\delta=1, \gamma=1, c_1, \dots, c_K, v)}{p(y_1, \dots, y_K | \delta=1, \gamma=1, c_1, \dots, c_K, v) p(\delta=1, \gamma=1, c_1, \dots, c_K) p(v) q(\delta^*=1, \gamma^*=1, c_1^*, \dots, c_K^*, v)} \right\} & \text{if } \delta^* + \gamma^* = 1 \\ \min & \left\{ 1, \frac{p(y_1, \dots, y_K | \delta^*=0, \gamma^*=0, c_1^*, \dots, c_K^*) p(\delta^*=0, \gamma^*=0, c_1^*, \dots, c_K^*) q(\delta=1, \gamma=1, c_1, \dots, c_K, v)}{p(y_1, \dots, y_K | \delta=1, \gamma=1, c_1, \dots, c_K, v) p(v^*) p(\delta=1, \gamma=1, c_1, \dots, c_K) q(\delta^*=0, \gamma^*=0, c_1^*, \dots, c_K^*)} \right\} & \text{if } \delta^* = \gamma^* = 0 \end{aligned}$$

Given a proposal  $q(v)$ , the algorithm can be made into a Gibbs sampler by choosing the proposals

$$q(\delta = 1, \gamma = 1, c_1, \dots, c_K | v) \propto \frac{p(y_1, \dots, y_K | \delta = 1, \gamma = 1, v) p(v) p(\delta = 1, \gamma = 1, c_1, \dots, c_K)}{q(v)}$$

$$q(\delta = 0, \gamma = 0, c_1, \dots, c_K | v) \propto p(y_1, \dots, y_K | \delta = 0, \gamma = 0, v) p(\delta = 0, \gamma = 0, c_1, \dots, c_K).$$

After the burn-in phase, we do not perform a full Gibbs sampler by computing the probability of every state but propose a new state  $(\delta^*, \gamma^*, c_1^*, \dots, c_K^*)$  from the approximation  $\hat{p}(\delta_{im}^s, \gamma_{im}^s, c_{imk}^s)$  described in the main text, which is accepted using the MH ratio defined above.

- Update  $T$

The precision matrix can be updated following the approach described by Li et al. (2019) and Wang (2012). Briefly, we introduce the variables  $\nu_{ij}$  and  $\xi$  as in Makalic and Schmidt

(2015) as

$$\left\{ \begin{array}{l} \omega_{ii} \propto \text{Exp}(\frac{\lambda}{2}) \\ \omega_{ij:i < j} \sim \text{N}(0, \lambda_{ij}^2 \tau^2) \\ \lambda_{ij:i < j}^2 \sim \text{IG}(\frac{1}{2}, \frac{1}{\nu_{ij}}) \\ \nu_{ij:i < j} \sim \text{IG}(\frac{1}{2}, 1) \\ \tau^2 \sim \text{IG}(\frac{1}{2}, \frac{1}{\xi}) \\ \xi \sim \text{IG}(\frac{1}{2}, 1) \end{array} \right. .$$

The rest of the parameters can be updated straightforwardly.

- Update  $p_s$

$$p_s \sim \text{Beta}(a_p + n_{p_s}, b_p + m_{p_s} - n_{p_s})$$

where  $n_{p_s} = \sum_{i,m,k} 1_{\delta_{im}^s + \gamma_{im}^s = 1, c_{imk}^s = 1}$  and  $m_{p_s} = \sum_{i,m,k} 1_{\delta_{im}^s + \gamma_{im}^s = 1}$

- Update  $q_s$

$$q_s \sim \text{Beta}(a_q + n_{q_s}, b_q + m_{q_s} - n_{q_s})$$

where  $n_{q_s} = \sum_{i,m,k} 1_{\delta_{im}^s = 0, \gamma_{im}^s = 0, c_{imk}^s = 2}$  and  $m_{q_s} = \sum_{i,m,k} 1_{\delta_{im}^s = 0, \gamma_{im}^s = 0}$

- Update  $\tilde{\mu}$

This parameter is update using a MH with random walk proposal

- Update  $\zeta_s$

$\zeta_s \sim \text{Beta}(a_\zeta + n_\zeta, b_\zeta + m_\zeta - n_\zeta)$  where  $n_\zeta = \sum_{i,m} 1_{\delta_{i,m}^s = 0, \gamma_{i,m}^s = 1}$  and  $m_\zeta = \sum_{i,m} 1_{\delta_{i,m}^s = 0}$ .

- Update  $\pi, \mu_0, n_0$

$$\pi \sim \text{Beta}(a_\pi + N_0, b_\pi + M_0 - N_0)$$

where  $M_0 = \sum_{s,i,m,k} 1_{c_{imk}^s = 0}$  and  $N_0 = \sum_{s,i,m,k} 1_{c_{imk}^s = 0, y_{imk}^s = 0}$ .

$\mu_0$  and  $n_0$  are updated using a MH with random walk proposal.

## 1.1 Full conditional of matrix normal distributions

**Proposition 1.1.** *Given  $x \sim MN(\mu, U, V)$ , the full conditional distribution  $x_{i,j}|x_{-(i,j)}$  has the form  $N(\tilde{\mu}_{(i,j)}, \tilde{\sigma}_{(i,j)}^2)$ , where  $\tilde{\mu}_{(i,j)}$  and  $\tilde{\sigma}_{(i,j)}^2$  can be computed according to the following algorithm:*

1. Compute  $x_1$  by solving  $U_{-i,-i}x_1 = U_{-i,i}$ .
2. Compute  $\tilde{\sigma}_{i,j}^2$  as  $V_{j,j}U_{i,i} - (V_{j,j} \otimes U_{-i,i})x_1 + U_{i,i}\tilde{V} - (U_{i,-i} \cdot x_1)\tilde{V}$ , where  $\tilde{V} = (V_{j,-j} \cdot V_{-j,-j}^{-1} \cdot V_{-j,j})$
3. Compute  $y_1$  by solving  $U_{-i,-i}y_1 = \frac{\mu_{(-i,j)} - \mu_{(-i,-j)}\tilde{V}_2}{V_{j,j} - \tilde{V}}$ , where  $\tilde{V}_2 = (V_{1,-1} \cdot V_{-1,-1}^{-1})$ .
4. Compute  $\tilde{\mu}_{i,j}$  as  $\mu_{(i,j)} + (V_{j,j} \otimes U_{-i,i})y_1 + \mu_{(i,-j)} \cdot \tilde{V}_2 - \tilde{V} \cdot U_{i,-i} \cdot y_1$ .

## 1.2 Proof

$\tilde{\mu}_{(i,j)}$  and  $\tilde{\sigma}_{(i,j)}^2$  satisfy the equations

$$\tilde{\mu}_{(i,j)} = \mu_{(i,j)} + \Sigma_{(i,j),-(i,j)} \Sigma_{-(i,j),-(i,j)}^{-1} \mu_{-(i,j),(i,j)}$$

and

$$\tilde{\sigma}_{(i,j)}^2 = \Sigma_{(i,j),(i,j)} - \Sigma_{(i,j),-(i,j)} \Sigma_{-(i,j),-(i,j)}^{-1} \Sigma_{-(i,j),(i,j)},$$

where  $\Sigma = V \otimes U$ .

W.l.o.g. we set  $i = j = 1$ . Let us define

$$\Sigma_{-(1,1),-(1,1)} = \begin{bmatrix} V_{11} \otimes U_{-1,-1} & V_{1,-1} \otimes U_{-1,\cdot} \\ V_{-1,1} \otimes U_{\cdot,-1} & V_{-1,-1} \otimes U \end{bmatrix} = \begin{bmatrix} \tilde{\Sigma}_{11} & \tilde{\Sigma}_{12} \\ \tilde{\Sigma}_{21} & \tilde{\Sigma}_{22} \end{bmatrix}$$

$$\Sigma_{-(1,1),(1,1)} = \begin{bmatrix} V_{11} \otimes U_{-1,1} \\ V_{-1,1} \otimes U_{\cdot,1} \end{bmatrix} = \begin{bmatrix} \tilde{b}_1 \\ \tilde{b}_2 \end{bmatrix} = \tilde{b}$$

$$\mu_{-(1,1),(1,1)} = \begin{bmatrix} \mu_{(-1,1)} \\ \mu_{(\cdot,-1)} \end{bmatrix} = \begin{bmatrix} \tilde{\mu}_1 \\ \tilde{\mu}_2 \end{bmatrix}$$

To compute  $\tilde{\sigma}_{(1,1)}^2$ , we need to compute  $\Sigma_{(1,1),-(1,1)} \underbrace{\Sigma_{-(1,1),-(1,1)}^{-1}}_x \Sigma_{-(1,1),(1,1)}$ . Computing  $x$  is equivalent to solving  $\Sigma_{-(1,1),-(1,1)}x = \tilde{b}$ . Using Schur complement, this is equivalent to solving  $\underbrace{(\tilde{\Sigma}_{11} - \tilde{\Sigma}_{12}\tilde{\Sigma}_{22}^{-1}\tilde{\Sigma}_{21})}_A x_1 = \tilde{b}_1 - \tilde{\Sigma}_{12}\tilde{\Sigma}_{22}^{-1}\tilde{b}_2$  and then  $\tilde{\Sigma}_{22}x_2 = \tilde{b}_2 - \tilde{\Sigma}_{21}x_1$ .

Simplifications are available.

$$A = V_{1,1} \otimes U_{-1,-1} - \underbrace{(V_{1,-1} \cdot V_{-1,-1}^{-1} \cdot V_{-1,1})}_{\tilde{V}} \otimes (U_{-1,\cdot} \cdot U^{-1} \cdot U_{\cdot,-1}) =$$

$$V_{1,1} \otimes U_{-1,-1} - \tilde{V} \otimes (U_{-1,-1}) = (V_{1,1} - \tilde{V}) \cdot U_{-1,-1}$$

and

$$\tilde{b}_1 - \tilde{\Sigma}_{12}\tilde{\Sigma}_{22}^{-1}\tilde{b}_2 = V_{11} \otimes U_{-1,1} - (V_{1,-1} \cdot V_{-1,-1}^{-1} \cdot V_{-1,1}) \otimes U_{-1,1} = (V_{1,1} - \tilde{V})U_{-1,1}$$

and therefore  $x_1$  can be computed by solving  $U_{-1,-1}x_1 = U_{-1,1}$ .

