

1

Université de Montréal

2

## **Spatially continuous identification of beta diversity hotspots using species distribution models**

3

4

By

5

**Gabriel Dansereau**

6

20147609

7

Département de sciences biologiques

8

Faculté des arts et des sciences

9

Advisory Committee Meeting

10

November 29, 2019

## **11    Contents**

12	Abstract . . . . .	3
13	Introduction . . . . .	4
14	Methods . . . . .	6
15	1. Data Collection . . . . .	6
16	2. Data Manipulation . . . . .	7
17	3. SDM – The BIOCLIM Method . . . . .	7
18	4. LCBD Calculation . . . . .	8
19	5. Prediction Validity . . . . .	9
20	6. Alternative methods . . . . .	10
21	7. Climate Change Scenarios and Temporal Beta Diversity . . . . .	11
22	Preliminary Results . . . . .	11
23	References . . . . .	14
24	Appendix . . . . .	17

25      **List of Tables**

26      1      Structure of the Warblers data in the eBird checklists for the countries used in	17
27             the analyses . . . . .	
28      2      Description of the WorldClim 2 climate variables used in the analyses . . . . .	18

29      **List of Figures**

30      1      Distribution of species richness in North America, defined as the number	
31      of Warblers species per site. The raw occurrence observations from eBird	
32      (fig. 1a) and the SDM predictions from the BIOCLIM model (fig. 1b) were both	
33      transformed into presence-absence data per species before calculating richness.	19
34      2      Distribution of the LCBD values in North America, calculated from the vari-	
35      ance of the community matrix Y and scaled to the maximum value observed.	
36      The Hellinger transformation was applied on the raw occurrence data (fig. 2a)	
37      before calculating the LCBD indices. SDM predictions (fig. 2b) were con-	
38      verted into presence-absence data, but no transformation was applied before	
39      calculating the LCBD indices. . . . .	20
40      3      Relationship between the species richness and the LCBD value of the each	
41      site for raw occurrence data (blue) and SDM predictions (orange). Species	
42      richness was calculated as the number of species in a site ( $\alpha$ ), divided by the	
43      total number of species ( $\gamma$ ). LCBD values were scaled to the maximum value	
44      observed. Hellinger transformation was applied on the raw occurrence data	
45      before calculating LCBD indices. . . . .	21

<sup>46</sup> **Abstract**

<sup>47</sup> Beta diversity is an essential measure to describe the organization of biodiversity in space.  
<sup>48</sup> The calculation of local contributions to beta diversity (LCBD), specifically, allows for the  
<sup>49</sup> identification of sites with exceptional diversity within a region of interest, which is useful  
<sup>50</sup> for both community ecology and conservation purposes. However, beta diversity implies a  
<sup>51</sup> comparison among the sites of a given region, thus, its use is restricted to sites with known  
<sup>52</sup> species composition, and to discontinuous spatial scales. We therefore propose a method  
<sup>53</sup> to calculate LCBD indices on continuous scales for a whole region of interest, including  
<sup>54</sup> unsampled sites. First, species distributions can be predicted on continuous scales using  
<sup>55</sup> species distribution models (SDM). These models, such as the BIOCLIM method, use the  
<sup>56</sup> environmental conditions at sampled sites to predict the presence or absence of each species at  
<sup>57</sup> unsampled locations. Second, LCBD statistics can then be computed on the SDM predictions.  
<sup>58</sup> We show that it is therefore possible to identify beta diversity hotspots on spatially continuous  
<sup>59</sup> and extended scales. Our results confirm that LCBD values are related to species richness,  
<sup>60</sup> and that species-poor sites contribute most to beta diversity.

