

1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13  
14  
15  
16  
17  
18  
19  
20  
21  
22  
23  
24  
25  
26

DR. RAHEL SOLLMANN (Orcid ID : 0000-0002-1607-2039)

DR. MITCHELL JOSEPH EATON (Orcid ID : 0000-0001-7324-6333)

Article type : Articles

Journal: Ecological Applications

Manuscript type: Article

Running Head: Dirichlet process occupancy model

A Bayesian Dirichlet process community occupancy model to estimate community structure and species similarity

R. Sollmann,<sup>1,7</sup> M.J. Eaton,<sup>2</sup> W.A. Link,<sup>3</sup> P. Mulondo,<sup>4</sup> S. Ayebare,<sup>4</sup> S. Prinsloo,<sup>4</sup> A.J. Plumptre,<sup>5,8</sup> and D.S. Johnson<sup>6</sup>

<sup>1</sup> University of California Davis, Department of Wildlife, Fish and Conservation Biology, 1088 Academic Surge, One Shields Ave, Davis, CA 95616, USA

<sup>2</sup> U.S. Geological Survey, Southeast Climate Adaptation Science Center, N.C. State University, Raleigh, NC, USA

<sup>3</sup> US Geological Survey Patuxent Wildlife Research Center, Laurel, MD, USA

<sup>4</sup> Wildlife Conservation Society, PO Box 7487, Kampala, Uganda.

<sup>5</sup> KBA Secretariat, c/o BirdLife International, David Attenborough Building, Pembroke Street,

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1002/EAP.2249](https://doi.org/10.1002/EAP.2249)

This article is protected by copyright. All rights reserved

27 Cambridge, UK.

28 <sup>6</sup> Alaska Fisheries Science Center, National Marine Fisheries Service, NOAA, Seattle, WA  
29 98115, USA.

30 <sup>7</sup> Corresponding author. Email: [rsollmann@ucdavis.edu](mailto:rsollmann@ucdavis.edu)

31 <sup>8</sup> Current address, Wildlife Conservation Society, 2300 Southern Boulevard, Bronx, NY

32

33 Manuscript received 12 June 2020; accepted 17 August 2020; final version received 29 October  
34 2020.

35

### 36 **Abstract**

37 Community occupancy models estimate species-specific parameters while sharing  
38 information across species by treating parameters as sampled from a common distribution. When  
39 communities consist of discrete groups, shrinkage of estimates towards the community mean can  
40 mask differences among groups. Infinite mixture models using a Dirichlet process (DP)  
41 distribution, in which the number of latent groups is estimated from the data, have been proposed  
42 as a solution. In addition to community structure, these models estimate species similarity, which  
43 allows testing hypotheses about whether traits drive species response to environmental  
44 conditions. We develop a community occupancy model (COM) using a DP distribution to model  
45 species-level parameters. Because clustering algorithms are sensitive to dimensionality and  
46 distinctiveness of clusters, we conducted a simulation study to explore performance of the DP-  
47 COM with different dimensions (i.e., different numbers of model parameters with species-level  
48 DP random effects) and under varying cluster differences. Because the DP-COM is  
49 computationally expensive, we compared its estimates to a COM with a normal random species  
50 effect. We further applied the DP-COM model to a bird dataset from Uganda. Estimates of the  
51 number of clusters and species cluster identity improved with increasing difference among  
52 clusters and increasing dimensions of the DP; but the number of clusters was always  
53 overestimated. Estimates of number of sites occupied and species and community level covariate  
54 coefficients on occupancy probability were generally unbiased with (near-) nominal 95%  
55 Bayesian Credible Interval coverage. Accuracy of estimates from the normal and the DP-COM  
56 were similar. The DP-COM clustered 166 bird species into 27 clusters regarding their affiliation  
57 with open or woodland habitat and distance to oil wells. Estimates of covariate coefficients were

58 similar between a normal and the DP-COM. Except sunbirds, species within a family were not  
59 more similar in their response to these covariates than the overall community. Given that  
60 estimates were consistent between the normal and the DP-COM, and considering the  
61 computational burden for the DP models, we recommend using the DP-COM only when the  
62 analysis focuses on community structure and species similarity, as these quantities can only be  
63 obtained under the DP-COM.

64

65 **Key words:** bird point-counts, clustering, community occupancy model, dimensionality,  
66 Dirichlet process, latent groups, infinite mixture models

67

## 68 **Introduction**

69 Occupancy models (MacKenzie et al. 2002) have rapidly gained popularity in wildlife research,  
70 because they offer a means of estimating ecologically relevant parameters (species occurrence,  
71 association with covariates, colonization/extinction rates) while accounting for imperfect  
72 detection (MacKenzie et al. 2017) using relatively inexpensive species detection/non-detection  
73 data. The basic single-season single-species occupancy model has seen many modifications,  
74 including the joint modeling of multiple species in a community modeling framework (Dorazio  
75 and Royle 2005). Community models share information across species while maintaining the  
76 ability to estimate species-specific parameters by assuming that all parameters come from a  
77 common distribution. This distribution, in turn, is governed by hyperparameters, which reflect  
78 community-level patterns or processes; this model formulation is equivalent to including a  
79 species-level random effect. The community modeling approach has been combined with single-  
80 season (e.g., Zipkin et al. 2009, Sollmann et al. 2017) and dynamic (e.g., Dorazio et al. 2010)  
81 occupancy models, as well as with other hierarchical modeling frameworks such as distance  
82 sampling (Sollmann et al. 2016), or N-mixture modeling (Yamaura et al. 2016).

83 Choice of the specific distribution used to model species level parameters entails assumptions  
84 about how the community is structured. A common choice is the normal distribution, postulating  
85 that variation in parameter values across species can be described using a bell-shaped curve  
86 (Sauer and Link 2002, Kéry and Royle 2008, Zipkin et al. 2009). Particularly for data-sparse  
87 species, parameter estimates are pulled closer to the overall mean. Although the ability to derive

88 more precise parameter estimates for rarely observed species is a significant benefit of  
89 community models, this shrinkage of parameters towards the mean can mask effects that are  
90 present only in subgroups of the entire community (Pacifci et al. 2014). This problem can be  
91 circumvented by grouping species a priori and analyzing groups, rather than entire communities.  
92 This approach, however, reduces overall sample size and thus, precision of parameter estimates.  
93 Additionally, results can be sensitive to a priori grouping of species (Pacifci et al. 2014).

94 Finite mixture models, in which species are assigned probabilistically to a pre-defined number of  
95 latent groups, are an alternative to a priori grouping, and have been employed in a community  
96 modeling context (Dunstan et al. 2011, 2013). Building on the idea that communities consist of  
97 latent groups of species, Johnson and Sinclair (2017) proposed an infinite mixture approach for  
98 the joint modeling of multi-species abundance data using a Dirichlet process (DP) prior. In this  
99 approach, the number of clusters into which species group is unknown and must be estimated.  
100 Briefly, the DP consists of a base distribution from which cluster-specific parameter values are  
101 generated, and a concentration parameter  $\alpha$ , which determines the amount of clustering. In the  
102 context of community models, species are allocated to clusters based on cluster probabilities,  
103 which are generated with an algorithm governed by  $\alpha$  (for details, see Methods). All species in a  
104 cluster share the same parameter value, which serves to reduce the number of model parameters  
105 (Escobar and West 1995). Compared to normally distributed random effects, this semiparametric  
106 approach also increases the flexibility to capture patterns in parameter distribution within the  
107 community of interest (Dorazio et al. 2008). In addition, the approach provides information on  
108 community structure (number of clusters in the community), as well as the degree of similarity of  
109 species (how often two species belong to the same cluster) (Johnson and Sinclair 2017). The  
110 ability to estimate the degree of similarity in how species occurrence responds to covariates  
111 holds potential to address questions of ecological and conservation interest: the degree of  
112 similarity among species with similar functional traits can be used to quantify a community's  
113 response diversity, defined as the variation of responses to environmental change, which is a key  
114 determinant of ecosystem resilience (Mori et al. 2013). Further, estimates of similarity in habitat  
115 use can be contrasted with phylogenetic relatedness to investigate questions of coexistence and  
116 niche partitioning among closely related species, a topic of ongoing debate in ecology  
117 (Hutchinson 1959, Gotelli 2000, Graham et al. 2004).

118 In clustering algorithms, the cluster identity of objects is estimated based on multivariate data  
119 measured for each object. Clustering algorithms identify cluster identity with greater accuracy  
120 when more dimensions (i.e., more variables) are used to describe objects, as long as added  
121 dimensions contain information about clusters (e.g., Azizyan et al. 2013). Further, clustering in  
122 high-dimensional data (with 100s or 1000s of dimensions) suffers from the “curse of  
123 dimensionality” – the fact that in high dimensional space, volume expands so rapidly that data  
124 appear sparse and dissimilar, causing common clustering algorithms to be inefficient (Bellman  
125 1957, Houle et al. 2010). Given the dependency of clustering algorithms on the dimensions of  
126 the data, the performance of a community model using a DP prior to cluster species likely also  
127 depends on the dimensions of the DP process. To our knowledge, the effect of dimensionality on  
128 the ability of the DP to recover information on clustering of and similarity among objects has not  
129 been explored in the context of ecological modeling.

130 In this study, we develop a community occupancy model (COM) with a multivariate DP  
131 distribution for species level parameters (DP-COM). Using a simulation study, we first assess the  
132 model’s ability to recover community structure (number of clusters and species similarity) and  
133 estimate parameters of ecological interest in occupancy modeling (number of sites occupied and  
134 coefficients describing the relationship between occupancy probability and environmental  
135 variables). We set up the simulation to test how the dimensionality of the DP and differences  
136 among clusters affect these estimates. We then apply the DP-COM to bird survey data from  
137 Murchison Falls National Park, Uganda, to illustrate the modeling approach and its ability to  
138 address questions of species similarity. Finally, because DP priors are computationally expensive  
139 (Johnson and Sinclair 2017), tradeoffs between their use and traditional normal random effects  
140 models should be considered. We therefore compared accuracy of estimates from the DP-COM  
141 with that of a COM using a customary normal random species level effect (normal COM).