Next,

$$x_2 = \tilde{\Sigma}_{22}^{-1}(\tilde{b}_2 - \tilde{\Sigma}_{21}x_1) = (V_{-1,-1}^{-1} \otimes U^{-1})(V_{-1,1} \otimes U_{\cdot,1} - (V_{-1,1} \otimes U_{\cdot,-1})x_1) =$$

$$(V_{-1,-1}^{-1} \cdot V_{-1,1}) \otimes (U^{-1} \cdot U_{\cdot,1}) - (V_{-1,-1}^{-1} \cdot V_{-1,1}) \otimes (U^{-1} \cdot U_{\cdot,-1})x_1$$

and therefore  $\tilde{b}_2x_2 = (V_{-1,1} \otimes U_{\cdot,1})x_2 = \tilde{V} \otimes (U_{\cdot,1} \cdot U^{-1} \cdot U_{\cdot,1}) - \tilde{V} \otimes (U_{\cdot,1} \cdot U^{-1} \cdot U_{\cdot,-1})x_1 = U_{1,1}\tilde{V} - (U_{1,-1} \cdot x_1)\tilde{V}$ .

To compute  $\tilde{\mu}_{(1,1)}$ , we need to compute  $\Sigma_{(1,1),-(1,1)} \underbrace{\Sigma_{-(1,1),-(1,1)}^{-1}}_y \mu_{-(1,1),(1,1)}$ . Therefore, we have to solve  $Ay_1 = \tilde{\mu}_1 - \underbrace{(V_{1,-1} \cdot V_{-1,-1}^{-1})}_{\tilde{V}_2} \otimes (U_{-1,\cdot} \cdot U^{-1})\tilde{\mu}_2 = \tilde{\mu}_1 - (\tilde{V}_2 \otimes e_1)\tilde{\mu}_2 = \tilde{\mu}_1 - (\tilde{\mu}_2)_1 \cdot \tilde{V}_2$ .

Since  $A = (V_{1,1} - \tilde{V}) \cdot U_{-1,-1}$ , we can find  $y_1$  by solving  $U_{-1,-1}y_1 = \frac{\tilde{\mu}_1 - (\tilde{\mu}_2)_1 \cdot \tilde{V}_2}{V_{1,1} - \tilde{V}}$

Next,  $y_2 = \tilde{\Sigma}_{22}^{-1}(\tilde{\mu}_2 - \tilde{\Sigma}_{21}y_1) = (V_{-1,-1}^{-1} \otimes U^{-1})(\tilde{\mu}_2 - (V_{-1,1} \otimes U_{\cdot,-1})y_1)$  and therefore

$$\tilde{\Sigma}_{12}\tilde{\Sigma}_{22}^{-1}(\tilde{\mu}_2 - \tilde{\Sigma}_{21}y_1) = (V_{1,-1} \cdot V_{-1,-1}^{-1}) \otimes (U_{1,\cdot} \cdot U^{-1})(\tilde{\mu}_2 - (V_{-1,1} \otimes U_{\cdot,-1})y_1) =$$

$$(\tilde{V}_2 \otimes e_1)\tilde{\mu}_2 - \tilde{V} \cdot U_{1,-1} \cdot y_1 = (\tilde{\mu}_2)_1 \cdot \tilde{V}_2 - \tilde{V} \cdot U_{1,-1} \cdot y_1$$

The last equation leads to the algorithm.

## 2 Prior settings

We use the following prior settings:

- $\sigma_\beta^2 = 1$
- We set  $\Sigma$  such that the  $(i, j)$  element  $\Sigma_{i,j} = e^{\left(-\frac{1}{2l_\Sigma^2}d(s_i, s_j)^2\right)}$ , where  $d(s_i, s_j)$  is the distance between site  $i$  and site  $j$  and  $l_\Sigma$  is a scale parameter modeling the spatial autocorrelation. We have set the scale parameters  $l_\Sigma = .05$ .
- $a_\zeta = 1, b_\zeta = 50$
- $\sigma_\mu = 1, \nu_s = 1$
- $a_p = 20, b_p = 1$
- $a_q = 1, b_q = 100$
- $\mu_r = 100, \sigma_r = 100$

## 3 Simulation study settings

### 3.1 Study design simulations

For the differences in DNA biomasses, we set:

- $S = 40$
- $\tau = .5, \sigma = .5, \sigma_u = 1$
- $\beta_0 = 0$
- $\lambda_s \sim N(7, 1)$
- $r_s = 100$
- $\phi_0 \sim N(-1.5, .001)$

- $\zeta^s = .02, \nu_s = 1$
- $p = .95, q = .05$
- $\mu_0 = 5, n_0 = 5, \pi = .9$
- $\tilde{\mu} = 100$

We used the same settings for the covariate coefficients, but we also selected

$$\beta_s^z = \begin{cases} 1 & i \text{ odd} \\ 0 & i \text{ even.} \end{cases} .$$

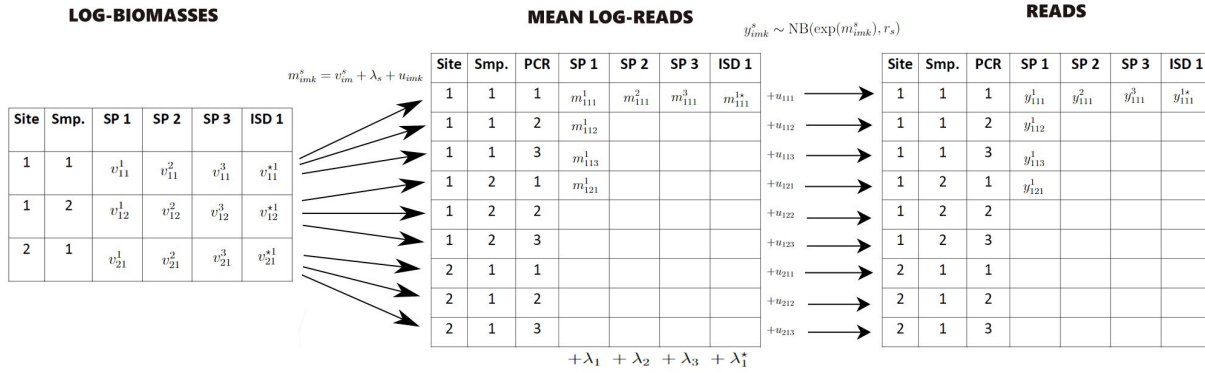
### 3.2 Spike-in simulations

For the differences in DNA biomasses, we set:

- $n = 100$
- $S = 10$
- $\sigma_u = 1$
- $\lambda_s \sim N(7, 1)$
- $r_s = 100$

For the covariate coefficients, the settings were analogous but we set  $n = 300$ .

## 4 Additional plots on sampling



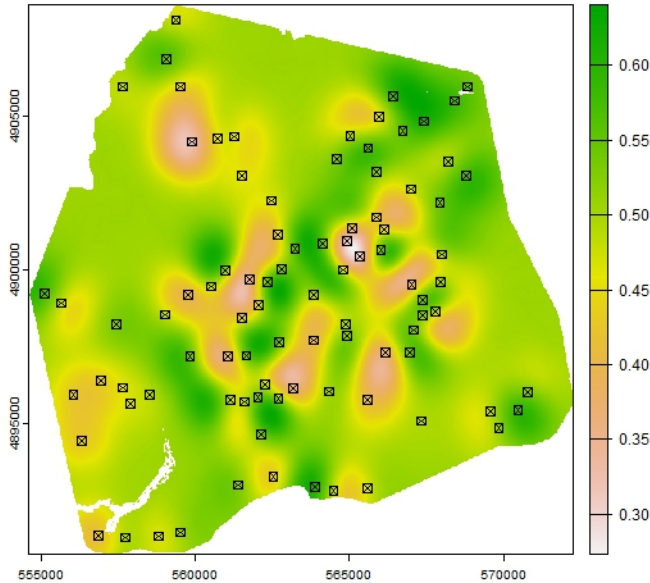
**Figure 1:** Representation of the biomass analysis stage. The number of reads (on the log scale) obtained for species  $s$ , in site  $i$ , sample  $m$ , and PCR  $k$ , is denoted by  $y_{imk}^s$  and is a function of the amount of log-biomass of that species in the corresponding sample ( $v_{im}^s$ ), of the species effect,  $\lambda_s$ , which is common across samples and PCR runs, and of PCR noise,  $u_{imk}$ , which is common across species.



## 5 Additional plots on case study

1	Arachnida Araneae Anyphaenidae Anyphaena Anyphaena pacifica
2	Arachnida Opiliones Phalangiidae Leptobunus Leptobunus parvulus
3	Chromadorea Rhabditida Aphelenchoididae Bursaphelenchus Bursaphelenchus abruptus
4	Collembola Entomobryidae Family`indet Genus`indet Species`indet
5	Insecta Blattodea Archotermopsidae Zootermopsis Zootermopsis angusticollis
6	Insecta Coleoptera Family`indet Genus`indet Species`indet
7	Insecta Coleoptera Cerambycidae Leptura Leptura obliterata
8	Insecta Coleoptera Elateridae Hypogamus Species`indet
9	Insecta Coleoptera Elateridae Megapenthes Megapenthes caprella
10	Insecta Coleoptera Scaptiidae Anaspis Anaspis rufa
11	Insecta Diptera Asilidae Genus`indet Species`indet
12	Insecta Diptera Tabanidae Hybomitra Species`indet
13	Insecta Diptera Tabanidae Hybomitra Hybomitra liorhina
14	Insecta Diptera Tachinidae Eucelatoria Species`indet
15	Insecta Diptera Scathophagidae Scathophaga Scathophaga furcata
16	Insecta Diptera Muscidae Phaonia Species`indet
17	Insecta Diptera Family`indet Genus`indet Species`indet
18	Insecta Diptera Syrphidae Hadromyia Hadromyia pulchra
19	Insecta Diptera Muscidae Spilogona Spilogona bifimbriata
20	Insecta Diptera Scathophagidae Microprosopa Species`indet
21	Insecta Diptera Empididae Genus`indet Species`indet
22	Insecta Diptera Cecidomyiidae Genus`indet Species`indet
23	Insecta Diptera Anthomyiidae Genus`indet Species`indet
24	Insecta Diptera Mycetophilidae Cordyla Species`indet
25	Insecta Diptera Phoridae Megaselia Species`indet
26	Insecta Diptera Mycetophilidae Genus`indet Species`indet
27	Insecta Diptera Keroplatidae Genus`indet Species`indet
28	Insecta Diptera Sciaridae Genus`indet Species`indet
29	Insecta Diptera Rhagionidae Genus`indet Species`indet
30	Insecta Diptera Muscidae Helina Helina troene
31	Insecta Hemiptera Reduviidae Zelus Zelus tetracanthus
32	Insecta Hemiptera Cicadellidae Osbornellus Osbornellus borealis
33	Insecta Hymenoptera Formicidae Lasius Lasius pallitarsis
34	Insecta Hymenoptera Formicidae Camponotus Camponotus modoc
35	Insecta Hymenoptera Formicidae Leptothorax Species`indet
36	Insecta Hymenoptera Vespidae Dolichovespula Dolichovespula maculata
37	Insecta Hymenoptera Vespidae Vespula Vespula alascensis
38	Insecta Hymenoptera Ichneumonidae Genus`indet Species`indet
39	Insecta Hymenoptera Vespidae Dolichovespula Dolichovespula alpicola
40	Insecta Hymenoptera Diapriidae Genus`indet Species`indet
41	Insecta Lepidoptera Geometridae Lambdina Lambdina fiscellaria
42	Insecta Lepidoptera Family`indet Genus`indet Species`indet
43	Insecta Lepidoptera Geometridae Neoalcis Neoalcis californiaria
44	Insecta Lepidoptera Geometridae Ceratodalia Ceratodalia gueneata
45	Insecta Lepidoptera Geometridae Eulithis Eulithis destinata
46	Insecta Neuroptera Hemerobiidae Hemerobius Species`indet
47	Insecta Orthoptera Rhabdiphoridae Pristoceuthophilus Pristoceuthophilus cercalis
48	Insecta Psocodea Caeciliusidae Valenzuela Species`indet
49	Insecta Psocodea Dasydemellidae Teliapsocus Teliapsocus conterminus
50	Insecta Psocodea Psocidae Loensia Loensia maculosa