61 **Introduction**

62 Beta diversity, defined as the variation in species composition among sites in a geographic  
63 region of interest (Legendre, Borcard, and Peres-Neto 2005), is an essential measure to  
64 describe the organization of biodiversity in space. Total beta diversity within a community  
65 can be partitioned into local contributions to beta diversity (LCBD) (Legendre and De Cáceres  
66 2013), which allows for the identification of sites with exceptional species composition,  
67 hence exceptional biodiversity. Such a method is useful for both community ecology and  
68 conservation biology, as it highlights sites that are most important for their research or  
69 conservation values. However, LCBD calculation methods require complete information  
70 on community composition, such as a community composition matrix  $Y$ , thus they are  
71 inappropriate for partially sampled or unsampled sites. To our knowledge, these methods  
72 have mostly been applied on community data from sampled sites, hence on discontinuous  
73 spatial scales, e.g. at intervals along a river stream (Legendre and De Cáceres 2013). This  
74 raises the following questions: 1) could LCBD indices be extended to continuous spatial  
75 scales, and 2) could this provide novel ecological insights in poorly sampled regions? We  
76 aim to answer these questions by combining the LCBD calculation methods with predictive  
77 biogeography approaches, and suggest that this would allow for the identification of hotspots  
78 with high conservation value in poorly sampled regions.

79 Species distribution models (SDMs) already allow to make predictions on continuous spatial  
80 scales, and these predictions could therefore be used to calculate LCBD indices. SDMs, also  
81 known as bioclimatic envelope models (Araújo and Peterson 2012), aim to predict species  
82 presence or absence based on previous observations of occurrence, and the environmental  
83 conditions at which these were made (Poisot et al. 2019). Examples of uses include climate  
84 change impact and invasion risk assessment, reserve selection and design, and discovery  
85 of new populations (Araújo and Peterson 2012). This way, they generate novel ecological  
86 insights for unsampled or lesser-known locations (Poisot et al. 2019), an approach yet to  
87 be applied to the LCBD framework. We believe that a predictive approach such as this one  
88 would bring a new perspective to biodiversity study and community ecology. By using SDMs,  
89 we would be able to expand community information already available, and thus work on a  
90 much larger community matrices than in typical LCBD studies, which might highlight new

91 diversity hotspots.

92 Climate and biodiversity data on extended spatial scales are increasingly available online.  
93 For instance, the Worldclim 2.0 database (Fick and Hijmans 2017) provides interpolated  
94 climate data for global land areas at very high spatial resolution, and the eBird platform  
95 (Sullivan et al. 2009) provides a rapidly growing, citizen-contributed database of worldwide  
96 bird observations. Both of these are commonly used in SDMs, and offer relevant information  
97 on extended spatial scales. Therefore, we believe that these datasets could be used to predict  
98 community composition and calculate LCBD indices on continuous spatial scales, and that  
99 the result would be representative of the true community structure.

100 The predictive approach we suggest would be especially useful in poorly sampled regions, or  
101 in regions with only sparse sampling. While it does not replace a full sampling within the  
102 community, predictions and exploratory analyses do provide relevant ecological insights that  
103 could be used in different ways. For instance, our method could help identify unsampled  
104 sites with potential conservation value which should be targeted as soon as possible in future  
105 studies. We believe that the method could also be combined with IPCC climate change  
106 scenarios, which provide projections for climate variables, in a way that would allow us to  
107 model beta diversity changes with climate change and to identify the sites where the changes  
108 in the community will be most important. Once again, this would prove very relevant in  
109 an informative approach, suggesting sites to prioritize for future conservation and more  
110 structured research.

111 In this document, we cover in more details the methods that we suggest for this M.Sc. research  
112 project. The preparation part of the project, including data collection and manipulation, has  
113 already been done. A workflow for the analyses, including code implementation, has been  
114 defined as well. We also detail preliminary analyses and results intended as proof-of-concept  
115 for the approach, which, of course, needs to be refined. Finally, we discuss methods that we  
116 intend to use in future analyses, and whose feasibility is not as clearly stated.

117 **Methods**

118 **1. Data Collection**

119 We decided to focus our analyses on bird species and collected the data available on eBird for  
120 the Warblers family. The complete database contains nearly 600 million observations. We  
121 chose to focus specifically on the Warblers family, as it is a diverse group, popular among  
122 birders, with over 30 million observations. Global citizen-contributed databases often present  
123 additional challenges compared to conventional datasets due to their lack of structure, as well  
124 as spatial and taxonomic biases (Johnston et al. 2019). For instance, there was a clear bias in  
125 our data towards the United States, where there were many more observations and sampling  
126 events (tbl. 1). However, eBird offers two advantages over other large scale datasets (Johnston  
127 et al. 2019): 1) the data is structured as checklists and users can explicitly specify their  
128 observations as “complete checklists” when all detected species were reported, which allows  
129 to infer information on species absences, and 2) the dataset is semi-structured and checklists  
130 are associated with metadata describing sampling effort, such as duration of search, distance  
131 travelled and number of observers, which can be used as controls in the analyses. Hence,  
132 model performance can be improved by inferring absences and subsampling checklists, while  
133 spatial bias can be compensated by including effort covariates in the model (Johnston et al.  
134 2019). Therefore, we believe the dataset can be appropriately used to achieve our objective of  
135 expanding measures of exceptional biodiversity through space.