142

## 143 **Methods**

### 144 *Model development*

145 Under the hierarchical formulation (Royle and Dorazio 2008) of single-species single-season  
146 occupancy models (MacKenzie et al. 2002), whether or not a site  $j$  is occupied by the species of  
147 interest,  $z_j$ , is a Bernoulli random variable governed by occupancy probability  $\psi$ , which can be  
148 modeled as a function of site-specific covariates on an appropriate link scale  $f$  (for example, logit

149 or probit):

$$z_j \sim \text{Bernoulli}(\psi_j)$$

$$f(\psi_j) = \mathbf{X}'_j \boldsymbol{\beta}$$

152 Here,  $\boldsymbol{\beta}$  is a vector of regression coefficients and  $\mathbf{X}'_j$  is a vector with measures of the  
153 corresponding site-level covariates for site  $j$ . Sites are visited on  $k$  occasions, and binary  
154 observations of the focal species,  $y_{jk}$ , are treated as Bernoulli random variables governed by the  
155 detection probability  $p$ , which is conditional on the latent true occupancy state  $z_j$  and either  
156 adopts the value  $p_{jk}$ , when  $z_j = 1$ , or a value of 0 when  $z_j = 0$ .

$$y_{jk} \sim \text{Bernoulli}(p_{jk} z_j)$$

158 Analogous to  $\psi$ ,  $p$  can be modeled as a function of both site and occasion specific covariates.

159 To extend this to a community occupancy model, the parameters and latent variables of the  
160 model described above are further indexed by species,  $i$ . Rather than treating species-level  
161 parameters as independent, we assume that parameters come from a common distribution,  
162 governed by community (or hyper-) parameters (Dorazio and Royle 2005, Dorazio et al. 2006).  
163 This model formulation constitutes a form of information sharing, which allows us to include  
164 species with sparse data into the analysis.

165 A normal distribution is a common choice to describe species level parameters; however, this  
166 entails parametric assumptions of unimodality and symmetry in the community. In contrast, the  
167 semi-parametric DP allows fitting infinite mixture models that treat species as belonging to latent  
168 clusters and lets the data govern the specific cluster structure of the community. The DP consists  
169 of a base distribution,  $G_0$ , which generates cluster-level parameters, and a concentration  
170 parameter  $\alpha$ , which governs the amount of clustering. Under this formulation, the probability  
171 distribution for species-level parameters is a random draw from a DP [for a formal description of  
172 the DP, see Sethuraman 1994]. There are multiple means of implementing a DP; we opted for the  
173 Stick Breaking Algorithm (Sethuraman 1994), because it can be readily implemented in JAGS  
174 (Ohlssen et al. 2007). In the Stick Breaking Algorithm, cluster probabilities are generated using a  
175 sequence of auxiliary variables  $v \sim \text{Beta}(1, \alpha)$ , with mean  $E(v) = 1/(1 + \alpha)$ . The variable  $v$  can be  
176 thought of as the proportion that is broken off a stick. The proportion  $v_1$  corresponds to the  
177 probability of cluster 1,  $\pi_1$ ;  $v_2$  is the proportion broken off the remaining stick, and can be

178 translated into  $\pi_2$  by scaling it back to the size of the original stick,  $\pi_2 = v_2(1 - v_1)$ , and so forth  
 179 for the remaining clusters. The  $n$  species are then assigned to a cluster using a Multinomial( $n, \boldsymbol{\pi}$ )  
 180 distribution. If  $\alpha$  is large, only small pieces are broken off, leading to many clusters  $K$  and a  
 181 distribution of species-level parameters that approximates  $G_0$ . If  $\alpha$  is small, large pieces are  
 182 broken off, resulting in few clusters and a distribution of species-level parameters that can look  
 183 very different from  $G_0$ .

184 It is common practice (though not exclusive) in community models to ascribe separate univariate  
 185 hyperdistributions to each set of species-specific parameters. To take advantage of the  
 186 relationship between the number of dimensions of multivariate data and the ability to identify  
 187 clusters in the data, we followed Johnson and Sinclair (2017) and specified  $G_0$  in the DP-COM as  
 188 a multivariate normal ( $MVN$ ) distribution. Here, the  $MVN$  means correspond to the community  
 189 hyperparameters  $\boldsymbol{\beta}$ , which determine the distribution of parameters across clusters. Rather than  
 190 estimating the  $MVN$  means directly, we estimated them as separate fixed parameters and  
 191 parameterized the  $MVN$   $G_0$  in terms of deviations from the community mean effect,  $\boldsymbol{\delta}_i^*$ . This  
 192 allowed us to center the  $MVN$  on 0 for identifiability (Johnson and Sinclair, 2017):

$$\begin{aligned} \boldsymbol{\delta}_i^* &\sim DP(G_0, \alpha) \\ G_0 &= MVN(\mathbf{0}, \Omega) \end{aligned}$$

195 Note that this is equivalent to  $\boldsymbol{\delta}_k \sim MVN(\mathbf{0}, \Omega)$ , where  $\boldsymbol{\delta}_k$  are cluster-level deviations from the  
 196 community means. Species-level coefficients  $\boldsymbol{\beta}_i^*$  can be derived as

$$\boldsymbol{\beta}_i^* = \boldsymbol{\beta} + \boldsymbol{\delta}_{k[g[i]]},$$

198 where  $g[i]$  is the cluster identity of species  $i$ , estimated using the cluster probabilities generated  
 199 under the Stick Breaking Algorithm.

200 The number of dimensions of the  $MVN$  and thus the DP is determined by the number of  
 201 parameters that are modeled with random species effects. As an example, when the intercept and  
 202 all coefficients for  $m$  covariates are modeled as having species-level random effects, then the  
 203 multivariate DP for  $\boldsymbol{\delta}_i^*$  has  $m + 1$  dimensions. Occupancy models are composed of an  
 204 observational (detection) and an ecological (occupancy) component, and researchers are likely  
 205 interested in understanding species similarities with respect to each component separately (i.e.,  
 206 which species are ecologically similar vs which species are detected similarly). We therefore

207 specified separate DPs for each model component. Though not necessary, this choice also allows  
208 for efficient priors for  $\delta_k$  (see Simulation study below).

209

### 210 *Simulation study*

211 To evaluate the effect of the dimensionality of the multivariate DP on the model's ability to  
212 recover community structure, we set up a simulation study. We simulated occupancy and  
213 detection data for a community of  $n = 30$  species, grouped into  $K = 5$  clusters (10, 8, 6, 4 and 2  
214 species per cluster) across  $J = 35$  sampling locations and  $T = 5$  sampling occasions. We held  
215 detection probability constant across species, sites and occasions at  $p = 0.24$  but allowed for  
216 cluster-specific intercepts and coefficients in the predictor of occupancy (not adding the DP  
217 structure to the detection component made the models run faster and thus made the simulation  
218 study viable). We considered 5 scenarios of dimensionality, using 0 to 4 predictor variables for  
219 occupancy, corresponding to  $m = 1$  to 5 regression parameters (intercept and coefficient(s)) and,  
220 therefore, dimensions of the multivariate DP. Predictor variables were simulated as independent  
221 random variables following a  $Normal(0,1)$  distribution and we modeled their effect on  
222 occupancy probability using a probit link function. We set community hyperparameters  
223 (intercept followed by covariate coefficients)  $\beta = \{0, 1, -0.5, 0.5, -1\}$ . Following Johnson and  
224 Sinclair (2017), we modeled cluster-specific deviations from community level parameters,  $\delta_k$ , as  
225 a  $MVN(0, \Omega)$ , where  $\Omega = \omega^2(\mathbf{H}'\mathbf{H})^{-1}$ ,  $\omega$  determines the amount of variation among cluster-  
226 specific coefficients and  $\mathbf{H}$  is a  $J \times m$  matrix of predictors measured at each sampling site  
227 (including an intercept term). This  $MVN$  corresponds to a g-prior (Tiao and Zellner 1964), which  
228 is often used for regression coefficients, because of its property that with a single parameter,  $\omega$ , it  
229 controls the scale of variance and covariance based on the variance and correlation of predictor  
230 variables.

231 Because it is intuitive and has been shown (Johnson and Sinclair, 2017) that the differences  
232 among clusters influence how well a DP model can reproduce community structure, we  
233 considered three levels of among-cluster variation,  $\omega = 1, 2$  and 5, for each dimensionality  
234 scenario, yielding a total of 15 scenarios. We generated 50 data sets for each scenario, fitting the  
235 generated data to the above described DP-COM using the same covariates as the data-generating  
236 model.



237 We fit models in a Bayesian framework using a  $Beta(1,1)$  prior for  $p$  and priors suggested by  
238 Johnson and Sinclair (2017) for parameters of the DP component of the model, namely:  
239 (1) for the DP concentration parameter  $\alpha$ , a  $Gamma(a, b)$  prior where  $a$  and  $b$  are chosen  
240 depending on  $n$  so that  $[k] \approx 1/k$ , thus favoring smaller number of clusters (i.e., a more  
241 parsimonious model);  
242 (2) for  $\beta$ , a  $MVN$  g-prior with  $\mu = 0$  and  $\Sigma = 10,000(\mathbf{X}\mathbf{X})^{-1}$ , where  $\mathbf{X}$  is the design matrix for  
243 community-level effects and the specific multiplicative factor ensures sufficient variance to  
244 create a vague prior for our specific data.  
245 (3) For  $\omega$ , a scaled half-T distribution with  $\phi=1$  and  $df=1$ , which corresponds to a half-Cauchy  
246 prior distribution.

247 We simulated and analyzed data using the software R version 3.5.1 (R Core Team 2018). We fit  
248 models in JAGS 4.3.0 (Plummer 2003), accessed through the R package jagsUI 1.5.0 (Kellner  
249 2019). We ran three parallel chains with 30,000 iterations of which we discarded 10,000 as burn-  
250 in. We thinned chains by 10 to reduce output size. We used the posterior mean as a point  
251 estimate, except for the number of clusters ( $K$ ) and the number of sites occupied by species  $i$  ( $N_i$ ,  
252 derived as  $\sum_{j=1}^J z_{ij}$ ), for which we used the posterior mode (a more representative quantity in  
253 skewed posterior distributions typical for positive integer variables with small values). From  
254 model output we derived species-specific occupancy coefficients  $\beta_i^*$ . We further calculated  
255 pairwise species clustering rates as the proportion of MCMC iterations in which two species  
256 were estimated to be in the same cluster. This  $n \times n$  matrix can also be viewed as a species  
257 similarity matrix with respect to occupancy coefficients. We used the similarity matrix to  
258 calculate true and false pairwise clustering rates: first, we constructed a species-by-species  
259 matrix from the simulated data, in which species pairs received an entry of 1 if they were in the  
260 same cluster, and an entry of 0 otherwise. Then, we calculated the average pairwise clustering  
261 rate from the model output for all true species pairs (i.e., pairs with an entry of 1 in the data  
262 matrix) as true clustering rate, and the average pairwise clustering rate for all false species pairs  
263 (pairs with an entry of 0 in the data matrix) as false clustering rate.