**Table 1:** Species used in the case study



**Figure 2:** Biodiversity map with sampling points.

## 6 Comparison with existing methods

We compare our model with two existing approaches, a jSDM using the package *gllvm* and the two-stage occupancy model of Griffin et al. (2020).

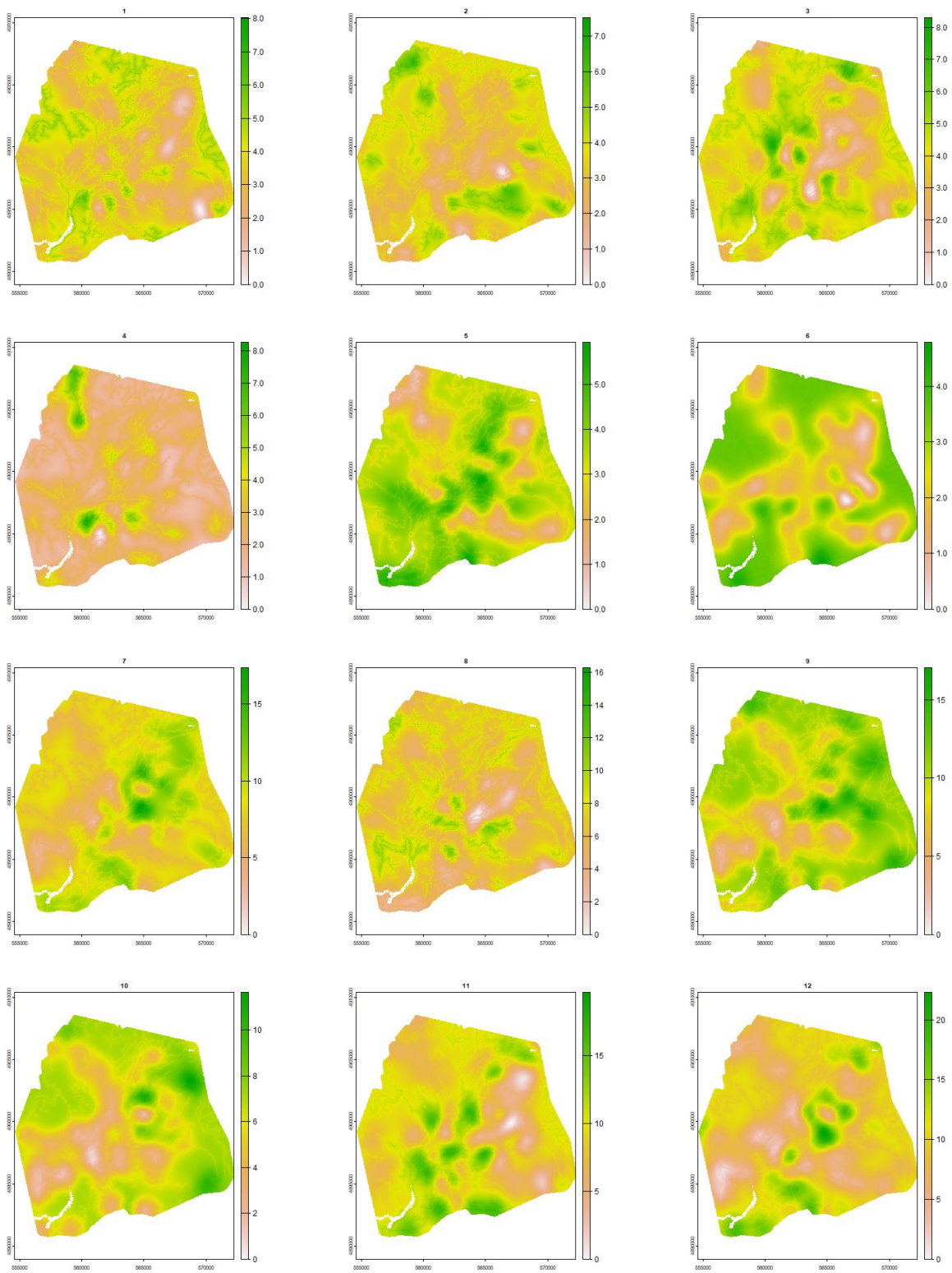
### 6.1 *gllvm*

We use the *gllvm* package to fit a jSDM model based on a factor representation of the form

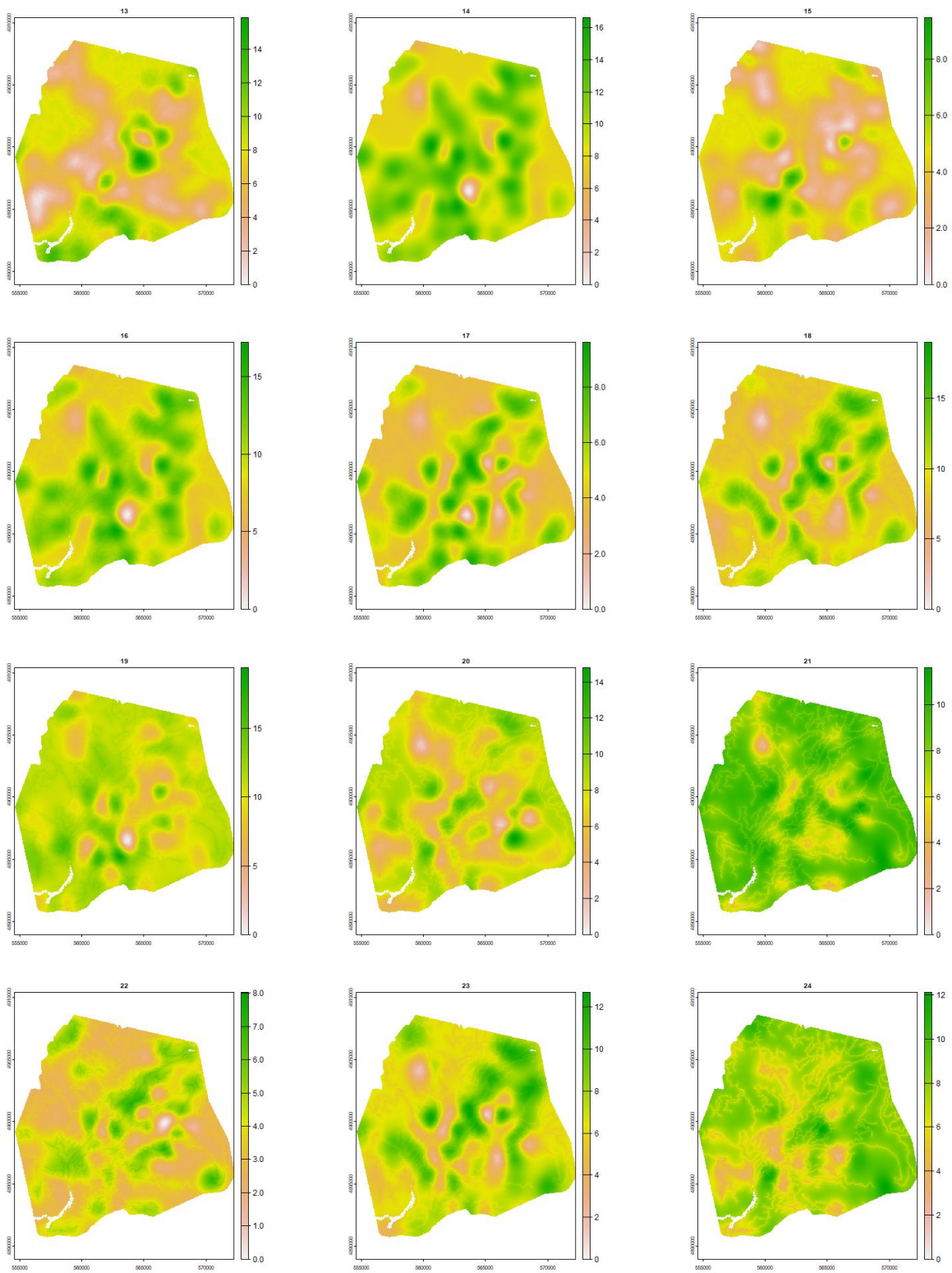
$$Y_i^s \sim \text{Pois}(\exp(X_i\beta_s + U_i.\Lambda_s))$$

where  $Y_i^s$  are the counts of species  $s$  in sample  $i$  aggregated over the 3 PCRs,  $\beta_s$  are species-specific covariate coefficients and an estimate of the between-species covariance matrix is given by  $\Lambda\Lambda^T$ .