136 We collected the data available in the WorldClim 2 database (Fick and Hijmans 2017) for  
137 North America, to which we decided to restrict our analyses. The WorldClim data consists  
138 of spatially interpolated monthly climate data for global areas, available for resolutions  
139 from 10 arc-minutes to 30 arc-seconds (around 18 km<sup>2</sup> and 1 km<sup>2</sup> at the equator). Since the  
140 release of the first version of the database in 2005 (Hijmans et al. 2005), it became the most  
141 common source of climate data for SDM studies (Booth et al. 2014). The variables we used  
142 were different measures of temperature and precipitation (tbl. 2), which very high global  
143 cross-validation coefficients (> 0.99 and 0.86 respectively) (Fick and Hijmans 2017). We chose  
144 to use the coarser 10 arc-minutes resolution in our preliminary analyses, as we believed it  
145 was sufficient for proof of concept of our method. However, Hijmans et al. (2005) showed  
146 high within-grid cell variation in the 10 arc-minutes data, and therefore recommended the

<sup>147</sup> use of the finer resolution, which hid less of the variation known to the model. Given this, we  
<sup>148</sup> might reconsider the resolution to use in our final analyses.

<sup>149</sup> We chose to restrict our analyses to North America given the high amount of data available in  
<sup>150</sup> eBird. We believed it represented a suitable scale for our models, large enough to cover a lot  
<sup>151</sup> of variation in environmental variables and community structure, as well as phenomena  
<sup>152</sup> such as species migration. We also expected such extent of the spatial scale to cover for  
<sup>153</sup> imprecision in estimated species ranges.

## <sup>154</sup> 2. Data Manipulation

<sup>155</sup> WorldClim variables and eBird occurrence data are provided in different formats, so they  
<sup>156</sup> required some manipulations to be combined together. WorldClim variables are provided in  
<sup>157</sup> a 2-dimensional grid format, useful for large scale analyses and visualization, where each  
<sup>158</sup> cell or pixel has a size corresponding to the resolution of 10 arc-minutes. Each of the 19  
<sup>159</sup> variables forms a different grid. On the other hand, eBird records are occurrence-based, so  
<sup>160</sup> each entry in the dataset corresponds to an observation of a single species at a given time  
<sup>161</sup> and location. These entries can easily be matched to the 2D grid format of the WorldClim  
<sup>162</sup> variables through their spatial coordinates, which we found more useful for large scale  
<sup>163</sup> analyses and visualization. Hence, for each species, we matched all occurrences in eBird to  
<sup>164</sup> the grid format of the WorldClim variables, and then created a presence-absence community  
<sup>165</sup> matrix  $Y$ , taking all the grid cells as sites. At the 10 arc-minutes resolution, we obtained  
<sup>166</sup> 39 024 sites with occurrences and 62 species in total. All data manipulations and further  
<sup>167</sup> analyses were realized in *Julia v1.2.0* (Bezanson et al. 2017), with the basic structure built  
<sup>168</sup> around the soon-to-be-released `SimpleSDMLayers.jl` package.<sup>1</sup>

## <sup>169</sup> 3. SDM – The BIOCLIM Method

<sup>170</sup> We predicted species distributions using the BIOCLIM method (Nix 1986), a climate-envelope  
<sup>171</sup> model, considered a classic in the field. This method simply relates a species' distribution  
<sup>172</sup> to the ranges of bioclimatic variables at known locations (Booth et al. 2014). It has long  
<sup>173</sup> been outperformed by other methods (Elith et al. 2006), but it is still commonly used for  
<sup>174</sup> its simplistic approach and ease of understanding, as well as its simple relation to niche

---

<sup>1</sup><https://github.com/EcoJulia/SimpleSDMLayers.jl>

175 theory (Booth et al. 2014; Hijmans et al. 2017). It is also primarily designed for presence-only  
176 data. Despite that, we chose this method for our preliminary analyses as it was easier to  
177 implement and because we believe it to be sufficient for proof-of-concept. We discuss possible  
178 alternatives in the “Alternative Methods” section below.