264 We assessed convergence of parallel chains using the Gelman-Rubin statistic, R-hat (Gelman  
265 and Hill 2006). However, this statistic was not devised for a DP-type mixture model in which  
266 cluster labels switch (i.e., cluster 1 does not have the same identity throughout all iterations), and

267 as a result, cluster level parameters also switch. We were therefore more liberal in our  
268 assessment of convergence. We considered that we had achieved convergence when all structural  
269 parameters ( $\alpha$ ,  $\omega$ ,  $p$ ,  $\beta$ ) as well as all species-level coefficients,  $\beta_i^*$ , had an R-hat value  $<1.5$  and  
270 excluded iterations that did not meet this criterion. We inspected chain plots for several cases of  
271  $1.1 < \text{Rhat} < 1.5$  and found that generally, parallel chains fluctuated around the same average value,  
272 but that mixing was poor. Because these models are time intensive, we opted against running  
273 chains for more iterations, as this would have made the simulation study unfeasible.

274 To evaluate the performance of the DP model under the different scenarios, we calculated  
275 median bias (absolute bias,  $\hat{x} - x$ , for  $\beta$  and  $\beta_i^*$ , because true values were often close/equal to 0;  
276 relative bias,  $(\hat{x} - x)/x$  for all other parameters), median coefficient of variation (CV; posterior  
277 standard deviation divided by point estimate), median true and false clustering rates, and 95%  
278 Bayesian Credible Interval (BCI) coverage (percentage of iteration in which the 95% BCI  
279 included the true parameter value; henceforth just coverage) across all iterations that reached  
280 convergence. We used the median across iterations rather than the mean, because for some  
281 parameters, particularly the number of clusters  $K$ , the distribution across iterations was highly  
282 skewed, most likely due to poor identifiability particularly in scenarios with low  $\omega$  and  $m$ .

283 To evaluate whether we lose accuracy in parameter estimates when using the normal-COM on a  
284 clustered community, we also ran a normal COM using the same data generated under the 15  
285 scenarios described above and compared median bias and CV of estimates of  $N_i$  and  $\beta_i^*$  between  
286 the two approaches. We used the same g-priors for  $\beta$  and  $\delta_k$  (which correspond to  $\delta_i^*$  in the  
287 normal COM where each species forms its own cluster) and half-Cauchy prior on  $\omega$ , the same  
288 MCMC settings and applied the same convergence criteria as for the DP-COM.

289

#### 290 *Application: Bird survey data from Uganda*

291 Avian point-count data were collected from Murchison Falls National Park (MFNP) in western  
292 Uganda. The park covers nearly 4,000 km<sup>2</sup> in East Africa's Albertine Rift Valley, an area  
293 containing the highest vertebrate biodiversity on the African continent (Plumptre et al. 2007).  
294 Elevations in MFNP range from 620 m at the shore of Lake Albert to nearly 1,300 m in the  
295 southeast. The park experiences two rainy seasons (March – June and August – November), with  
296 an average annual rainfall of 1,100 mm.

297 Between 2010 and 2011, the Wildlife Conservation Society conducted bird surveys at elevations  
298 ranging between 650 and 720 m. The Ugandan government recently granted access to MFNP for  
299 oil exploration, and bird survey transects were established relative to the location of oil drilling  
300 platforms with the goal of evaluating the effects of drilling activities on bird populations. The  
301 survey area contained a mosaic of grasslands, dense and open borassus palm (*Borassus*  
302 *aethiopum*) woodland, dense and open woodland, and bush habitat. Transects measuring 2000 m  
303 were located in an easterly or westerly direction on either side of four oil-well pads (Appendix  
304 S1: Figure S1). Twenty-one point-count locations were established along each transect. The first  
305 point was located adjacent to the pad perimeter fence and subsequent points were spaced every  
306 100 m. Transects were visited on average once every 2.5 days. Following a 2-minute rest period  
307 upon arrival at a point-count location, the survey team leader (accompanied by 1 or 2 assistants)  
308 recorded all birds seen or heard over a 5-minute period, within an estimated radius of 500 m.  
309 Data collected included time of day, number of each bird species detected, estimated distance to  
310 observer, elevation of point-count location and habitat type. Surveys took place between  
311 February and September of 2010 and March to June of 2011. We selected a subset of the data  
312 that included 62 survey dates between 22 February to 4 May 2010 (corresponding to the early  
313 wet season). During that time, 149 points were visited at least once, with a mean of 23.3 (SD =  
314 3.2) visits per point, resulting in 3464 visits across all points. We assumed that bird populations  
315 were demographically closed during this period.

316 For each sampling location, we classified habitats into a binary variable of either open habitat  
317 (grassland, bush; 72 locations) or woodland (Borassus and other woodland; 77 locations) and  
318 determined the distance to the nearest oil well. In addition, for each visit, we had information on  
319 observer experience. This was evaluated qualitatively by the lead field investigator (AJP) based  
320 on years of experience, ability to identify species by call and accuracy in determining number of  
321 individuals and distance from observation point. Although all observers were competent in  
322 species identification, there was variation in experience and lead observers were ranked from 1  
323 to 3, as most to least experienced, respectively.

324 To construct a species level detection-non-detection matrix, we considered each visit a sampling  
325 occasion and reduced observations to binary species-level detection-non-detection data. We  
326 excluded species from analysis that had fewer than 5 observations, resulting in a data set of 166

327 species. Species were encountered, on average, during 121 (SD 206) visits, at 39 (SD 35)  
328 sampling locations.

329 We included the binary habitat information (open versus closed) and scaled distance to oil well  
330 as covariates on occupancy probability. Detection probability was modeled as a function of the  
331 experience of the survey team leader; because implementing the DP community occupancy  
332 model was very time consuming and our dataset had many occasions, we calculated the average  
333 experience score of a site across all visits to avoid having to model detection probability as  
334 varying by occasion. The resulting values were almost binary (either 2 or >2); we therefore  
335 included average experience as a binary covariate on detection probability (2 = intermediate  
336 experience; >2 = high experience). We modeled occupancy intercept and regression coefficients  
337 as species specific, with a multivariate DP (see below). We modeled the detection intercept with  
338 a normal random effect and the effect of observer experience on detection as fixed across all  
339 species. We opted for a normal random effect in the detection component because our  
340 simulations indicated that a univariate DP performed poorly at estimating the cluster structure of  
341 a community (see Results). Our model ignored the potential spatial autocorrelation in occupancy  
342 stemming from the surveys recording birds up to 500 m from survey points spaced 100 m apart.  
343 In practice, 94% of all observations in our datasets were within 200 m from the survey point, and  
344 78% were within 100 m. As this case study serves to demonstrate the DP-COM, rather than as an  
345 in-depth analysis of bird community ecology, we felt comfortable with the choice to ignore  
346 spatial autocorrelation.

347 For parameters of the occupancy component, we used the same priors as described for the  
348 simulation study, except that we set the multiplicative factor for the g-prior on  $\beta$  to 100,000 (to  
349 avoid overly low values in the prior variance-covariance matrix). We used a *Normal*(0, 10) prior  
350 on the mean and a *Gamma*(0.01, 0.01) prior on the standard deviation of the normal random  
351 effect on the intercept of  $\text{probit}(p)$ . To improve computational speed, we used an upper bound of  
352 100 for  $K$ . Imposing an upper bound on  $K$  is an accepted approximation of the infinite-mixture  
353 DP as long as it is set sufficiently high (Reich and Bondell 2011). Upper 95 BCI limits for the  
354 estimate of  $K$  were well below 100 (see Results), suggesting our choice of this upper limit did  
355 not affect estimates.

356 We implemented the models using the same software as for the simulation study. We ran three  
357 parallel chains with 50,000 iterations, of which we discarded 20,000 as burn-in. We thinned the  
358 remaining iterations by 30 (to avoid unwieldy model output). All model parameters except 5  $\delta_k$   
359 and 1  $\beta$  had  $R_{\text{hat}} < 1.5$  and in spite of these convergence issues, all species-specific  $\beta_i^*$  had  $R_{\text{hat}}$   
360  $\leq 1.1$ . As we focus on species-level parameters and species similarity, these convergence  
361 problems should not impact our inference. Running this model took about 5.5 days on an IBM  
362 HS22 virtual BladeCenter server with an allocation of 3 logical cores using Intel Xenon E5645  
363 processors at 2.4GHz and 1 GB RAM running ESXi. Further, we fit a normal-COM with both  
364 covariates to the data and compared estimates of  $\beta_i^*$  and  $N_i$ .

365 Finally, we explored the information provided by the DP-COM on bird community structure:  
366 First, to provide context for the amount of clustering suggested by the DP models, we compared  
367 the estimated number of clusters as well as the average pairwise clustering rate across the  
368 community to the respective expected values if species clustered at random. We generated these  
369 numbers by simulating draws from a Multinomial distribution with  $K=100$  categories and equal  
370 cell probabilities ( $\pi = 1/K$ ). The number of categories with at least one species corresponds to  
371 the number of random clusters. For each simulated set of cluster identities, we constructed a  
372 pairwise species clustering matrix, as described above. We simulated 3,000 such multinomial  
373 draws. We present the mean, SD, and range for the number of clusters; and the average (across  
374 all species) proportion of iterations that two species fell in the same cluster. Further, we  
375 contrasted average pairwise clustering rate of all families with at least 5 member species against  
376 community-wide average pairwise clustering rate, to investigate whether closely related taxa  
377 tended to respond to covariates more or less similarly than the entire community.

378

## 379 **Results**

380 Species in the simulated data sets occupied, on average (across species, iterations and scenarios),  
381 17.06 of the 35 sites (average range: 9.23 – 24.87). They were detected, on average, 20.43 times  
382 (average range: 8.75 – 34.17) and at 12.71 sites (average range: 5.98 – 20.01). Across scenarios,  
383 for the DP-COM we excluded between 1 and 13 of the 50 iterations due to convergence  
384 problems; the number of excluded iterations increased with increasing number of parameters  $m$   
385 and decreasing among-cluster variation  $\omega$ . In comparison, for the normal COM we excluded

386 between 0 and 9 iterations due to convergence issues (Appendix S2: Table S1). When using the  
387 customary cut-off of  $R_{hat} \leq 1.1$  for the normal COM, this number rose to between 0 and 23  
388 iterations (Appendix S3).