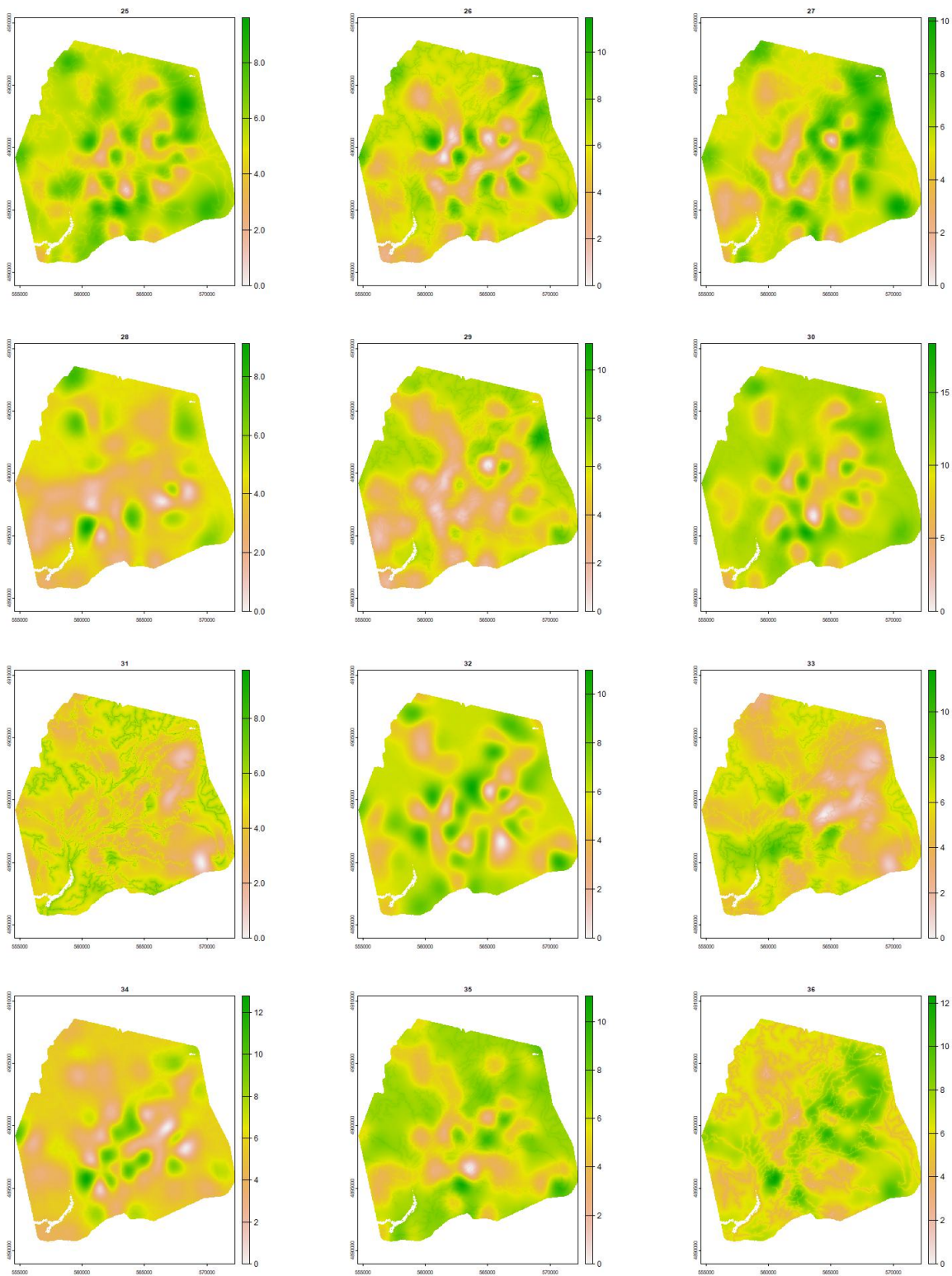
We use the same covariates as the case study: log-elevation and log-distance to road and summarise the estimated coefficients in In Fig. 8, while in Fig. 9 we report the residual covariance matrix of the observations.



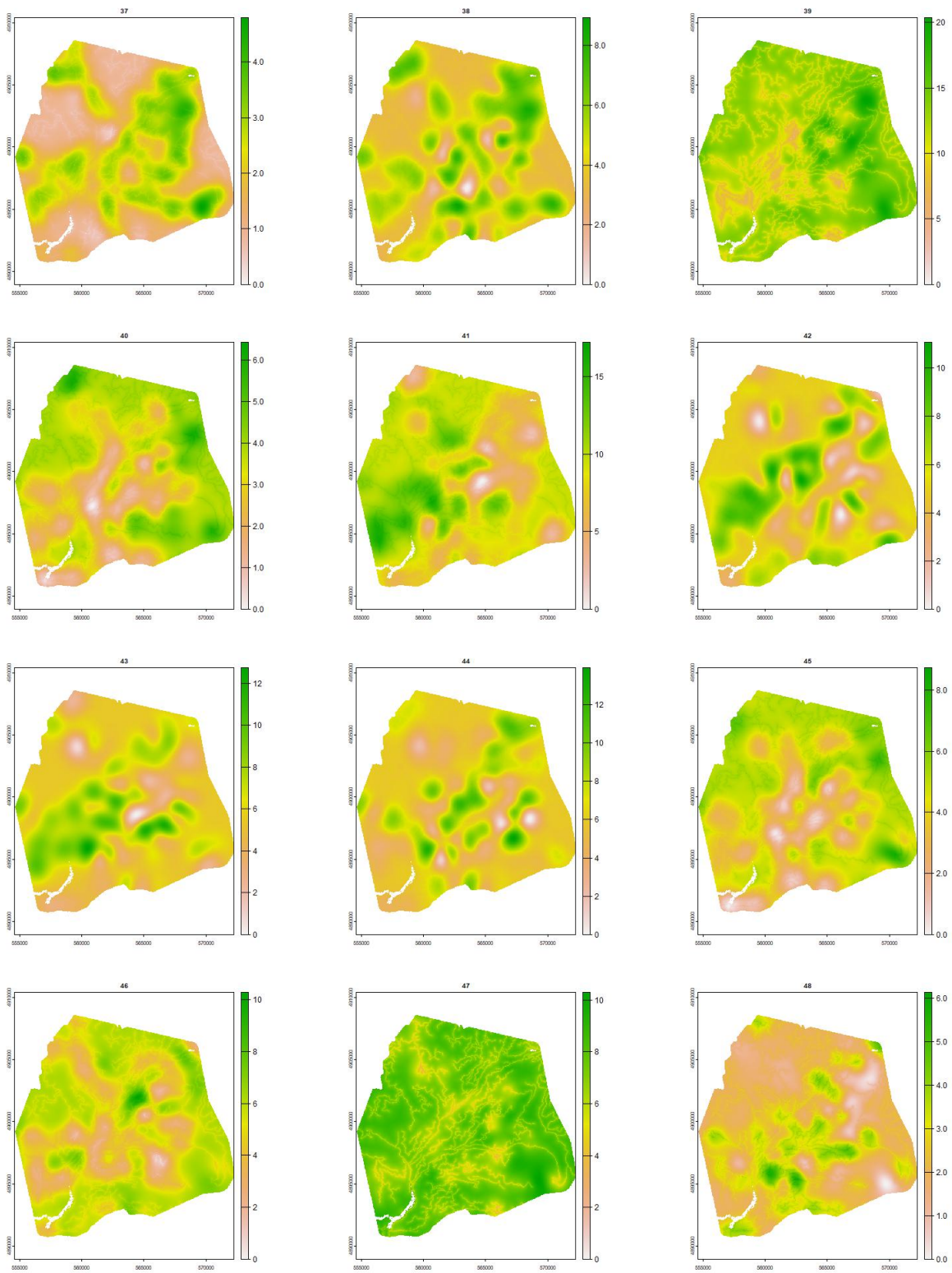
**Figure 3:** Maps of posterior mean log-biomass for species 1 to 12.



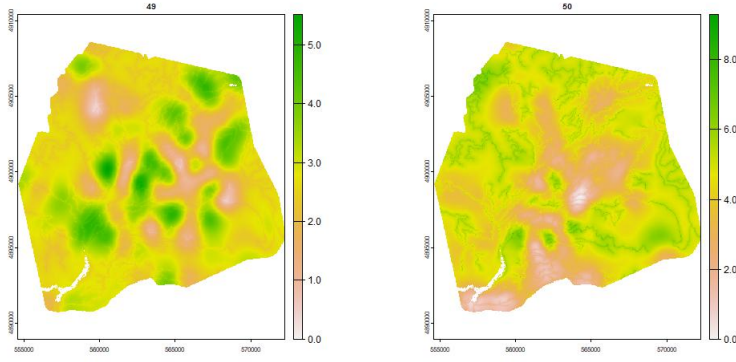
**Figure 4:** Maps of posterior mean log-biomass for species 13 to 24.



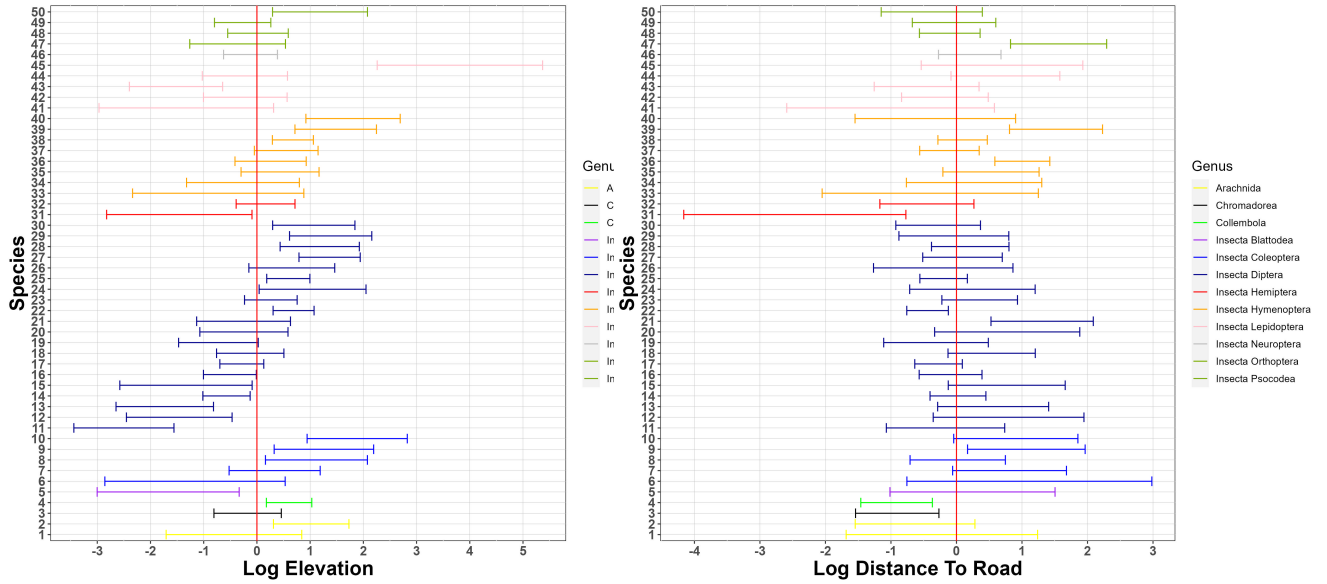
**Figure 5:** Maps of posterior mean log-biomass for species 25 to 36.