179 The BIOCLIM method defines species potential ranges as a multidimensional environmental  
180 hypervolume bounded by the minimum and maximum values for all occurrences (Franklin  
181 2010). For each species, we established the percentile distribution of each environmental  
182 variable at the known locations of occurrence (Hijmans et al. 2017). All sites were then  
183 compared to those percentile distributions and given a score per variable according to  
184 their ranking between 0.0 (1st percentile) and 1.0 (100th percentile). The median or 50th  
185 percentile was considered the most suitable value of the variable, and values larger than  
186 0.5 were subtracted from 1. Therefore, both tails were considered the same. The minimum  
187 percentile score across all environmental variables was then selected as the predicted value  
188 for each site. Values were multiplied by 2 and could therefore be interpreted as probabilities  
189 of species occurrence (Hijmans et al. 2017). Predictions of 1 should be rare by definition,  
190 as they require a perfectly median site on all variables, and values of 0 should be frequent,  
191 since they occur whenever an environmental value is outside the range of the observed ones  
192 (Hijmans et al. 2017).

193 The final step was to convert the probabilities into presence-absence data, so they could  
194 be compared with the raw occurrence data. We transformed the probabilities into zeros  
195 and ones by converting all values greater than zero to one. Although it might tend to  
196 overestimate species ranges, such a transformation is common in SDMs and can be accounted  
197 for during result validation with specific methods (Franklin 2010). We also considered  
198 applying a threshold determined by sensitivity analysis, but we haven’t done it yet. In any  
199 case, converting into presence-absence data allowed easier calculation of the richness and  
200 beta diversity metric.

#### 201 4. LCBD Calculation

202 We calculated the LCBD statistics through the total variance of the matrix  $Y$  for both the raw  
203 data and SDM predictions. Legendre and De Cáceres (2013) showed that LCBD coefficients

204 can be calculated directly through the total variance of matrix  $Y$ , or through a matrix of  
205 dissimilarities among sampling units. We chose the first approach as it also allows to compute  
206 species contributions to beta diversity (SCBD), which could also prove useful for conservation  
207 purposes, but we did not investigate these for now. Before computing the LCBD statistics,  
208 the presence-absence matrix  $Y$  had to be transformed in an appropriate way (Legendre and  
209 De Cáceres 2013). We chose to apply the Hellinger transformation to the raw data and no  
210 transformation on the SDM predictions for now, although we did not investigate these in  
211 detail. The most appropriate transformation still needs to be determined, especially for  
212 the SDM predictions. We then computed a matrix  $S$  of squared deviations from column  
213 means and summed all the values of  $S$  to obtain the total sum of squares ( $SS$ ) of the species  
214 composition data (Legendre and De Cáceres 2013). LCBD coefficients are then computed  
215 as  $LCBD_i = SS_i/SS_{Total}$ , where  $SS_i$  is the sum of squares of a sampling unit  $i$ . Finally, since  
216 our matrix  $Y$  is very large, the LCBD coefficients are very small, so we scaled them to the  
217 maximum value observed.

## 218 **5. Prediction Validity**

219 The exact way of testing the validity of the predictions remains to be determined, and will  
220 also depend on the exact methods used to make the SDM predictions. A key element to note  
221 is that both SDM predictions and LCBD values will have to be validated, hence they might  
222 require different methods. Metrics that measure the accuracy of categorical or probabilistic  
223 predictions in SDMs are well documented, and take various forms. Some require absence  
224 data to test against, and can be used on probabilistic predictions directly (area-under-curve,  
225 AUC) or after a conversion of the predictions to binary presence-absence using a given  
226 threshold (Kappa index, measuring the difference between observed and chance agreement  
227 in a confusion matrix) (Franklin 2010). Other methods are appropriate for presence-only  
228 data, such as the Boyce Index. In any case, measuring prediction error is only one part  
229 of the validation. Finding appropriate data for evaluation is also critical (Franklin 2010),  
230 especially since we aim to describe community structure. Separating the data into training  
231 and testing datasets, with 70% and 30% of the observations for instance, is an approach  
232 common in machine learning methods. However, all of the available observations might  
233 be needed in some cases (Franklin 2010). An interesting approach, suggested by Elith et al.