### 389 *Community structure, species similarity in simulated communities*

390 Bias in estimates of  $K$  generally decreased with increasing  $\omega$  (i.e., with increasing variation  
391 among clusters) and  $m$  (i.e., dimensions of the DP) (Figure 1a). For  $\omega = 1$ ,  $K$  (true value of 5)  
392 was consistently underestimated, with a median estimate of 1 ( $m = 1$ ) to 2 ( $m > 1$ ) clusters. For  
393 almost all other scenarios,  $K$  was overestimated, up to  $\hat{K} = 12$ . At  $\omega = 5$  and  $m \geq 2$ , the median  
394 estimate of  $K$  was consistently at 7, but variability in estimates across iterations decreased with  
395 increasing  $m$ . Precision of estimates of  $K$  increased with increasing  $\omega$  (from a maximum CV of  
396 2.38 at  $\omega = 1$  to a minimum of 0.29 at  $\omega = 5$ ). There was no evident relationship of the CV with  
397  $m$ . Finally, coverage was nominal or near nominal (at least 93%) for all scenarios with  $\omega < 5$  but  
398 dropped to between 80 and 88% for  $\omega = 5$  and  $m \geq 2$  (Appendix S2: Table S2).

399 The rate at which two species were correctly estimated as being in the same cluster (true  
400 clustering rate) ranged from 42% to 76% across scenarios (Figure 1b; Appendix S2: Table S2).  
401 The highest true clustering rate was attained at  $\omega = 1$  and  $m = 1$ ; however, the rate at which two  
402 species were incorrectly estimated to be in the same cluster (false clustering rate) for that  
403 scenario was almost as high (72%), consistent with an average estimate of  $K=1$  for this scenario.  
404 For most other scenarios, true clustering rate was  $<60\%$ . Only for  $\omega = 5$  did the true clustering  
405 rate tend to increase with increasing  $m$ , and within  $\omega = 5$ , only for  $m \geq 3$  did the correct  
406 clustering rate exceed 60%. False clustering rate decreased with increasing  $\omega$ , ranging from 36%  
407 to 72% at  $\omega = 1$ , from 21% to 38% at  $\omega = 2$ , and from 6% to 18% at  $\omega = 5$ . Only for  $\omega = 5$  did  
408 the false clustering rate continuously decrease with increasing  $m$ .

### 409 *Occupancy in simulated communities*

411 Across all scenarios, the number of sites occupied by a given species was estimated without bias  
412 (Figure 2a), though in some rare species-iteration combinations, bias reached 100%. The median  
413 CV of the number of sites occupied ranged from 9% to 15%; values decreased slightly with  
414 increasing  $m$  and decreasing  $\omega$ . (Figure 2a) The incidence of extreme CVs (at or above 100%)  
415 for specific species-iteration combinations increased with increasing  $\omega$ . Coverage was nominal

416 for all scenarios (Appendix S2: Table S3).

417 Estimates of community level regression coefficients  $\beta$  showed low to moderate absolute bias.

418 For example, depending on the scenario, the median estimate of the community intercept (true  
419 value of 0) ranged from -0.05 to 0.15, with most scenarios having median estimates  $<|0.10|$ .

420 There were no apparent patterns in bias with respect to  $m$ , but bias tended to increase with

421 increasing  $\omega$ . Estimates were less precise (i.e., had higher CVs) with increasing  $m$ , except for the

422 community intercept. Coverage was nominal for all parameters and scenarios (Appendix S2:

423 Table S4).

424 Median bias (across species and iterations) in species-specific regression coefficients  $\beta_i^*$  was

425 low to moderate. Median bias, as well as the incidence (i.e., particular species-iteration

426 combinations) of strong bias, increased with increasing  $\omega$  and increasing  $m$ , though the latter was

427 less pronounced. CVs increased with increasing  $\omega$ , with the exception of the intercept, where

428 CVs decreased with increasing  $\omega$ . There was no discernable relationship between CVs and  $m$ .

429 Coverage ranged from 89% to 97% and increased with increasing  $\omega$  (see Appendix S2: Figure

430 S1 for an example plot and Table S5 for details of simulation results).

431

#### 432 *Comparison with normal COM in simulated communities*

433 Bias and CV for estimates of  $N_i$  were very similar across the DP and the normal COM (Figure 2,

434 Appendix S2: Table S6); across scenarios, the DP model tended to have lower median CVs but

435 only by 1 or 2 percentage points. For  $\beta$ , median bias was similar between both models across

436 parameters and scenarios, but parameters from the DP model had considerably higher CVs

437 (Appendix S2: Figure S2, Table S7). For  $\beta_i^*$ , both bias and CVs were very similar between the

438 two modeling approaches (Appendix S2: Figure S3, Table S8). These patterns were the same

439 when applying the  $R_{hat} \leq 1.1$  cut-off to the normal COM results (Appendix S3).

440

#### 441 *Other parameters: $\omega$ , $\alpha$ and $p$ in simulated communities*

442 Detection probability  $p$  was estimated with minimal bias (-2% to 1%), 4-5% CV and BCI

443 coverage between 86% and 97% (Appendix S2: Table S9). Median estimates of the DP

444 concentration parameter  $\alpha$  ranged from 4.29 ( $\omega = 5$  and  $m = 5$ ) to 12.64 ( $\omega = 2$ ,  $m = 2$ ) (Appendix

445 S2: Table S10).

446 Estimates of  $\omega$  were most biased for  $\omega = 1$  (-28% to -84%). For all other scenarios, relative bias  
447 was low to moderate, ranging from -2% to 16%. The CV of  $\omega$  increased with increasing  $\omega$  and  
448  $m$ . Coverage ranged from 87% to 100%, except for  $\omega = 1$  and  $m = 1$ , where coverage was 0%  
449 (Appendix S2: Table S11).

450

#### 451 *Bird case study*

452 For 5 out of 300  $\delta_k$  and one  $\beta$ ,  $R\text{-hat} > 1.5$ ; however, all (derived) species-specific regression  
453 coefficients  $\beta_i^*$  had  $R\text{-hat} < 1.1$ . We visually checked chains for the non-converged  $\delta_k$  and  $\beta$ ,  
454 which appeared to be strongly autocorrelated but oscillated around the same average value; we  
455 therefore felt confident to use the estimates.

456 For the occupancy component of the model (with a multivariate DP for the coefficients of the  
457 probit-linear predictor of occupancy probability), species comprised 27 clusters (SD = 4.16, 95%  
458 BCI 22 - 37; Figure 3). Probabilities of two species clustering together ranged from 0 to 0.92.  
459 The estimate of  $\omega$  for the full model was 8.60 (SD 1.04, 95BCI 6.85 – 10.98), indicating  
460 considerable variation in regression coefficients among clusters.

461 The data set contained ten families with at least 5 species, comprising 90 species total. When  
462 looking at pairwise clustering rates for these families, we found that most families showed  
463 clustering probabilities similar to those of the entire community. However, the sunbirds  
464 (Nectariniidae, 5 species) had considerably higher clustering probabilities, whereas the  
465 Cisticolidae (12 species) and the bee-eaters (Meropidae, 5 species) had lower clustering  
466 probabilities (Figure 4).

467 Species were estimated to occupy between 1 and 147 of the 149 sample sites. We observed  
468 strong effects (i.e., with 95% BCI not overlapping 0) of woodland habitat for 57 species, with 24  
469 negative and 33 positive coefficients. For distance from oil well, 14 species showed strong  
470 negative and 12 species showed strong positive effects (Figure 5). Species with positive  
471 associations with woodland habitat tended to have positive associations with distance to oil as  
472 well (52 species), and vice versa (72 species).

473 When comparing estimates of  $\beta_i^*$  and  $N_i$  between the DP-COM and the normal COM, both  
474 modeling approaches produced very similar results (Appendix S1: Figure S2).



475

## 476 **Discussion**

477 In wildlife research, DP priors have been used to model genetic population structure (Reich and  
478 Bondell 2011), spatial variation in abundance (Dorazio et al. 2008, Dorazio 2009), spatial  
479 clustering of population trends (Johnson et al. 2013), and clustering of species with respect to  
480 habitat coefficients in the context of community distribution models (Johnson and Sinclair 2017).  
481 Our simulation study showed that a community occupancy model with a DP, instead of the  
482 customary normal random species effect, was able to retrieve aspects of community structure  
483 when differences among clusters and the number of parameters making up the multinomial DP  
484 were sufficient. Applied to data for a bird community, the model led to a considerable reduction  
485 in the number of parameters estimated, grouping 166 species into 27 clusters. This suggests that  
486 detection/non-detection data contain information on the similarity of species in a community that  
487 can be exploited with a DP model. Major shortcomings of the approach were its computational  
488 expense, poor mixing and difficulty with convergence of MCMC chains due to label switching  
489 among clusters, and its reduced performance in retrieving community structure when cluster  
490 parameters were similar and/or few parameters were used in the DP. These drawbacks may  
491 appear particularly off-putting given that there are no a priori tests that would indicate whether  
492 the existence of, and differences among, clusters warrant exploring a “costly” DP-COM.  
493 Moreover, the customary model with a normal random effect performed similarly to the DP-  
494 COM, even when applied to data from a clustered community, suggesting that a normal random  
495 effect is flexible enough to capture variation in parameters that do not follow a normal  
496 distribution. For analyses focused on community and species-level responses in occurrence  
497 (and/or detection) to covariates, or simply the estimation of occupancy probabilities in the  
498 absence of covariates, we recommend the more efficient normal COM. Only the DP-COM,  
499 however, returns estimates of community structure and species similarity in their response to  
500 covariates; for analyses aimed at testing hypotheses regarding these measures, the additional  
501 time investment needed to fit a DP-COM seems worthwhile.

502

### 503 *Factors affecting the performance of the DP-COM*

504 We found that both the variability among clusters and the dimensionality of the DP affected the  
505 ability of the model to retrieve information on community structure. Median bias in  $K$ , the

506 incidence of large bias and the incidence of large CVs all declined with increasing number of  
507 dimensions of the DP; when variation among clusters was high ( $\omega = 5$ ), increased dimensionality  
508 also led to higher true and lower false clustering rates. Across levels of among-cluster variation,  
509 univariate DPs did the worst in terms of clustering rates and estimating  $K$ . All of this indicates  
510 improved ability of the model to identify cluster identity of species with increased  
511 dimensionality. Estimates of community structure may not be reliable when only a single  
512 dimension is considered. As such, the DP-COM may be more useful for data sets with sufficient  
513 replication to support modeling of multiple covariates. It is possible, however, that if variation  
514 among clusters is stronger than what we considered in the simulation, a univariate DP may be  
515 able to identify clusters. Even though the effect of dimensionality on the performance of  
516 clustering algorithms is known (e.g., “curse of dimensionality”; Bellmann, 1957) and the DP is a  
517 widely used Bayesian clustering algorithm outside of wildlife research, to our knowledge this is  
518 the first study to demonstrate that the performance of the DP model is dependent on the number  
519 of dimensions of the base distribution.