**Figure 6:** Maps of posterior mean log-biomass for species 37 to 48.



**Figure 7:** Maps of posterior mean log-biomass for species 49 to 50.



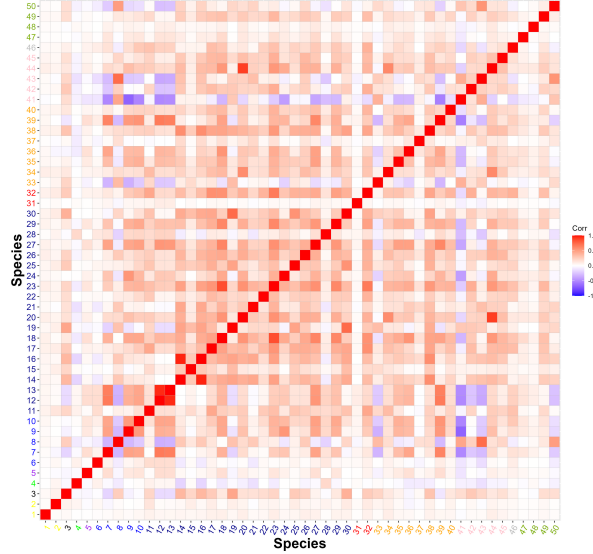
**Figure 8:** Results from gllvm: 95% PCI of the species-specific coefficients of log elevation (left) and distance to road (right) in the model for log-biomass. Species are grouped taxonomically.

## 6.2 Two stage occupancy model

We fit the two-stage occupancy model of Griffin et al. (2020):

$$\left\{ \begin{array}{l} z_i^s \sim \text{Be}(\psi_i^s) \\ \text{logit}(\psi_i^s) = X_i \beta_s \\ \delta_{im}^s \sim \text{Be}(z_i^s \theta_{im}^s + (1 - z_i^s) \theta_0) \\ \text{logit}(\theta_i^s) = \beta_s^\theta \\ y_{imk}^s \sim \text{Be}(\delta_{im}^s p_s + (1 - \delta_{im}^s) q_s) \end{array} \right.$$

where:



**Figure 9:** Results from gllvm: Correlation plot of all species. Red represents positive correlations while blue represents negative correlations. Species are grouped taxonomically.

- $\psi$  is the occupancy probability;
- $\theta$  is the true positive probability at Stage 1;
- $\theta_0$  is the false positive probability at Stage 1;
- $p$  is the true positive probability at Stage 2;
- $q$  is the false positive probability at Stage 2;

As before, we used as covariates the log-elevation and log-distance to road and summarise the estimated coefficients in Fig. 10.

## 7 Proof of study design results

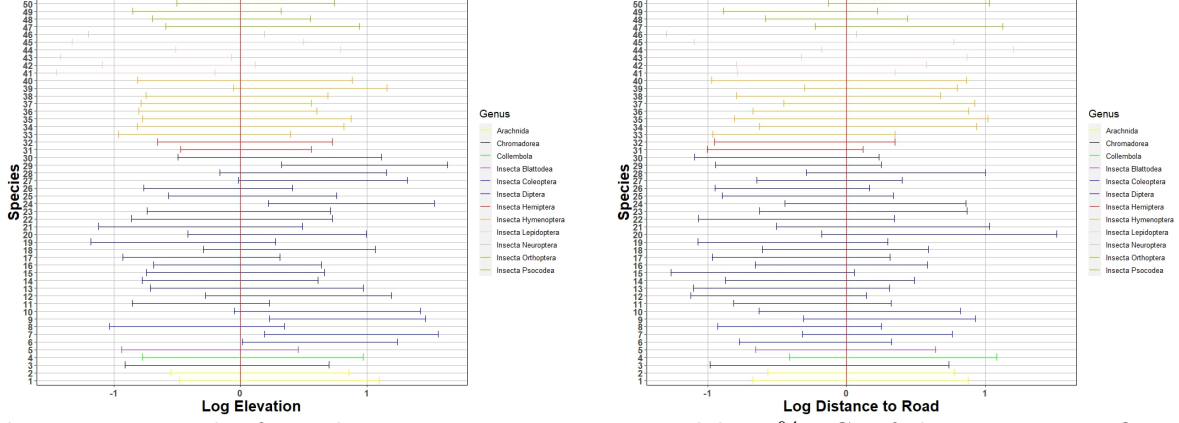
### 7.1 Lemmas

We are going to denote by  $I_n$  the identity matrix of dimension  $n$ , the  $n \times k$  matrix filled with 1's by  $A_{n,k}$  and as  $\mathbf{1}_{n,k}$  the  $n \times nk$  matrix having 1's in the positions  $(i, k(i-1) + 1, \dots, k(i-1) + k)$ ,  $i = 1, \dots, n$  and 0 everywhere else.

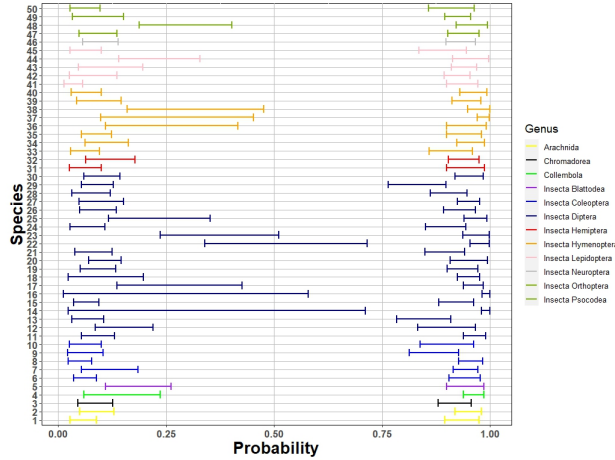
**Lemma 7.1.**

$$(a_1 I_n + a_2 A_{n,n})^{-1} = \frac{1}{a_1} I_n - \frac{a_2}{a_1(a_1 + na_2)} A_{n,n}$$





**Figure 10:** Results from the two-stage occupancy model: 95% PCI of the species-specific coefficients of log elevation (left) and distance to road (right) in the model for log-biomass. Species are grouped taxonomically.



**Figure 11:** Results from the two-stage occupancy model: false positives and false negatives rate at the lab stage. Species are grouped taxonomically.

**Lemma 7.2.** *Let*

$$\begin{cases} y_i \sim N(x, \sigma^2) & i = 1, \dots, n \\ x \sim N(\mu, \tau^2) \end{cases}$$

*Then  $(y_1, \dots, y_n) \sim N(\mu, Q^{-1})$ , where  $Q = a_Q I + b_Q A_n$ , with  $a_Q = \frac{1}{\sigma^2}$ ,  $b_Q = -\frac{\tau^2}{\sigma^2(\tau^2 n + \sigma^2)}$ .*

**Lemma 7.3.** *Let  $y_i$  be a  $k$  vector and  $x$  a scalar and*

$$\begin{cases} y_i \sim N(\mathbf{1}x, \Sigma) & i = 1, \dots, n \\ x \sim N(\mu, \tau^2) \end{cases}$$

*Then  $(y_1, \dots, y_n) \sim N(\mu, Q^{-1})$ , where  $Q = I_n \otimes \Sigma^{-1} - A_{n,n} \otimes \Sigma^{-1} a_\tau A_{k,k} \Sigma^{-1}$ , where  $a_\tau =$*

$$\frac{1}{n \sum_{i,k} \Sigma_{i,k}^{-1} + \frac{1}{\tau^2}}.$$

**Lemma 7.4.** Let  $\Sigma_i$  be a  $k \times k$  matrix and  $x$  a  $k$  vector, if

$$\begin{cases} y_i \sim N(x, \Sigma_i) \\ x \sim N(\mu, \tau^2 I) \end{cases}$$

$$\text{then } (y_1, \dots, y_n) \sim N(\mu, Q^{-1}), \text{ where } Q = \begin{bmatrix} \Sigma_1^{-1} & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \Sigma_n^{-1} \end{bmatrix} - \begin{bmatrix} \Sigma_1^{-1} \Lambda_\tau \Sigma_1^{-1} & \cdots & \Sigma_n^{-1} \Lambda_\tau \Sigma_1^{-1} \\ \vdots & \ddots & \vdots \\ \Sigma_n^{-1} \Lambda_\tau \Sigma_n^{-1} & \cdots & \Sigma_n^{-1} \Lambda_\tau \Sigma_n^{-1} \end{bmatrix}$$

$$\text{and } \Lambda_\tau^{-1} = (\sum_i \Sigma_i^{-1} + \frac{1}{\tau^2} I_k)$$

**Lemma 7.5.** Let  $\lambda$  be a  $1 \times n$  vector and  $\Sigma$  a  $nk \times nk$  matrix. If

$$\begin{cases} y \sim N(\lambda \mathbf{1}_{n,nk}, \Sigma) \\ \lambda \sim N(\mu, \tau^2 I) \end{cases}$$

then  $y \sim N(0, Q^{-1})$ , where  $Q = \Sigma^{-1} - \Sigma^{-1} \mathbf{1}^T \Lambda_\tau^{-1} \mathbf{1} \Sigma^{-1}$ , with  $\Lambda_\tau = (\mathbf{1} \Sigma^{-1} \mathbf{1}^T + \frac{1}{\tau^2} I_n)$ .