234 (2006) for SDMs, would be to find independent, well-structured presence-absence datasets  
235 for validation, on which both SDM predictions and beta diversity metrics could be tested.  
236 This approach has the advantage that the testing data is truly independent of the training  
237 one, hence it could be used with certain tests of significance. Although it might not cover the  
238 entire extent of the predictions in a single test, this method would bring a closer comparison  
239 to the way LCBD metrics are used in most studies. Therefore, it would provide interesting  
240 perspectives if combined with other, full-extent validation methods.

## 241 **6. Alternative methods**

242 Many methods generally outperform BIOCLIM for the predictions, as shown by Elith et  
243 al. (2006). In our case, better predictions will come by two different means: 1) approaches  
244 that are better than BIOCLIM to model the relationship between species presence-absence  
245 (or even abundance) and environmental variables, and 2) approaches that account for other  
246 drivers of species distributions, such as ecological interactions and species migration. Machine  
247 learning methods, especially, would be interesting alternatives to consider. MAXENT (Phillips,  
248 Anderson, and Schapire 2006), another presence-only method, has come to be one of the most  
249 widely used methods in SDM studies, often with WorldClim variables (Booth et al. 2014).  
250 Similarly, Random Forests are simple to put in place, take into account both presence and  
251 absence data, allow for quantification of the variables importance in explaining variation, and  
252 offer intrinsic testing metrics (Franklin 2010). However, while those methods might return  
253 more accurate predictions, they do not implicitly model other drivers of species distribution,  
254 among which species interactions and functional niche. Integrating those factors might  
255 prove more difficult given our dataset and our focus on Warblers species, as no appropriate  
256 information on their interaction is available. Joint species distribution models (JSDMs) might  
257 be an interesting way to encompass those, as they attempt to model species co-occurrence,  
258 rather than the distribution of single species (Pollock et al. 2014). Also, a different taxonomic  
259 group and dataset with more details on interactions could simply be used. On the other hand,  
260 a method that could be applied to any taxonomic group, especially those well represented  
261 in large citizen-contributed datasets, would be most useful for research and conservation  
262 purposes.

263 **7. Climate Change Scenarios and Temporal Beta Diversity**

264 We aim to apply our method to environmental conditions from IPCC climate change scenarios.  
265 First, community compositions after climate change could be modelled on continuous scales  
266 through SDMs. Second, we could identify the sites where the community has changed in the  
267 most exceptional ways. This identification can be done by looking at the variation in LCBD  
268 values, but also through the use of temporal beta diversity indices (TBI) (Legendre 2019).  
269 TBI indices allow to study changes in community composition through time from repeated  
270 surveys at given sites. Whereas LCBD values essentially measure the contribution to beta  
271 diversity of one site compared to all others, TBI measure changes in community composition  
272 site-wise between two surveys. Moreover, TBI indices can be decomposed into species losses  
273 and gains, and can be tested for significance using a permutation test (Legendre 2019). An  
274 approach similar to that of Legendre and Condit (2019) would be interesting to follow in  
275 our case. First, they computed LCBD indices and compared the location of the sites with  
276 exceptional compositions between two surveys 30 years apart. The comparison showed that  
277 important changes seemed to have occurred in a specific swamp region. Then, they used TBI  
278 indices to confirm the sites with significant changes, decompose these changes into losses  
279 and gains, and identify the species that had changed the most. An approach such as this one  
280 could be highly informative with our data, although the permutation tests and corrections  
281 required might cause some problems given the number of sites in our study.

282 The possibility of using climate change scenarios in the SDMs also needs to be assessed. We  
283 did not try to download nor find the appropriate data for now. However, interpolated climate  
284 change variables are sometimes different than the ones in WorldClim. Therefore, the SDM  
285 models to use and the resulting predictions might have to be different too, and potentially  
286 less reliable. Nonetheless, we believe it will be possible to do some kind of time analysis  
287 linking beta diversity, climate change and species distribution modelling, and that it could  
288 return highly informative results for conservation purposes.