520 Not surprisingly, we found that the variability among clusters strongly affected the ability of the  
521 DP-COM to estimate the number of clusters in the community, as well as pairwise species  
522 clustering rates. While increasing  $\omega$  resulted in higher true clustering rates, lower false clustering  
523 rates and lower bias in  $K$ , it also resulted in increased bias and CV in most  $\beta$  and  $\beta_i^*$  and higher  
524 incidences of extreme bias and CVs in  $N_i$ . There appears to be a trade-off between improvements  
525 in estimation of community structure and species similarities as a function of cluster  
526 discrimination and the accuracy of other parameters of ecological interest. Regardless, coverage  
527 of these parameters was nominal or near nominal across scenarios.

528 At  $\omega = 1$ , the DP-COM was essentially unable to detect cluster structure and, in most iterations,  
529 estimated that all species belonged to the same cluster (regardless, estimates of  $\beta_i^*$  and  $N_i$  were  
530 largely unbiased). Further, in our simulation, the actual number of clusters was, on average, not  
531 estimated well (median bias was mostly  $>40\%$ ), and coverage of the true value was  $<90\%$  for  
532 scenarios that estimated  $K$  with the lowest bias (i.e.,  $\omega = 5$  and  $m \geq 2$ ). Both findings contradict  
533 results by Johnson and Sinclair (2017), whose proof of concept simulation for a community  
534 Poisson regression resulted in accurate estimates of  $K$  for various values of  $\omega$ , as long as  $\omega > 0.5$ .  
535 We implemented the DP on parameters of the occupancy component of the DP-COM, which is

536 binary and partially latent (only for sites where the species is detected, do we observe occupancy  
537 state). It is conceivable that the differences between clusters need to be more pronounced, and/or  
538 that it is generally more difficult for the DP algorithm to retrieve community structure for a  
539 binary partially latent process. Based on these results, we suggest interpreting the absolute  
540 estimated number of clusters with caution and focus instead on estimates of species similarity.

541 We only explored two factors likely to affect the performance of the DP-COM, though many  
542 other factors may be influential. Particularly, we imagine that the total community size and the  
543 cluster-to-size ratio (i.e., whether communities consist of many small or few large clusters) may  
544 affect the estimation of community structure: we would expect that more clusters should improve  
545 estimation of parameters governing  $G_0$ , and more species per cluster should improve estimates of  
546 cluster-level parameters. Additional simulations with communities of 60 species distributed  
547 across 5 or 10 clusters (i.e., representing a scenario with a higher species-to-cluster ratio, and one  
548 with the same ratio as in our main simulation but with more data) somewhat support these  
549 expectations, with community  $\beta$  coefficients having slightly lower CVs in the scenario with more  
550 clusters, and species-level coefficients (which are derived from cluster-level estimates) having  
551 slightly lower CVs when there were more species per cluster (Appendix2: Table S12). Having a  
552 larger community reduced bias in community and species coefficients, regardless of the  
553 community structure. Neither scenario, however, suggested that using the DP over a normal  
554 COM led to greater improvements in either precision or bias of estimates when communities  
555 were larger (Appendix S2: Figure S4 and S5). Factors of study design, such as spatial and  
556 temporal repeats, as well as the amount of data available for each species have been shown to  
557 affect performance of occupancy models (MacKenzie and Royle 2005, Pacifici et al. 2014) and  
558 may affect the DP-COM as well. Due to the computational cost of the DP-COM, however, we  
559 were unable to explore these additional dimensions in more depth.

560

#### 561 *Accuracy of parameter estimates*

562 Whereas estimates of typical parameters of interest (number of sites occupied, coefficients of the  
563 probit-linear predictor of occupancy) were, on average, unbiased under both modeling  
564 approaches, bias and CV were high in some individual species-iteration combinations,  
565 particularly in estimates of species-specific coefficients (Appendix2: Figure S1). Even though  
566 the DP-COM adequately reflected the clustered nature of the simulated communities, it did not

567 consistently improve bias and precision of estimates. We performed some exploratory post-hoc  
568 analyses (results not shown) that showed that specific species-iteration combinations had  
569 consistently high CV and bias across both modeling approaches, suggesting that some  
570 characteristic of the data was responsible for poor estimates. We investigated whether instances  
571 of large CVs and bias were associated with sparse data, but patterns were inconclusive. We do  
572 not have data on bias and precision of parameter estimates under the two modeling approaches  
573 fit to data generated under a normal COM (i.e., a non-clustered community), but we suspect that  
574 the incidences of high CVs and bias are related to the clustered structure of the community.

575

#### 576 *Sensitivity to prior*

577 It has been shown that the estimate of the concentration parameter  $\alpha$ , which determines the  
578 number of clusters, is sensitive to its prior (Dorazio 2009). Following the principle of preferring  
579 parsimonious models, we adopted the prior by Johnson and Sinclair (2017), which allows for a  
580 wide range of values of  $K$  but favors smaller values and did not appear to affect estimates of  $K$  in  
581 their simulation. Nonetheless, under most scenarios, we observed positive bias in  $\hat{K}$ . Because of  
582 the time-intensive nature of the DP model, thorough testing of sensitivity of  $\hat{\alpha}$  and, by extension,  
583  $\hat{K}$ , to priors was beyond the scope of this study. For a small subset of simulations, however, we  
584 explored whether a negative-exponential prior, which puts even more weight on fewer clusters,  
585 would reduce the positive bias in  $\hat{K}$ , but found no improvements. Dorazio (2009) suggested a  
586 *Gamma*( $a, b$ ) prior where  $a$  and  $b$  are chosen depending on  $n$  (the number of species in the data  
587 set), so that the prior on  $\alpha$  reflects a *discrete-Uniform*( $0, n$ ) prior on  $K$ . On the other hand, West  
588 et al. (1994) suggest a static *Gamma*(3.5, 0.5) prior allowing for a wide range of possible values  
589 for  $K$ , with low probability at 0 and  $n$ . Studies employing the DP-COM should evaluate the  
590 influence of the choice of prior for  $\alpha$  on main quantities of interest.

591

#### 592 *Structure and habitat associations in the MFNP bird community*

593 We found considerable structure within the MFNP bird community, with 166 species clustering  
594 into 27 groups regarding their associations with habitat type and distance to oil well. Some  
595 species pairs showed similarity scores  $>0.90$ , being in the same cluster virtually all the time.

596 Across the community, we found that occupancy of more species was significantly related to

597 habitat type (open versus woodland) than influenced by distance from oil wells. It is conceivable  
598 that the effect of oil drilling operations on bird occurrence may be temporally limited to when  
599 wells are active (Fuda et al. 2018). Our analysis of occupancy across multiple months may mask  
600 any such temporal effects and only show effects of this factor in cases of strong species  
601 responses. Coefficients of the two predictors of occupancy were positively correlated, indicating  
602 that an increasing preference of woodland habitat corresponded with greater avoidance of oil  
603 wells. No species that had strong negative associations with distance to oil wells had strong  
604 positive associations with woodland habitat; similarly, none of the species strongly preferring  
605 woodland habitat had significant negative associations with distance to oil well. This suggests  
606 that birds with a preference for more closed habitat tend to be more sensitive to habitat  
607 disturbance, a conclusion reached for birds and other taxa in a recent meta-analysis (Keinath et  
608 al. 2017).

609 The mechanisms determining how closely related species, possibly with similar morphologies  
610 and diets, can coexist has been an ongoing debate in ecology for decades (Hutchinson 1959,  
611 Gotelli 2000, Graham et al. 2004). Many studies have found that sister taxa commonly occupy  
612 different ecological niches and that co-occurring species are generally more distantly related  
613 (e.g., Silva et al. 2014), while others argue that phylogeny begets morphological similarity and,  
614 thus, higher likelihood of niche overlap (Gonçalves-Souza et al. 2014). These patterns are of  
615 interest in conservation biology as well, with recent studies exploring the effect of phylogeny  
616 and other traits on species susceptibility to disturbance (Nowakowski et al. 2017). When  
617 comparing within-family clustering rates – a measure of how similarly species in the present  
618 study use space – to average similarity of the community, only the sunbirds stood out as more  
619 similar than average. The tropical sunbirds are largely nectivorous but also consume fruit and  
620 insects and, thus, generally considered to be forest/woodland species where their specialized  
621 food sources are likely more plentiful (Cheke et al. 2019). The five species represented in our  
622 sample demonstrated significant niche conservatism, with consistently strong positive  
623 associations with woodland habitat and distance to oil (only one species, the Marico sunbird  
624 *Cinnyris mariquensis*, had 95BCI overlapping 0 for the latter). Thus, based on our findings the  
625 sunbirds represent an example of where “phylogeny begets niche overlap”, and possibly of low  
626 response diversity with respect to anthropogenic influence (oil wells). Of course, habitat  
627 partitioning among these species may very well happen on scales other than the one measured in

628 this study. Even though species-specific coefficient estimates were generally very similar  
629 between the DP and the normal COM, sunbird coefficients were more similar to each other under  
630 the DP-COM than the normal COM, suggesting that the approach is better able to represent  
631 similarities among species (Appendix S1: Figure S2).

632 In contrast, members of the Cisticolidae and Meropidae (bee-eaters), with 12 and 5  
633 representative species respectively, demonstrated clustering probabilities that were lower than  
634 average, and only slightly higher than expected under random clustering. Thus, they exemplify  
635 the “niche differentiation among closely related taxa” argument. Bee-eaters are considered  
636 habitat generalists, occupying both forests, edge and open habitat; their aerial behavior is  
637 generally independent of vegetation type (Fry 2019). Cisticolidae is a broad taxon that includes  
638 Old World warblers and other allies that occupy a range of habitats including forest, open  
639 woodland, scrub and grassland (Ryan 2019). This example illustrates the potential usefulness of  
640 the DP-COM for addressing ecological questions of species coexistence, estimating similarity of  
641 species while fully accounting for imperfect species detection and uncertainty in coefficient  
642 estimates.