**Lemma 7.6.** Let

$$A = \begin{bmatrix} a_1 I_{n_1} & 0 & 0 \\ 0 & a_2 I_{n_2} & 0 \\ 0 & 0 & a_3 I_{n_3} \end{bmatrix} - \begin{bmatrix} b_{11} A_{n_1, n_1} & b_{12} A_{n_1, n_2} & b_{13} A_{n_1, n_3} \\ b_{21} A_{n_2, n_1} & b_{22} A_{n_2, n_2} & b_{23} A_{n_2, n_3} \\ b_{31} A_{n_3, n_1} & b_{32} A_{n_3, n_2} & b_{33} A_{n_3, n_3} \end{bmatrix},$$

Then,

$$A^{-1} = \begin{bmatrix} c_1 I_{n_1} & 0 & 0 \\ 0 & c_2 I_{n_2} & 0 \\ 0 & 0 & c_3 I_{n_3} \end{bmatrix} - \begin{bmatrix} d_{11} A_{n_1, n_1} & d_{12} A_{n_1, n_2} & d_{13} A_{n_1, n_3} \\ d_{21} A_{n_2, n_1} & d_{22} A_{n_2, n_2} & d_{23} A_{n_2, n_3} \\ d_{31} A_{n_3, n_1} & d_{32} A_{n_3, n_2} & d_{33} A_{n_3, n_3} \end{bmatrix},$$

## 7.2 Proof

In this subsection 6.1, we prove result (5) and (6) of the paper. First, we state the result.

**Proposition.** *Consider the model*

$$\left\{ \begin{array}{lll} v_{im}^s \sim N(l_i^s, \sigma^2) & i = 1, \dots, n & m = 1, \dots, M \quad s = 1, \dots, S \\ v_{im}^s = 0 & i = 1, \dots, n & m = 1, \dots, M \quad s = S + 1, \dots, S + S^* \\ u_{imk} \sim N(0, \sigma_u^2) & i = 1, \dots, n & m = 1, \dots, M \quad k = 1, \dots, K \\ \lambda_s \sim N(0, \sigma_\lambda^2) & & s = 1, \dots, S + S^* \\ y_{imk}^s \sim N(u_{imk} + \lambda_s + v_{im}^s, \sigma_y^2) & i = 1, \dots, n & k = 1, \dots, K \quad s = 1, \dots, S + S^* \end{array} \right.$$

Then

$$\text{Var}(l_1^s - l_2^s | y) = \frac{1}{M} \left( \sigma^2 + \frac{\sigma_y^2}{K} \left( 1 + \frac{\frac{\sigma_u^2}{\sigma_y^2}}{\frac{\sigma_u^2}{\sigma_y^2} S^* + 1} \right) \right). \quad (1)$$

If we assume that

$$l_i^s \sim N(X_i \beta, \tau^2) \quad i = 1, \dots, n \quad s = 1, \dots, S$$

and  $\sigma_\lambda^2 \gg \max\{\sigma_u^2, \sigma^2, \sigma_y^2\}$ , to obtain

$$\text{Var}(\beta | y) = \frac{1}{n-1} \left( \tau^2 + \frac{1}{M} \left( \sigma^2 + \frac{\sigma_y^2}{K} \right) \right) \left( 1 + \frac{\sigma_u^2}{\sigma_y^2 + (M\tau^2 + \sigma^2)K(1 + S^* \frac{\sigma_u^2}{\sigma_y^2}) + \sigma_u^2(S + S^* - 1)} \right). \quad (2)$$

## 7.3 Proof

First, we prove result (1). Given  $\lambda = (\lambda_1, \dots, \lambda_S)$ ,  $u = (u_{imk})$ ,  $l^S = (l_1^S, \dots, l_n^S)$ , we have  $p(y|l^S) =$

$$\underbrace{\int p(\lambda) \int p(u) \left( \underbrace{\int \prod_{s=1}^{S-1} p(l_i^s) \int \prod_{s=1}^S \left[ p(v_{im}^s | l_i^s) \prod_{i,m,k} p(y_{imk}^s | \lambda_s, v_{im}^s, u_{imk}) \right] dv dl^{-S}}_{(a)} \right)}_{(c)} \underbrace{\left( \prod_{s=S+1}^{S+S^*} \prod_{i,m,k} p(y_{imk}^s | \lambda_s, u_{imk}) \right)}_{(b)} du d\lambda$$

We first focus on (a). Using Lemma 7.2

$$\int p(v_{im}^s | l_i^s) \prod_k p(y_{imk}^s | u_{imk}, \lambda_s, v_{im}^s) dv_{im}^s = N(y_{im\cdot}^s | l_i^s + u_{im\cdot} + \lambda_s, Q_1^{-1})$$

where  $Q_1 = a_{Q_1} I_K + b_{Q_1} A_K$  with  $a_{Q_1} = \frac{1}{\sigma_y^2}$  and  $b_{Q_1} = -\frac{\sigma_v^2}{\sigma_y^2(\sigma_v^2 K + \sigma_y^2)}$ , so that

$$(a) = \left( \prod_{s=1}^{S-1} \prod_{i=1}^n \int p(l_i^s) \prod_{m=1}^M N(y_{im\cdot}^s | l_i^s + u_{im\cdot} + \lambda_s, Q_1^{-1}) dl_i^s \right) \prod_{i=1}^n \prod_{m=1}^M N(y_{im\cdot}^S | l_i^S + u_{im\cdot} + \lambda_s, I_K)$$

Using Lemma 7.3,  $\int p(l_i^s) \prod_{m=1}^M N(y_{im\cdot}^s | l_i^s + u_{im\cdot} + \lambda_s, Q_1^{-1}) dl_i^s = N(y_{i\cdot}^s | u_{i\cdot} + \lambda_s, Q_2^{-1})$ , where  $Q_2 = I_M \otimes Q_1 - A_{M,M} \otimes a_l(Q_1 A_{K,K} Q_1)$ , with  $a_l$  defined in the Lemma. Therefore,

$$\begin{aligned} (a) &= \left( \prod_{s=1}^{S-1} \prod_{i=1}^n N(y_{i\cdot}^s | u_{i\cdot} + \lambda_s, Q_2^{-1}) \right) \prod_{i=1}^n N(y_{i\cdot}^S | l_i^S u_{i\cdot} + \lambda_s, \underbrace{I_M \otimes I_K}_{\hat{Q}_2^{-1}}) \\ &= \left( \prod_{s=1}^{S-1} N(y_{\cdot\cdot}^s | u_{\cdot\cdot} + \lambda_s, \underbrace{I_n \otimes Q_2^{-1}}_{Q_3^{-1}}) \right) N(y_{\cdot\cdot}^S | l^S + u_{\cdot\cdot} + \lambda_s, \underbrace{I_n \otimes \hat{Q}_2^{-1}}_{\hat{Q}_3^{-1}}) \end{aligned}$$

while, (b) can be written as  $\prod_{s=S+1}^{S+S^*} \prod_{i,m,k} p(y_{i,m,k}^s | \lambda_s + u_{i,m,k}) = N(\lambda_s + u | \underbrace{I_n \otimes I_M \otimes (Q_1^*)^{-1}}_{(Q_2^*)^{-1}}, \underbrace{\quad}_{(Q_3^*)^{-1}})$ , with  $Q_1^* =$

$$\frac{1}{\sigma_y^2} I_K = a_{Q_1^*} I_K.$$

Therefore, (c) =

$$\int p(u) \left( \prod_{s=1}^{S-1} N(y_{\cdot\cdot}^s | u_{\cdot\cdot} + \lambda_s, Q_3^{-1}) \right) N(y_{\cdot\cdot}^S | l^S + u_{\cdot\cdot} + \lambda_s, \hat{Q}_3^{-1}) \left( \prod_{s=S+1}^{S+S^*} N(y_{\cdot\cdot}^s | u + \lambda_s, (Q_3^*)^{-1}) \right) du$$

Using Lemma 7.4, this is equal to  $N((0, l^S, 0) + \lambda, Q_4^{-1})$ , where

$$Q_4 = \begin{bmatrix} I_{S-1} \otimes Q_3 & 0 & 0 \\ 0 & \hat{Q}_3 & 0 \\ 0 & 0 & I_{S^*} \otimes Q_3^* \end{bmatrix} -$$

$$\begin{bmatrix} A_{S-1,S-1} \otimes (Q_3 \Lambda_u^{-1} Q_3) & A_{S-1,1} (Q_3 \Lambda_u^{-1} \hat{Q}_3) & A_{S-1,S^*} \otimes (Q_3 \Lambda_u^{-1} Q_3^*) \\ A_{1,S-1} \otimes (\hat{Q}_3 \Lambda_u^{-1} Q_3) & (\hat{Q}_3 \Lambda_u^{-1} \hat{Q}_3) & A_{1,S^*} \otimes (\hat{Q}_3 \Lambda_u^{-1} Q_3^*) \\ A_{S^*,S-1} \otimes (Q_3^* \Lambda_u^{-1} Q_3) & A_{S^*,1} (\hat{Q}_3 \Lambda_u^{-1} Q_3) & A_{S^*,S^*} \otimes (Q_3^* \Lambda_u^{-1} Q_3^*) \end{bmatrix}$$

and  $\Lambda_u = ((S-1)Q_3 + \hat{Q}_3 + S^*Q_3^* + \frac{1}{\sigma_u^2} I_{nMK})^{-1}$ . Therefore,  $p(y | l^S) \propto$

$$\int p(\lambda) N(y_{1\cdot}, \dots, y_{S+S^*\cdot} | (0, l^S, 0) + \lambda, Q_4^{-1}) d\lambda$$

and finally, using Lemma 7.5,  $p(y|l^S) \sim N((0, \mathbf{1}_{n,nMK}^T l^S, 0), Q_5^{-1})$ , where

$$Q_5 = Q_4 - Q_4 \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T \Lambda_\lambda^{-1} \mathbf{1}_{(S+S^*),nMK(S+S^*)} Q_4,$$

and

$$\Lambda_\lambda = (\mathbf{1}_{(S+S^*),nMK(S+S^*)} Q_4 \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T + \frac{1}{\sigma_\lambda^2} I_{(S+S^*)})^{-1}.$$

Having derived an expression for  $p(y|l^S)$ , we can easily derive  $Var(l^S|y)$  but it's convenient to derive an expression for  $Q_5$  in terms of  $M$ ,  $K$ ,  $n$  and  $S^*$ .