289 **Preliminary Results**

290 Our preliminary results consisted of comparisons between the raw occurrence data and  
291 the SDM predictions for the following elements: species richness (fig. 1), LCBD coefficients

292 (fig. 2), as well as the relationship between the species richness and LCBD coefficients (fig. 3).  
293 Two main results emerged from them: 1) the models provided seemingly valid community  
294 composition predictions for poorly sampled regions, both expected species-poor and species-  
295 rich, and 2) the relationship between species richness and species distribution models was  
296 in line with previous studies for species poor sites, but the SDM models captured a new  
297 association for very rich sites.

298 First, species richness, defined as the number of species present per site, showed a clear  
299 latitude gradient, with the poorest sites to the North and the richest to the South (fig. 1).  
300 A form of altitude gradient could also be observed, with the Rockies and other mountains  
301 well delimited by their lower values. In both cases, the results make intuitive sense and  
302 highlight the models ability to predict species presence despite poor or no sampling. Mexico,  
303 for example, has much sparser sampling and fewer observations, but the models predict  
304 higher species richness than on the highly sampled Atlantic Coast nonetheless, which make  
305 sense for a more southern location. We believe these to be valid insights on poorly sampled  
306 locations, but there is still a need for an appropriate method of validation to confirm our  
307 intuition, as well as a thoughtful consideration of factors such as species migration.

308 Second, our preliminary LCBD results seemed to confirm the association between species  
309 richness and LCBD coefficients from previous studies, but the SDM predictions captured  
310 a new association for extremely rich sites. Indeed, raw occurrence data showed a negative  
311 relationship between species richness and LCBD coefficients (fig. 3), as observed previously  
312 by Heino and Grönroos (2017), with no clear geographic pattern (fig. 2a). On the other  
313 hand, SDM predictions showed a clear geographic pattern, with the highest values to the  
314 northern and southern extremes (fig. 2b). Moreover, the richness-LCBD relationship showed  
315 a quadratic form, with the LCBD coefficients re-increasing beyond a richness of 0.6 (fig. 3),  
316 which approximately corresponds to the maximum richness observed in the raw occurrence  
317 data. Therefore, the SDM predictions captured a new association that could not be seen  
318 in the occurrence data, possibly because there were no rich enough sites to display it. By  
319 definition, LCBD indices should highlight the most exceptional species compositions, both  
320 species poor or species rich. Thus, the quadratic relationship we observed makes sense, on  
321 the condition that extremely rich sites can realistically exist. These richest sites should likely

322 have been in the southernmost locations such as Mexico, which is heavily undersampled.  
323 For instance, all but two species were seen there at least once, but the maximum number  
324 of species recorded in a single checklist was lower than in the US and in Canada (tbl. 1),  
325 which was a little surprising. Hence, some extremely rich communities might not have been  
326 sampled sufficiently to reveal their true community structure. On the other hand, our models  
327 might have been too optimistic in predicting the existence of such rich sites, but, in any case,  
328 our method did provide interesting ecological insights. The concurrence of our results, both  
329 occurrence and SDM based, with those of Heino and Grönroos (2017) for intermediate and  
330 species poor sites is promising. The possibility of finding new associations should therefore  
331 only encourage to push its use even further.

332 Finally, one disappointing aspect of our method is that the result failed to identify patterns  
333 on finer scales. The trends shown by the SDMs for both the species richness and LCBD  
334 coefficients were large-scale, latitude-related patterns. Except for mountains, few exceptional  
335 sites are actually shown in the middle of the landscape. While it might have been unrealistic  
336 to expect such results from a coarse analysis like ours, it would be useful for conservation  
337 purposes to be able to identify precise sites within smaller regions. This might be achieved  
338 by using a finer resolution, which we should probably reconsider in light of these results, or by  
339 using a different technique, such as training the models and predicting species distributions  
340 on large scales, but computing and scaling LCBD values on finer local ones, which might  
341 highlight regional differences in a new way.