## 644 **Conclusion**

645 Dirichlet process distributions provide a flexible tool to model latent structure in wildlife  
646 communities and populations. Our DP-COM is a straight-forward extension of popular  
647 community occupancy models (e.g., Zipkin et al. 2009, Ruiz-Gutiérrez et al. 2010, Sollmann et  
648 al. 2017) and can be implemented in JAGS, a software that has become increasingly popular  
649 among ecologists and wildlife researchers (Kéry 2010, Kéry and Schaub 2012). Implementing  
650 the DP-COM in JAGS was computationally much more expensive than the normal COM – for  
651 the bird data set, the difference was on the scale of hours (normal COM) versus >5 days (DP-  
652 COM). Based on our simulation study, run time increases non-linearly with the addition of  
653 species to the data set (from 15 minutes for a 30-species community to 1.5 hours for a 60-species  
654 community). Mixing of chains was slow, suggesting that longer chains, and thus more  
655 computation time, would be beneficial. Whereas implementation of these models can be  
656 accelerated by using a custom MCMC algorithm, and likely also by using the reversible jump  
657 MCMC capacities of NIMBLE (de Valpine et al. 2017), they still remain computationally  
658 involved (Johnson and Sinclair 2017). This complicates thorough evaluation of model

659 performance under different conditions via simulations and makes models less accessible to  
660 practitioners. Even though the DP model has fewer parameters than the normal COM, the  
661 improvement in precision of estimates was marginal or non-existent, and in spite of the distinctly  
662 clustered simulated community, the normal COM returned estimates of ecological parameters  
663 that were, for the most part, as accurate and precise as those of the DP-COM. For studies where  
664 estimates of occupancy and associated covariate coefficients are the main focus, our results thus  
665 suggest the much faster and better-mixing normal COM provides reliable results. We did not test  
666 whether joint prediction of community occupancy at new sites benefits from the DP-COM, and  
667 this warrants further investigation for studies where prediction is a key objective. The DP-COM  
668 may be the better approach in situations where researchers would otherwise resort to a priori  
669 grouping of species. Especially for sparse data species, inference on the species level is affected  
670 by how groups are defined (Pacifci et al. 2014); under a DP-COM, estimates of parameters for  
671 such species will represent the average over possible group associations and thus avoid  
672 subjectivity in choosing a certain grouping. The main advantage of the DP-COM is the  
673 information about community structure and species similarity with respect to occupancy  
674 predictors that the normal model cannot provide directly. We present an example of how this  
675 information can be used to address questions of ecological relevance with the Uganda bird  
676 example.

677 Our model development only considers the simplest case of a DP model, in which no  
678 information on species cluster membership is available. The DP model can be extended to  
679 include covariates that can inform the probability of cluster membership (Johnson et al. 2013). In  
680 the context of community occupancy models, inclusion of potential clustering covariates enables  
681 testing whether species attributes such as taxonomy or functional traits explain community  
682 structure. As such, in spite of its drawbacks, the semi-parametric DP-COM holds potential as a  
683 flexible modelling approach in situations where community structure and species similarities are  
684 of primary interest.

685

## 686 **Acknowledgments**

687 Data collection was made by the Wildlife Conservation Society and financed by the  
688 USAID/WILD Program. We are particularly grateful to the ornithologists who collected the data,  
689 notably Hamlet Mugabe, Dennis Tumuhame, and Taban Bruhan. We are also grateful for

690 permission to undertake the research by the Uganda Wildlife Authority. We further thank Paul  
691 Conn (NOAA) for suggesting the conceptual DP-COM. Data and code storage on Dryad is  
692 sponsored by UC Davis. Any use of trade, firm, or product names is for descriptive purposes  
693 only and does not imply endorsement by the U.S. Government. Author contributions: PM, SA  
694 and SP oversaw data collection in the field, led the field teams or made regular visits to check on  
695 the data collection and to supply field teams with food and organize other logistical  
696 arrangements. AJP and SP conceived of the project, secured funding, and designed the study.  
697 RS, MJE and WAL conceived of the analytical approach for the paper, RS, DSJ and WAL led  
698 the model and simulation study development, RS performed the data analysis, RS and DSJ led  
699 simulation results interpretation, while RS and MJE led case study interpretation and writing of  
700 the manuscript, with input from all coauthors.

701

## 702 **Supporting Information**

703 Additional supporting information may be found online at: [link to be added in production]

704

## 705 **Data Availability**

706 Data and code can be accessed on the Dryad Digital Repository:

707 <https://doi.org/10.25338/B8GG8P>

708

## 709 **References**

710 Azizyan, M., A. Singh, and L. Wasserman. 2013. Minimax theory for high-dimensional gaussian  
711 mixtures with sparse mean separation. Pages 2139–2147 Proceedings of the 26th  
712 International Conference on Neural Information Processing Systems - Volume 2. Lake  
713 Tahoe, Nevada.

714 Bellman, R. 1957. Dynamic programming. Princeton University Press. Princeton, NJ.

715 Cheke, R., C. Mann, and A. Bonan. 2019. Sunbirds (Nectariniidae). Page *in* J. del Hoyo, A.  
716 Elliott, J. Sargatal, D. A. Christie, and E. de Juana, editors. Handbook of the Birds of the  
717 World Alive. Lynx Edicions, Barcelona, Spain.

718 Dorazio, R. M. 2009. On selecting a prior for the precision parameter of Dirichlet process  
719 mixture models. Journal of Statistical Planning and Inference 139:3384–3390.

720 Dorazio, R. M., M. Kéry, J. A. Royle, and M. Plattner. 2010. Models for inference in dynamic



- 721 metacommunity systems. *Ecology* 91:2466–2475.
- 722 Dorazio, R. M., B. Mukherjee, L. Zhang, M. Ghosh, H. L. Jelks, and F. Jordan. 2008. Modeling  
723 unobserved sources of heterogeneity in animal abundance using a Dirichlet process prior.  
724 *Biometrics* 64:635–644.
- 725 Dorazio, R. M., and J. A. Royle. 2005. Estimating size and composition of biological  
726 communities by modeling the occurrence of species. *Journal of the American Statistical*  
727 *Association* 100:389–398.
- 728 Dorazio, R. M., J. A. Royle, B. Söderström, and A. Glimskär. 2006. Estimating species richness  
729 and accumulation by modeling species occurrence and detectability. *Ecology* 87:842–  
730 854.
- 731 Dunstan, P. K., S. D. Foster, and R. Darnell. 2011. Model based grouping of species across  
732 environmental gradients. *Ecological Modelling* 222:955–963.
- 733 Dunstan, P. K., S. D. Foster, F. K. Hui, and D. I. Warton. 2013. Finite mixture of regression  
734 modeling for high-dimensional count and biomass data in ecology. *Journal of*  
735 *agricultural, biological, and environmental statistics* 18:357–375.
- 736 Escobar, M. D., and M. West. 1995. Bayesian density estimation and inference using mixtures.  
737 *Journal of the american statistical association* 90:577–588.
- 738 Fry, H. 2019. Bee-eaters (Meropidae). Page *in* J. del Hoyo, A. Elliott, J. Sargatal, D.A. Christie,  
739 and E. de Juana, editors. *Handbook of the Birds of the World Alive*. Lynx Edicions,  
740 Barcelona, Spain.
- 741 Fuda, R. K., S. J. Ryan, J. B. Cohen, J. Hartter, and J. L. Frair. 2018. Assessing the impacts of oil  
742 exploration and restoration on mammals in Murchison Falls Conservation Area, Uganda.  
743 *African Journal of Ecology* 56:804–817.
- 744 Gelman, A., and J. Hill. 2006. *Data Analysis Using Regression and Multilevel/Hierarchical*  
745 *Models*. First edition. Cambridge University Press, New York, USA.
- 746 Gonçalves-Souza, T., J. A. F. Diniz-Filho, and G. Q. Romero. 2014. Disentangling the  
747 phylogenetic and ecological components of spider phenotypic variation. *PloS one*  
748 9:e89314.
- 749 Gotelli, N. J. 2000. Null model analysis of species co-occurrence patterns. *Ecology* 81:2606–  
750 2621.
- 751 Graham, C. H., S. R. Ron, J. C. Santos, C. J. Schneider, and C. Moritz. 2004. Integrating

752 phylogenetics and environmental niche models to explore speciation mechanisms in  
753 dendrobatid frogs. *Evolution* 58:1781–1793.

754 Houle, M. E., H.-P. Kriegel, P. Kröger, E. Schubert, and A. Zimek. 2010. Can shared-neighbor  
755 distances defeat the curse of dimensionality? Pages 482–500 *International Conference on*  
756 *Scientific and Statistical Database Management*. Springer.

757 Hutchinson, G. E. 1959. Homage to Santa Rosalia or why are there so many kinds of animals?  
758 *The American Naturalist* 93:145–159.

759 Johnson, D. S., R. R. Ream, R. G. Towell, M. T. Williams, and J. D. L. Guerrero. 2013.  
760 Bayesian clustering of animal abundance trends for inference and dimension reduction.  
761 *Journal of Agricultural, Biological, and Environmental Statistics* 18:299–313.

762 Johnson, D. S., and E. H. Sinclair. 2017. Modeling joint abundance of multiple species using  
763 Dirichlet process mixtures. *Environmetrics* 28:e2440.

764 Keinath, D. A., D. F. Doak, K. E. Hodges, L. R. Prugh, W. Fagan, C. H. Sekercioglu, S. H.  
765 Buchart, and M. Kauffman. 2017. A global analysis of traits predicting species sensitivity  
766 to habitat fragmentation. *Global Ecology and Biogeography* 26:115–127.

767 Kellner, K. 2019. jagsUI: A Wrapper Around “rjags” to Streamline “JAGS” Analyses.  
768 <https://CRAN.R-project.org/package=jagsUI>

769 Kéry, M. 2010. *Introduction to WinBUGS for ecologists: A Bayesian approach to regression,*  
770 *ANOVA and related analyses*. Academic Press, Burlington, MA.

771 Kéry, M., and J. A. Royle. 2008. Hierarchical Bayes estimation of species richness and  
772 occupancy in spatially replicated surveys. *Journal of Applied Ecology* 45:589–598.

773 Kéry, M., and M. Schaub. 2012. *Bayesian population analysis using WinBUGS: a hierarchical*  
774 *perspective*. Academic Press, Burlington, MA.

775 MacKenzie, D. I., J. D. Nichols, G. B. Lachman, S. Droege, J. Andrew Royle, and C. A.  
776 Langtimm. 2002. Estimating site occupancy rates when detection probabilities are less  
777 than one. *Ecology* 83:2248–2255.