We have  $Q_2 = I_M \otimes Q_1 - A_{M,M} \otimes a_l(Q_1 A_{K,K} Q_1)$ , which after expanding can be written as  $a_{Q_2}(I_M \otimes I_K) + b_{Q_2}(I_M \otimes A_{K,K}) + c_{Q_2}(A_{M,M} \otimes A_{K,K})$  and similar relationship can be obtained for  $\hat{Q}_2$  and  $Q_2^*$ .

Next, since

$$\begin{aligned} \Lambda_u &= I_n \otimes ((S-1)Q_2 + \hat{Q}_2 + S^*Q_2^* + \frac{1}{\sigma_u^2} I_{nMK})^{-1} \\ &= (S-1)(a_{Q_2}(I_M \otimes I_K) + b_{Q_2}(I_M \otimes A_{K,K}) + c_{Q_2}(A_{M,M} \otimes A_{K,K})) \\ &\quad + (a_{\hat{Q}_2}(I_M \otimes I_K) + b_{\hat{Q}_2}(I_M \otimes A_{K,K}) + c_{\hat{Q}_2}(A_{M,M} \otimes A_{K,K})) \\ &\quad + S^*(a_{Q_2^*}(I_M \otimes I_K) + b_{Q_2^*}(I_M \otimes A_{K,K}) + c_{Q_2^*}(A_{M,M} \otimes A_{K,K})) + \frac{1}{\sigma_u^2} I_{MK}, \end{aligned}$$

we can write  $\Lambda_u^{-1}$  as  $I_n \otimes \underbrace{(a_u(I_M \otimes I_K) + b_u(I_M \otimes A_K) + c_u(A_M \otimes A_K))}_{(\Lambda_u^{-1})_0}$ .

We can write  $Q_4$  as  $P - R$ , where

$$P = \begin{bmatrix} I_{S-1} \otimes Q_3 & 0 & 0 \\ 0 & \hat{Q}_3 & 0 \\ 0 & 0 & I_{S^*} \otimes Q_3^* \end{bmatrix},$$

$$R = \begin{bmatrix} A_{S-1,S-1} \otimes I_n \otimes R_{11} & A_{S-1,1} \otimes I_n \otimes R_{12} & A_{S-1,S^*} \otimes I_n \otimes R_{13} \\ A_{1,S-1} \otimes I_n \otimes R_{21} & I_n \otimes R_{22} & A_{1,S^*} \otimes I_n \otimes R_{23} \\ A_{S^*,S-1} \otimes I_n \otimes R_{13} & A_{S^*,1} \otimes I_n \otimes R_{32} & A_{S^*,S^*} \otimes I_n \otimes R_{33} \end{bmatrix},$$

and  $R_{ij} = G_i(\Lambda_u^{-1})_0 G_j$  (where  $G_1 = Q_2$ ,  $G_2 = \tilde{Q}_2$ ,  $G_3 = Q_2^*$ ), which can be written as  $a_{r_{ij}} I_{MK} + b_{r_{ij}}(I_M \otimes A_{K,K}) + c_{r_{ij}} A_{MK,MK}$ . Therefore,

$$\mathbf{1}_{(S+S^*),nMK(S+S^*)} P \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T = \begin{bmatrix} (n \sum Q_2) I_S & 0 & 0 \\ 0 & (n \sum \hat{Q}_2) & 0 \\ 0 & 0 & (n \sum Q_2^*) I_{S^*} \end{bmatrix}$$

and

$$\mathbf{1}_{(S+S^*),nMK(S+S^*)} R \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T = \begin{bmatrix} (n \sum R_{11}) A_{S-1,S-1} & (n \sum R_{12}) A_{S-1,1} & (n \sum R_{13}) A_{S-1,S^*} \\ (n \sum R_{21}) A_{1,S-1} & n \sum R_{22} & (n \sum R_{23}) A_{1,S^*} \\ (n \sum R_{31}) A_{S^*,S-1} & (n \sum R_{32}) A_{S^*,1} & (n \sum R_{33}) A_{S^*,S^*} \end{bmatrix}$$

and hence  $\Lambda_\lambda$  can be written as

$$\begin{bmatrix} a_1^\lambda I_S & 0 & 0 \\ 0 & a_2^\lambda & 0 \\ 0 & 0 & a_3^\lambda I_{S^*} \end{bmatrix} - \begin{bmatrix} b_{11}^\lambda A_{S-1,S-1} & b_{12}^\lambda A_{S-1,1} & b_{13}^\lambda A_{S-1,S^*} \\ b_{21}^\lambda A_{1,S-1} & b_{22}^\lambda & b_{23}^\lambda A_{1,S^*} \\ b_{31}^\lambda A_{S^*,S-1} & b_{32}^\lambda A_{S^*,1} & b_{33}^\lambda A_{S^*,S^*} \end{bmatrix}$$

with  $a_1^\lambda = (n \sum Q_2 + \frac{1}{\sigma_\lambda^2})$ ,  $a_2^\lambda = (n \sum \hat{Q}_2 + \frac{1}{\sigma_\lambda^2})$ ,  $a_3^\lambda = (n \sum Q_3^* + \frac{1}{\sigma_\lambda^2})$ , and  $b_{ij}^\lambda = n \sum R_{ij}$ . Therefore,

$$\Lambda_\lambda^{-1} = \underbrace{\begin{bmatrix} c_1^\lambda I_S & 0 & 0 \\ 0 & c_2^\lambda & 0 \\ 0 & 0 & c_3^\lambda I_{S^*} \end{bmatrix}}_{K_1} - \underbrace{\begin{bmatrix} d_{11}^\lambda A_{S-1,S-1} & d_{12}^\lambda A_{S-1,1} & d_{13}^\lambda A_{S-1,S^*} \\ d_{12}^\lambda A_{1,S-1} & d_{22}^\lambda & d_{23}^\lambda A_{1,S^*} \\ d_{13}^\lambda A_{S^*,S-1} & d_{32}^\lambda A_{S^*,1} & d_{33}^\lambda A_{S^*,S^*} \end{bmatrix}}_{K_2}$$

with coefficients defined as in Lemma 7.6.

Next,  $Q_4 \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T = P \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T - R \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T$ , where

$$H_1 := P \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T = \begin{bmatrix} d_{Q_2} \mathbf{1}_{S-1,nMK(S-1)}^T & 0 & 0 \\ 0 & d_{\hat{Q}_2} \mathbf{1}_{1,nMK}^T & 0 \\ 0 & 0 & d_{Q_2^*} \mathbf{1}_{S^*,nMKS^*}^T \end{bmatrix}$$

where  $d_{Q_2} = a_{Q_2} + K b_{Q_2} + c_{Q_2} MK$  and similarly for  $d_{Q_2^*}$  and  $d_{\hat{Q}_2}$ , and

$$H_2 := R \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T = \begin{bmatrix} d_{R_{11}} A_{nMK(S-1),S-1} & d_{R_{12}} A_{nMKS,1} & d_{R_{13}} A_{nMKS,S^*} \\ d_{R_{21}} A_{nMK,S} & d_{R_{22}} A_{nMK,1} & d_{R_{23}} A_{nMK,S^*} \\ d_{R_{31}} A_{nMKS^*,S} & d_{R_{32}} A_{nMKS^*,1} & d_{R_{33}} A_{nMKS^*,S^*} \end{bmatrix}$$

where  $d_{R_{ij}} = a_{R_{ij}} + K b_{R_{ij}} + MK c_{R_{ij}}$  and similarly for the others. Therefore,  $Q_4 \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T \Lambda_\lambda^{-1} = \left( P \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T - R \mathbf{1}_{(S+S^*),nMK(S+S^*)}^T \right) (K_1 - K_2) = (H_1 - H_2)(K_1 - K_2)$ . We have

$$H_1 K_1 = \begin{bmatrix} c_1^\lambda d_{Q_2} \mathbf{1}_{S-1,nMKS}^T & 0 & 0 \\ 0 & c_2^\lambda d_{\hat{Q}_2} \mathbf{1}_{1,nMK}^T & 0 \\ 0 & 0 & c_3^\lambda d_{Q_2^*} \mathbf{1}_{S^*,nMKS^*}^T \end{bmatrix},$$

$$H_1 K_2 = \begin{bmatrix} c_1^\lambda d_{R_{11}} A_{nKS,S} & c_2^\lambda d_{R_{12}} A_{nKS,S^*} & c_3^\lambda d_{R_{13}} A_{nKS,S^*} \\ c_1^\lambda d_{R_{21}} A_{nKS^*,S} & c_2^\lambda d_{R_{22}} A_{nKS^*,S^*} & c_3^\lambda d_{R_{23}} A_{nKS^*,S^*} \\ c_1^\lambda d_{R_{31}} A_{nKS^*,S} & c_2^\lambda d_{R_{32}} A_{nKS^*,S^*} & c_3^\lambda d_{R_{33}} A_{nKS^*,S^*} \end{bmatrix},$$

$$H_2 K_1 = \begin{bmatrix} d_{Q_2} d_{11}^\lambda A_{nK(S-1),S-1} & d_{Q_2} d_{12}^\lambda A_{nK,S^*} & d_{Q_2} d_{13}^\lambda A_{nK(S-1),S^*} \\ d_{\tilde{Q}_2} d_{11}^\lambda A_{nK(S-1),S^*} & d_{\tilde{Q}_2} d_{12}^\lambda A_{nK,S^*} & d_{\tilde{Q}_2} d_{13}^\lambda A_{nK(S-1),S^*} \\ d_{Q_1^*} d_{21}^\lambda A_{nKS^*,S} & d_{Q_1^*} d_{22}^\lambda A_{nKS^*,S^*} & d_{Q_1^*} d_{22}^\lambda A_{nKS^*,S^*} \end{bmatrix},$$

and

$$H_2 K_2 = \begin{bmatrix} M_{S^*,1,1}^* & M_{S^*,1,2}^* & M_{S^*,1,3}^* \\ M_{S^*,2,1}^* & M_{S^*,2,2}^* & M_{S^*,2,3}^* \\ M_{S^*,3,1}^* & M_{S^*,3,3}^* & M_{S^*,3,3}^* \end{bmatrix} \text{ where } M_{T^*,i,j}^* = A_{nKS,T}(d_{R_{i1}} d_{1j}^\lambda S + d_{R_{i2}} d_{2j}^\lambda + d_{R_{i3}} d^{\lambda_{3j}}) S^*.$$