342    **References**

- 343    Araújo, Miguel B., and A. Townsend Peterson. 2012. "Uses and Misuses of Bioclimatic  
344       Envelope Modeling." *Ecology* 93 (7): 1527–39. <https://doi.org/10.1890/11-1930.1>.
- 345    Bezanson, Jeff, Alan Edelman, Stefan Karpinski, and Viral B. Shah. 2017. "Julia: A Fresh  
346       Approach to Numerical Computing." *SIAM Review* 59 (1): 65–98. <https://doi.org/10.1137/141000671>.
- 348    Booth, Trevor H., Henry A. Nix, John R. Busby, and Michael F. Hutchinson. 2014. "BIOCLIM:  
349       The First Species Distribution Modelling Package, Its Early Applications and Relevance to  
350       Most Current MaxEnt Studies." *Diversity and Distributions* 20 (1): 1–9. <https://doi.org/10.1111/ddi.12144>.
- 352    Elith, Jane, Catherine H. Graham, Robert P. Anderson, Miroslav Dudík, Simon Ferrier, Antoine  
353       Guisan, Robert J. Hijmans, et al. 2006. "Novel Methods Improve Prediction of Species'  
354       Distributions from Occurrence Data." *Ecography* 29 (2): 129–51. <https://doi.org/10.1111/j.2006.0906-7590.04596.x>.
- 356    Fick, Stephen E., and Robert J. Hijmans. 2017. "WorldClim 2: New 1-Km Spatial Resolution  
357       Climate Surfaces for Global Land Areas." *International Journal of Climatology* 37 (12):  
358       4302–15. <https://doi.org/10.1002/joc.5086>.
- 359    Franklin, Janet. 2010. "Moving Beyond Static Species Distribution Models in Support of  
360       Conservation Biogeography: Moving Beyond Static Species Distribution Models." *Diversity  
361       and Distributions* 16 (3): 321–30. <https://doi.org/10.1111/j.1472-4642.2010.00641.x>.
- 362    Heino, Jani, and Mira Grönroos. 2017. "Exploring Species and Site Contributions to Beta  
363       Diversity in Stream Insect Assemblages." *Oecologia* 183 (1): 151–60. <https://doi.org/10.1007/s00442-016-3754-7>.
- 365    Hijmans, Robert J., Susan E. Cameron, Juan L. Parra, Peter G. Jones, and Andy Jarvis. 2005.  
366       "Very High Resolution Interpolated Climate Surfaces for Global Land Areas." *International  
367       Journal of Climatology* 25 (15): 1965–78. <https://doi.org/10.1002/joc.1276>.
- 368    Hijmans, Robert J., Steven Phillips, John Leathwick, and Jane Elith. 2017. *Dismo: Species  
369       Distribution Modeling*. <https://CRAN.R-project.org/package=dismo>.

- 370 Johnston, A., W. M. Hochachka, M. E. Strimas-Mackey, V. Ruiz Gutierrez, O. J. Robinson,  
371 E. T. Miller, T. Auer, S. T. Kelling, and D. Fink. 2019. "Best Practices for Making Reli-  
372 able Inferences from Citizen Science Data: Case Study Using eBird to Estimate Species  
373 Distributions." *bioRxiv*, March, 574392. <https://doi.org/10.1101/574392>.
- 374 Legendre, Pierre. 2019. "A Temporal Beta-Diversity Index to Identify Sites That Have  
375 Changed in Exceptional Ways in Space-Time Surveys." *Ecology and Evolution* 9 (6): 3500–  
376 3514. <https://doi.org/10.1002/ece3.4984>.
- 377 Legendre, Pierre, Daniel Borcard, and Pedro R. Peres-Neto. 2005. "Analyzing Beta Diver-  
378 sity: Partitioning the Spatial Variation of Community Composition Data." *Ecological  
379 Monographs* 75 (4): 435–50. <https://doi.org/10.1890/05-0549>.
- 380 Legendre, Pierre, and Richard Condit. 2019. "Spatial and Temporal Analysis of Beta Diversity  
381 in the Barro Colorado Island Forest Dynamics Plot, Panama." *Forest Ecosystems* 6 (1): 7.  
382 <https://doi.org/10.1186/s40663-019-0164-4>.
- 383 Legendre, Pierre, and Miquel De Cáceres. 2013. "Beta Diversity as the Variance of Community  
384 Data: Dissimilarity Coefficients and Partitioning." *Ecology Letters* 16 (8): 951–63. <https://doi.org/10.1111/ele.12141>.
- 386 Nix, Henry A. 1986. "A Biogeographic Analysis of Australian Elapid Snakes." *Atlas of Elapid  
387 Snakes of Australia* 7: 4–15.
- 388 Phillips, Steven J., Robert P. Anderson, and Robert E. Schapire. 2006. "Maximum Entropy  
389 Modeling of Species Geographic Distributions." *Ecological Modelling* 190 (3): 231–59.  
390 <https://doi.org/10.1016/j.ecolmodel.2005.03.026>.
- 391 Poisot, Timothée, Richard LaBrie, Erin Larson, Anastasia Rahlin, and Benno I. Simmons. 2019.  
392 "Data-Based, Synthesis-Driven: Setting the Agenda for Computational Ecology." *Ideas in  
393 Ecology and Evolution* 12 (July). <https://doi.org/10.24908/iee.2019.12.2.e>.
- 394 Pollock, Laura J., Reid Tingley, William K. Morris, Nick Golding, Robert B. O'Hara, Kirsten  
395 M. Parris, Peter A. Vesk, and Michael A. McCarthy. 2014. "Understanding Co-Occurrence  
396 by Modelling Species Simultaneously with a Joint Species Distribution Model (JSDM)." *Methods in Ecology and Evolution* 5 (5): 397–406. <https://doi.org/10.1111/2041-210X.121>