778 MacKenzie, D. I., J. D. Nichols, J. A. Royle, K. H. Pollock, L. L. Bailey, and J. E. Hines. 2017.  
779 *Occupancy estimation and modeling: inferring patterns and dynamics of species*  
780 *occurrence*. 2nd edition. Elsevier, Amsterdam.

781 MacKenzie, D. I., and J. A. Royle. 2005. Designing occupancy studies: general advice and  
782 allocating survey effort. *Journal of Applied Ecology* 42:1105–1114.

- 783 Mori, A. S., T. Furukawa, and T. Sasaki. 2013. Response diversity determines the resilience of  
784 ecosystems to environmental change. *Biological Reviews* 88:349–364.
- 785 Nowakowski, A. J., J. I. Watling, S. M. Whitfield, B. D. Todd, D. J. Kurz, and M. A. Donnelly.  
786 2017. Tropical amphibians in shifting thermal landscapes under land-use and climate  
787 change. *Conservation Biology* 31:96–105.
- 788 Ohlssen, D. I., L. D. Sharples, and D. J. Spiegelhalter. 2007. Flexible random-effects models  
789 using Bayesian semi-parametric models: applications to institutional comparisons.  
790 *Statistics in Medicine* 26:2088–2112.
- 791 Pacifici, K., E. F. Zipkin, J. A. Collazo, J. I. Irizarry, and A. DeWan. 2014. Guidelines for a  
792 priori grouping of species in hierarchical community models. *Ecology and Evolution*  
793 4:877–888.
- 794 Plummer, M. 2003. JAGS: A program for analysis of Bayesian graphical models using Gibbs  
795 sampling. Pages 20–22 *Proceedings of the 3rd International Workshop on Distributed*  
796 *Statistical Computing (DSC 2003)*.
- 797 Plumtre, A. J., T. R. Davenport, M. Behangana, R. Kityo, G. Eilu, P. Ssegawa, C. Ewango, D.  
798 Meirte, C. Kahindo, M. Herremans, and others. 2007. The biodiversity of the Albertine  
799 Rift. *Biological conservation* 134:178–194.
- 800 R Core Team. 2018. R: A language and environment for statistical computing. R Foundation for  
801 Statistical Computing, Vienna, Austria. <http://www.R-project.org/>
- 802 Reich, B. J., and H. D. Bondell. 2011. A spatial Dirichlet process mixture model for clustering  
803 population genetics data. *Biometrics* 67:381–390.
- 804 Royle, J. A., and R. M. Dorazio. 2008. *Hierarchical modeling and inference in ecology*.  
805 Academic Press, London, UK.
- 806 Ruiz-Gutiérrez, V., E. F. Zipkin, and A. A. Dhondt. 2010. Occupancy dynamics in a tropical bird  
807 community: unexpectedly high forest use by birds classified as non-forest species.  
808 *Journal of Applied Ecology* 47:621–630.
- 809 Ryan, P. 2019. Cisticolas and allies (Cisticolidae). Page *in* J. del Hoyo, A. Elliott, J. Sargatal,  
810 D.A. Christie, and E. de Juana, editors. *Handbook of the Birds of the World Alive*. Lynx  
811 Edicions, Barcelona, Spain.
- 812 Sauer, J. R., and W. A. Link. 2002. Hierarchical modeling of population stability and species  
813 group attributes from survey data. *Ecology* 83:1743–1751.

- 814 Sethuraman, J. 1994. A constructive definition of Dirichlet priors. *Statistica Sinica* 4:639–650.
- 815 Silva, D. P., B. Vilela, P. De Marco Jr, and A. Nemesio. 2014. Using ecological niche models  
816 and niche analyses to understand speciation patterns: the case of sister neotropical orchid  
817 bees. *PLoS One* 9:e113246.
- 818 Sollmann, R., B. Gardner, K. A. Williams, A. T. Gilbert, and R. R. Veit. 2016. A hierarchical  
819 distance sampling model to estimate abundance and covariate associations of species and  
820 communities. *Methods in Ecology and Evolution* 7:529–537.
- 821 Sollmann, R., A. Mohamed, J. Niedballa, J. Bender, L. Ambu, P. Lagan, S. Mannan, R. C. Ong,  
822 A. Langner, B. Gardner, and Wilting, Andreas. 2017. Quantifying mammal biodiversity  
823 co-benefits in certified tropical forests. *Diversity and Distributions* 23:317–328.
- 824 Tiao, G. C., and A. Zellner. 1964. Bayes's theorem and the use of prior knowledge in regression  
825 analysis. *Biometrika* 51:219–230.
- 826 de Valpine, P., D. Turek, C. J. Paciorek, C. Anderson-Bergman, D. T. Lang, and R. Bodik. 2017.  
827 Programming with models: writing statistical algorithms for general model structures  
828 with NIMBLE. *Journal of Computational and Graphical Statistics* 26:403–413.
- 829 Yamaura, Y., M. Kery, and J. A. Royle. 2016. Study of biological communities subject to  
830 imperfect detection: bias and precision of community N-mixture abundance models in  
831 small-sample situations. *Ecological Research* 31:289–305.
- 832 Zipkin, E. F., A. DeWan, and A. J. Royle. 2009. Impacts of forest fragmentation on species  
833 richness: a hierarchical approach to community modelling. *Journal of Applied Ecology*  
834 46:815–822.

### 835 Figure legends

836 Figure 1: Estimated number of clusters  $K$  (a) and pairwise clustering rates (b) from a Dirichlet  
837 Process (DP) community occupancy model used to analyze data simulated under different levels  
838 of cluster distinctiveness ( $\omega$ ) and different number of coefficients in the probit-linear predictor of  
839 occupancy (corresponding to dimensions of the multivariate DP),  $m$ . For a), violin plots depict  
840 posterior modes of  $K$  across iterations; red line shows the data generating value. For b), violins  
841 show the average number of MCMC iterations during which two species were estimated to be in  
842 the same cluster when in the simulated data they were in the same cluster (blue) and when they  
843 were in different clusters (orange). In both panels, dots represent the median across iterations.

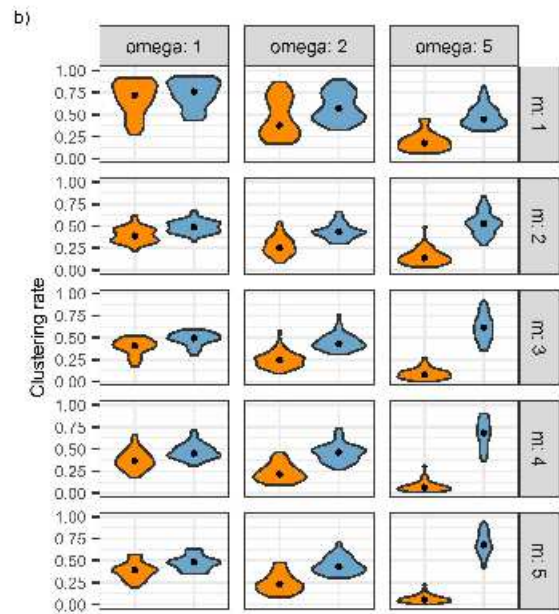
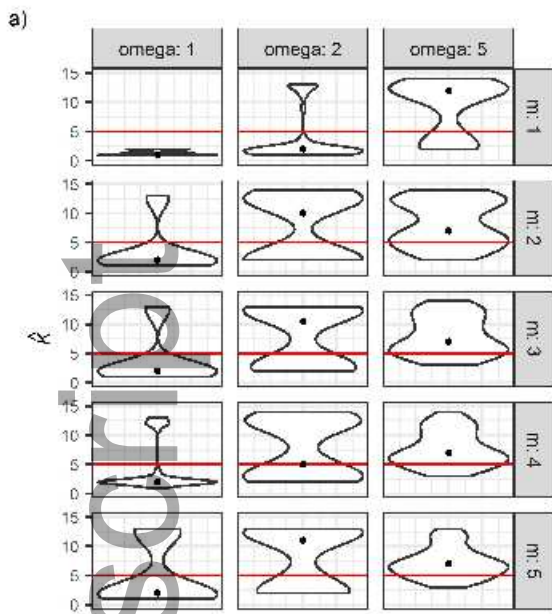
844

845 Figure 2: Bias (a) and coefficient of variation, CV, (b) of estimated number of sites occupied by  
846 species from community occupancy models using either a Dirichlet Process (DP) or a normal  
847 species level random effect. Models were used to analyze data simulated under different levels of  
848 cluster distinctiveness ( $\omega$ ) and different number of coefficients in the probit-linear predictor of  
849 occupancy (corresponding to dimensions of the multivariate DP in the DP COM),  $m$ . Violins  
850 represent estimates across species and iterations in a given scenario. Plot y-axes capped at -1/1  
851 (a) and 0/1 (b) for aesthetic reasons.

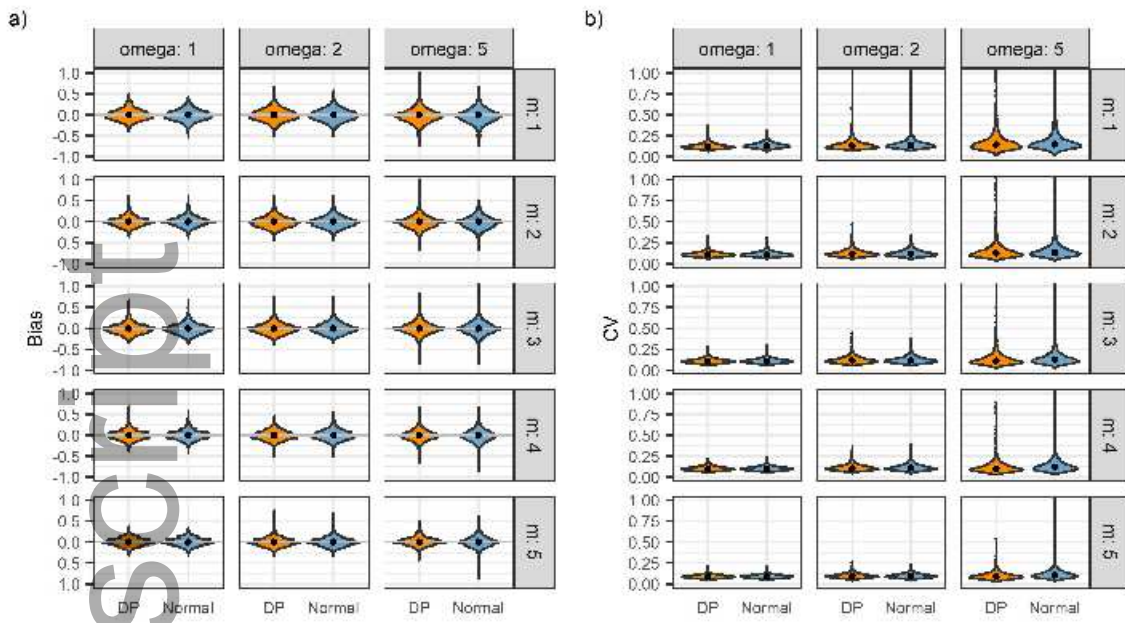
852  
853 Figure 3: Probability of joint cluster membership for 166 birds in Murchison Falls National Park,  
854 Uganda, estimated from a Dirichlet Process community occupancy model, based on coefficients  
855 in the probit-linear predictor of occupancy probability, including the effect of habitat type (open  
856 versus woodland) and distance from oil well. Both axes represent species identity and color  
857 gradient expresses the probability of joint cluster membership.