Implying that

$$Q_4 \mathbf{1}_{S,nKS}^T \Lambda_\lambda^{-1} = \underbrace{\begin{bmatrix} \mathbf{1}_{S,nKS}^T c_{\eta_1} & 0 & 0 \\ 0 & \mathbf{1}_{S^*,nKS^*}^T c_{\eta_2} & 0 \\ 0 & 0 & \mathbf{1}_{S^*,nKS^*}^T c_{\eta_3} \end{bmatrix}}_{J_1} - \underbrace{\begin{bmatrix} d_{\eta_{11}} A_{nKS,S} & d_{\eta_{12}} A_{nkS,S^*} & d_{\eta_{13}} A_{nkS,S^*} \\ d_{\eta_{21}} A_{nKS,S} & d_{\eta_{22}} A_{nkS,S^*} & d_{\eta_{23}} A_{nkS,S^*} \\ d_{\eta_{31}} A_{nKS,S} & d_{\eta_{32}} A_{nkS,S^*} & d_{\eta_{33}} A_{nkS,S^*} \end{bmatrix}}_{J_2}.$$

And finally,

$$J_1 H_1^T = \begin{bmatrix} (c_{\eta_1} d_{Q_2}) I_{S-1} \otimes A_{nMK,nMK} & 0 & 0 \\ 0 & (c_{\eta_2} d_{\tilde{Q}_2}) A_{nMK,nMK} & 0 \\ 0 & 0 & (c_{\eta_3} d_{Q_2^*}) I_{S^*} \otimes A_{nMK,nMK} \end{bmatrix},$$

$$J_1 H_2^T + J_2 H_1^T = \begin{bmatrix} N_{S-1,S-1,c_{\eta_1},R_{11},d_{\eta_{11}},Q_2}^* & N_{S-1,1,c_{\eta_1},R_{12},d_{\eta_{12}},\tilde{Q}_2}^* & N_{S-1,S^*,c_{\eta_1},R_{13},d_{\eta_{13}},Q_2^*}^* \\ N_{1,S-1,c_{\eta_2},R_{21},d_{\eta_{21}},Q_2}^* & N_{1,1,c_{\eta_2},R_{22},d_{\eta_{22}},\tilde{Q}_2}^* & N_{1,S^*,c_{\eta_2},R_{23},d_{\eta_{23}},\tilde{Q}_2}^* \\ N_{S^*,S-1,c_{\eta_3},R_{31},d_{\eta_{31}},Q_2}^* & N_{S^*,1,c_{\eta_3},R_{32},d_{\eta_{22}},\tilde{Q}_2}^* & N_{S^*,S^*,c_{\eta_3},R_{33},d_{\eta_{33}},Q_2^*}^* \end{bmatrix}$$

where  $N_{S_1,S_2,c_\eta,R,d_\eta,Q}^* = (c_\eta d_R + d_\eta d_Q) A_{nMKS_1,nMKS_2}$ ,

$$J_2^T H_2^T = \begin{bmatrix} A_{nMK(S-1),nMK(S-1)} t_{11} & A_{nMK(S-1),nMK} t_{12} & A_{nMK(S-1),nMKS^*} t_{13} \\ A_{nMKS^*,nMKS} t_{21} & A_{nMK^*,nMK} t_{22} & A_{nMK^*,nMK^*} t_{23} \\ A_{nMKS^*,nMKS^*} t_{31} & A_{nMKS^*,nMK} t_{32} & A_{nMKS^*,nMKS^*} t_{33} \end{bmatrix}$$

where  $t_{ij} = (S-1)d_{\eta_{i1}}d_{R_{1j}} + d_{\eta_{i2}}d_{R_{2j}} + S^*d_{\eta_{i3}}d_{R_{3j}}$ . This implies that

$$\begin{aligned}
& Q_4 \mathbf{1}_{S,nMK}^T \Lambda_\lambda^{-1} \mathbf{1}_{S,nMK} Q_4 \\
&= \begin{bmatrix} (c_{\eta_1} d_{Q_2}) I_S \otimes A_{nMK,nMK} & 0 & 0 \\ 0 & (c_{\eta_2} d_{\bar{Q}_2}) I_{S^*} \otimes A_{nMK,nMK} & 0 \\ 0 & 0 & (c_{\eta_3} d_{Q_2^*}) I_{S^*} \otimes A_{nMK,nMK} \end{bmatrix} \\
&- \begin{bmatrix} f_{11} A_{nMK(S-1),nMK(S-1)} & f_{12} A_{nMK(S-1),nMK} & f_{13} A_{nMK(S-1),nMK} S^* \\ f_{21} A_{nMK,nMK(S-1)} & f_{22} A_{nMK,nMK} & f_{23} A_{nMK,nMK} S^* \\ f_{31} A_{nMK} S^*, nMK(S-1) & f_{32} A_{nMK} S^*, nMK & f_{33} A_{nMK} S^*, nMK} S^* \end{bmatrix}.
\end{aligned}$$

From the expression for  $p(y|l^S)$ , we obtain

$$Var(l^S|y) = (\mathbf{1}_{n,nMK} (Q_5)_{nMK(S-1)+1, \dots, nMK; nMK(S-1)+1, \dots, nMK} \mathbf{1}_{n,nMK}^T)^{-1}$$

and, since

$$\begin{aligned}
& \mathbf{1}_{n,nMK} (Q_5)_{nMK(S-1)+1, \dots, nMK; nMK(S-1)+1, \dots, nMK} \mathbf{1}_{n,nMK}^T \\
&= I_n \underbrace{(MK(a_{\bar{Q}_2} - a_{R_{22}}) + MK^2(b_{\bar{Q}_2} - b_{R_{22}}) + M^2K^2(c_{\bar{Q}_2} - c_{R_{22}}))}_{a_{Q_3}} + A_{n,n} \underbrace{(-M^2K^2(c_{\eta_2} d_{\bar{Q}_2} - f_{22}))}_{b_{Q_3}},
\end{aligned}$$

we obtain  $Var(l^S|y) = c_{Q_3} I_n + d_{Q_3} A_{n,n}$ , where  $c_{Q_3} = \frac{1}{a_{Q_3}}$  and so

$$Var(l_{1S} - l_{2S}) = 2(Var(l_{1S}) - Cov(l_{1S}, l_{2S})) = 2c_{Q_3}.$$

To prove result (2), we need to compute  $p(y|X\beta^S)$ . We have that

$$\int p(\lambda) \int p(u) \underbrace{\left( \int \prod_{s=1}^{S-1} p(l_i^s) \int \prod_{s=1}^S \left[ p(v_{im}^s) \prod_{i,m,k} p(y_{imk}^s | \lambda_s, v_{imk}^s, u_{imk}) \right] dv dl^{-S} \right)}_{(a)} \underbrace{\left( \prod_{s=S+1}^{S+S^*} \int \int \prod_{i,m,k} p(y_{imk}^s | \lambda_s, u_{imk}) \right)}_{(b)} d\lambda du$$

With similar calculations to before, we obtain  $y|X\beta^S \sim N((0, (\mathbf{1}_{n,nMK}^T X\beta_s)^T, 0)^T, Q_5^{-1})$ , where  $Q_5 = Q_4 - Q_4 \mathbf{1}_{S,nMK}^T \Lambda_\lambda^{-1} \mathbf{1}_{S,nMK} Q_4$ .

Now,  $\Lambda_u^{-1} = I_n \otimes (SQ_2 + S^*Q_2^* + \frac{1}{\sigma_u^2} I_{MK})^{-1}$ . We have

$$\begin{aligned}
SQ_2 + S^*Q_2^* + \frac{1}{\sigma_u^2} I_{nMK} &= S(a_{Q_2}(I_M \otimes I_K) + b_{Q_2}(I_M \otimes A_{K,K}) + c_{Q_2}(A_{M,M} \otimes A_{K,K})) \\
&+ S^*(a_{Q_2^*}(I_M \otimes I_K) + b_{Q_2^*}(I_M \otimes A_{K,K}) + c_{Q_2^*}(A_{M,M} \otimes A_{K,K})) + \frac{1}{\sigma_u^2} I_{MK}
\end{aligned}$$



and therefore  $\Lambda_u^{-1} = I_n \otimes \underbrace{(a_u(I_M \otimes I_K) + b_u(I_M \otimes A_K) + c_u(A_M \otimes A_K))}_{(\Lambda_u^{-1})_0}$ .

From the expression for  $p(y|\beta)$ , we obtain

$$\text{Var}(\beta|y) = (X^T \mathbf{1}_{n,nMK} (Q_5)_{nMK(S-1)+1, \dots, nMK; nMK(S-1)+1, \dots, nMK} \mathbf{1}_{n,nMK}^T X)^{-1}$$

and since, with similar derivations as before, we can write  $\mathbf{1}_{n,nMK} (Q_5)_{nMK(S-1)+1, \dots, nMK; nMK(S-1)+1, \dots, nMK} \mathbf{1}_{n,nMK}^T$  as  $a_{Q_3} I_{n,n} + b_{Q_3} A_{n,n}$ , we have  $\text{Var}(\beta|y) = \frac{1}{(X^T X)_{a_{Q_3}} + (\sum X)^2 b_{Q_3}}$ . The final result can be obtained by noting that if  $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} N(0, 1)$ , then  $\mathbf{E}[X^T X] = \mathbf{E}[(\sum X_i)^2] = n$ .

## References

- J. E. Griffin, E. Matechou, A. S. Buxton, D. Bormpoudakis, and R. A. Griffiths. Modelling environmental DNA data; Bayesian variable selection accounting for false positive and false negative errors. *Journal of the Royal Statistical Society: Series C (Applied Statistics)*, 69(2):377–392, 2020.
- Y. Li, B. A. Craig, and A. Bhadra. The graphical horseshoe estimator for inverse covariance matrices. *Journal of Computational and Graphical Statistics*, 28(3):747–757, 2019.
- E. Makalic and D. F. Schmidt. A simple sampler for the horseshoe estimator. *IEEE Signal Processing Letters*, 23(1):179–182, 2015.
- N. G. Polson, J. G. Scott, and J. Windle. Bayesian inference for logistic models using Pólya–Gamma latent variables. *Journal of the American Statistical Association*, 108(504):1339–1349, 2013.
- H. Wang. Bayesian graphical lasso models and efficient posterior computation. *Bayesian Analysis*, 7(4):867–886, 2012.