858  
859 Figure 4: Pairwise probabilities of joint cluster membership (similarity), estimated from a  
860 Dirichlet Process community occupancy model, for 10 bird families with at least 5 species  
861 observed during a survey in Murchison Falls National Park, Uganda (number of species given  
862 above error bars). Dots: average probabilities of joint cluster membership across species; error  
863 bars: 5<sup>th</sup> and 95<sup>th</sup> percentiles; black line/grey rectangle: mean and 5<sup>th</sup> and 95<sup>th</sup> percentile of  
864 probabilities of joint cluster membership for entire community; red line: maximum clustering  
865 probability observed when simulating random clustering.

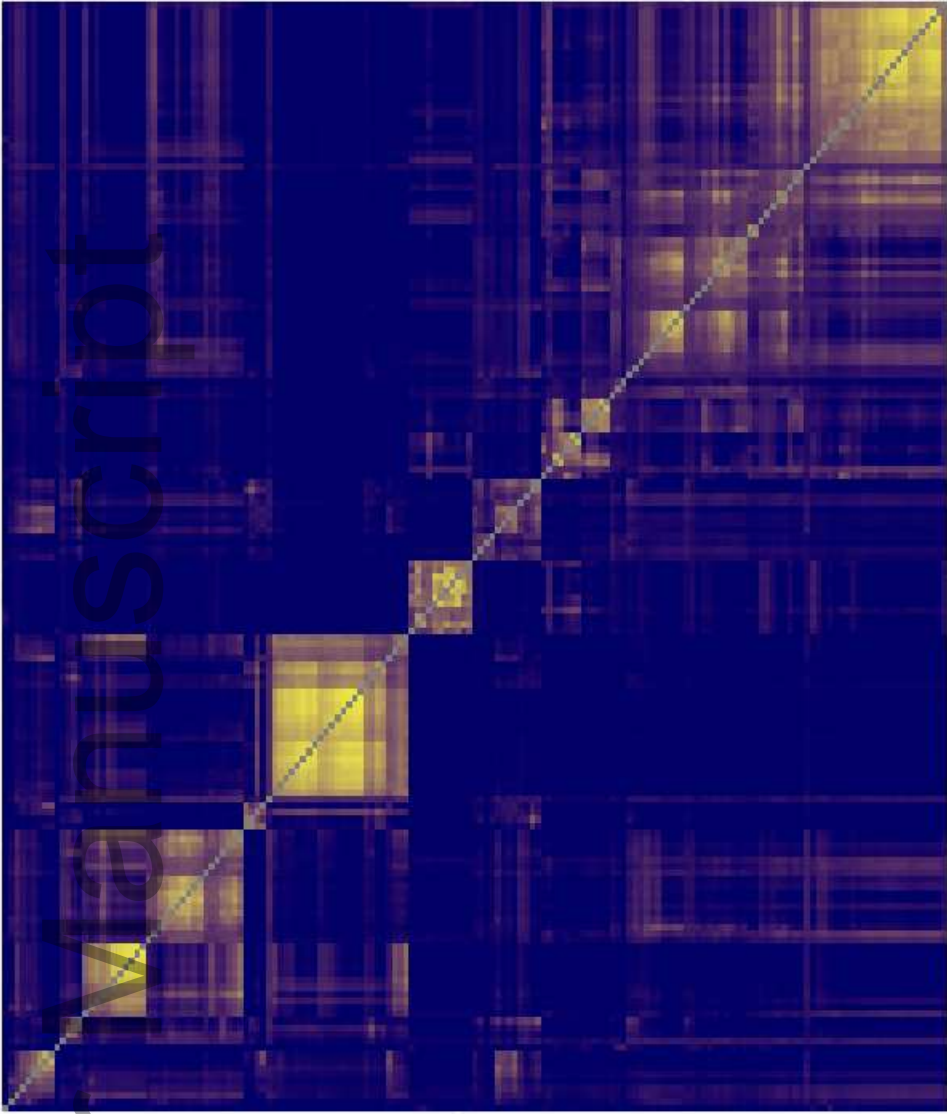
866  
867 Figure 5: Beta coefficients for effect of woodland habitat,  $\beta(\text{habitat})$ , and distance to oil well,  
868  $\beta(\text{Oil})$ , on occupancy probability for 166 birds surveyed in Murchison Falls National Park,  
869 Uganda, estimated with a Dirichlet Process community occupancy model. Effects considered  
870 strong when 95% Bayesian Credible Intervals did not overlap 0.



eap\_2249\_f1.jpg

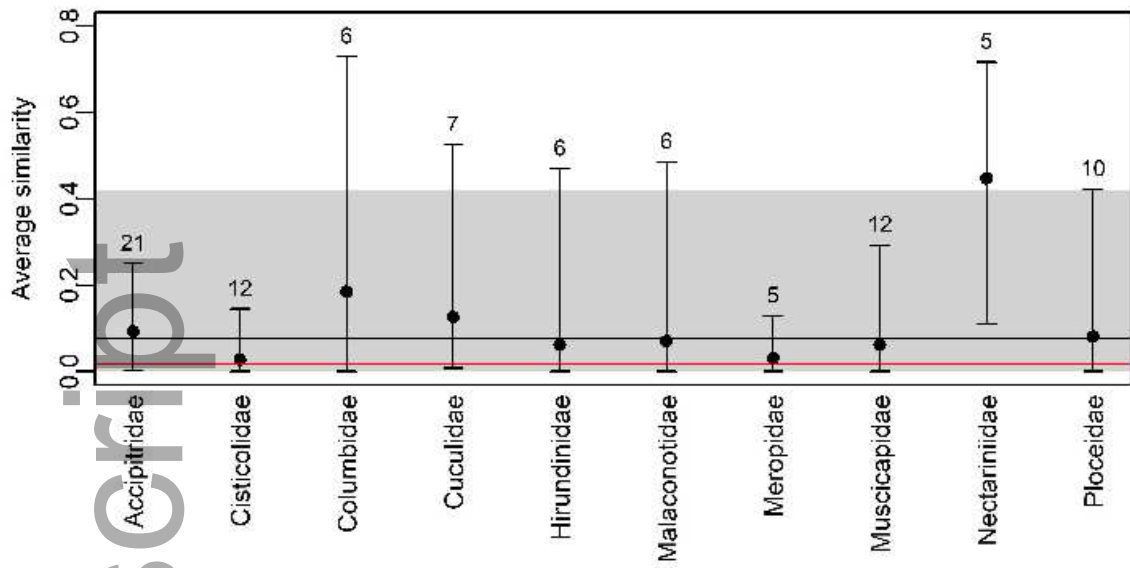


eap\_2249\_f2.jpg

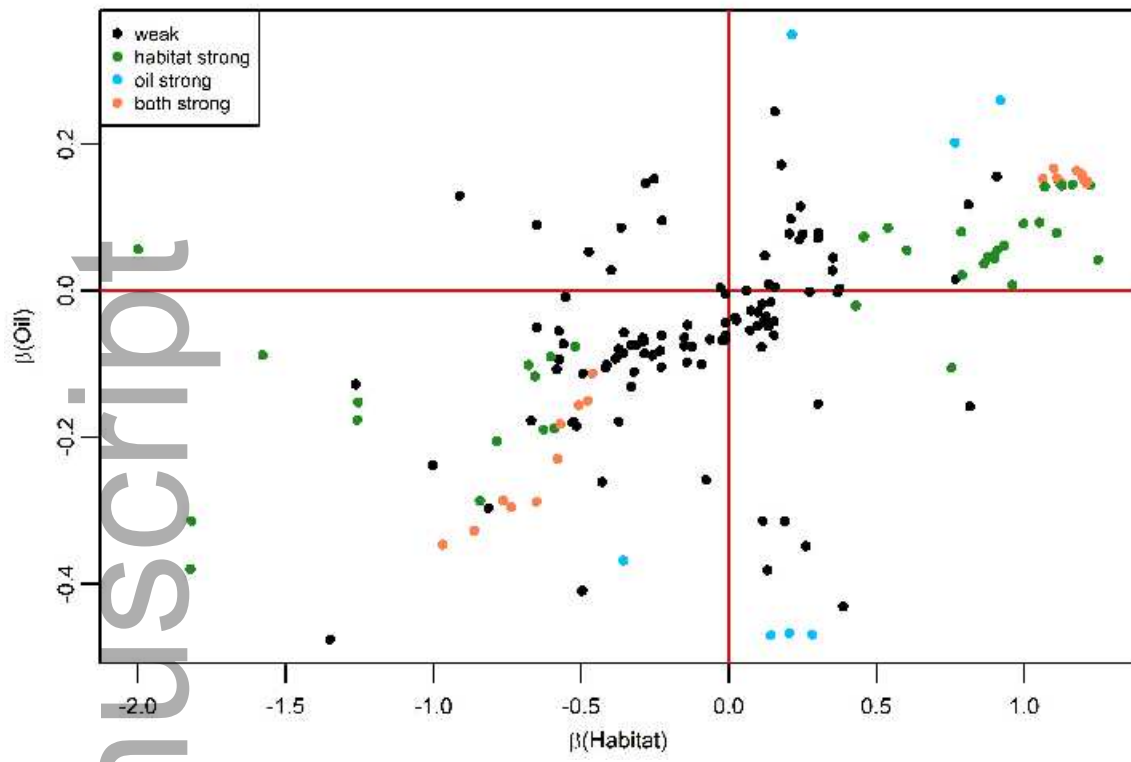


eap\_2249\_f3.jpg





eap\_2249\_f4.jpg



eap\_2249\_f5.jpg