399 Sullivan, Brian L., Christopher L. Wood, Marshall J. Iliff, Rick E. Bonney, Daniel Fink, and  
400 Steve Kelling. 2009. "eBird: A Citizen-Based Bird Observation Network in the Biological  
401 Sciences." *Biological Conservation* 142 (10): 2282–92. <https://doi.org/10.1016/j.biocon.2009.05.006>.

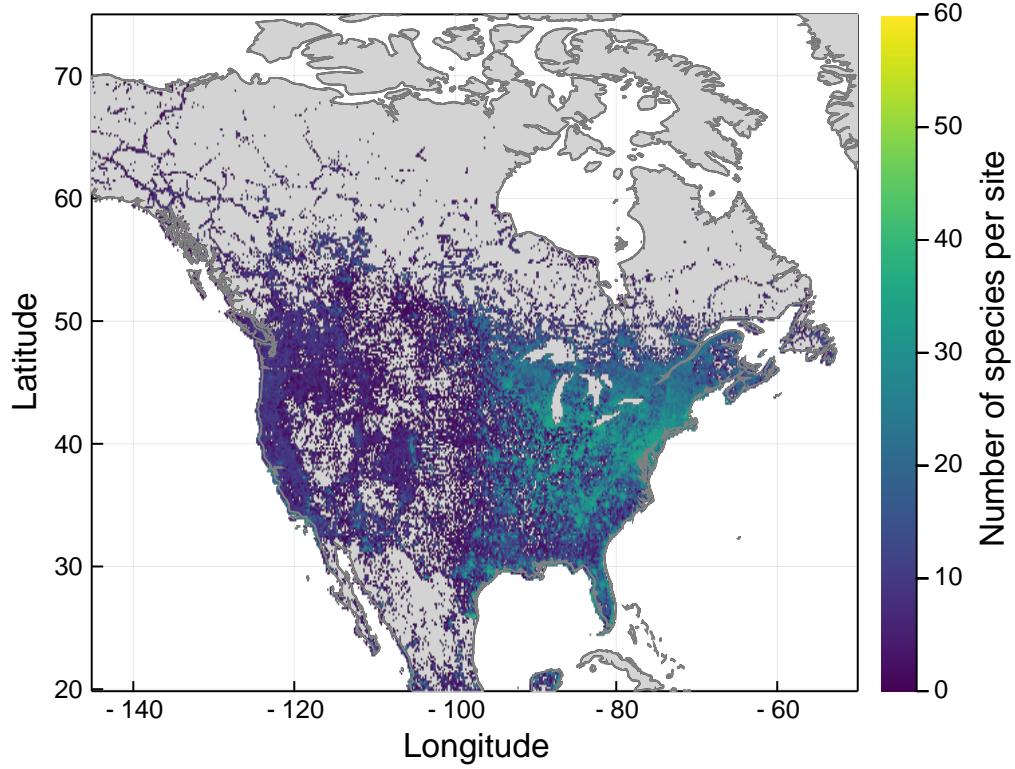
<sup>403</sup> **Appendix**

Table 1: Structure of the Warblers data in the eBird checklists for the countries used in the analyses

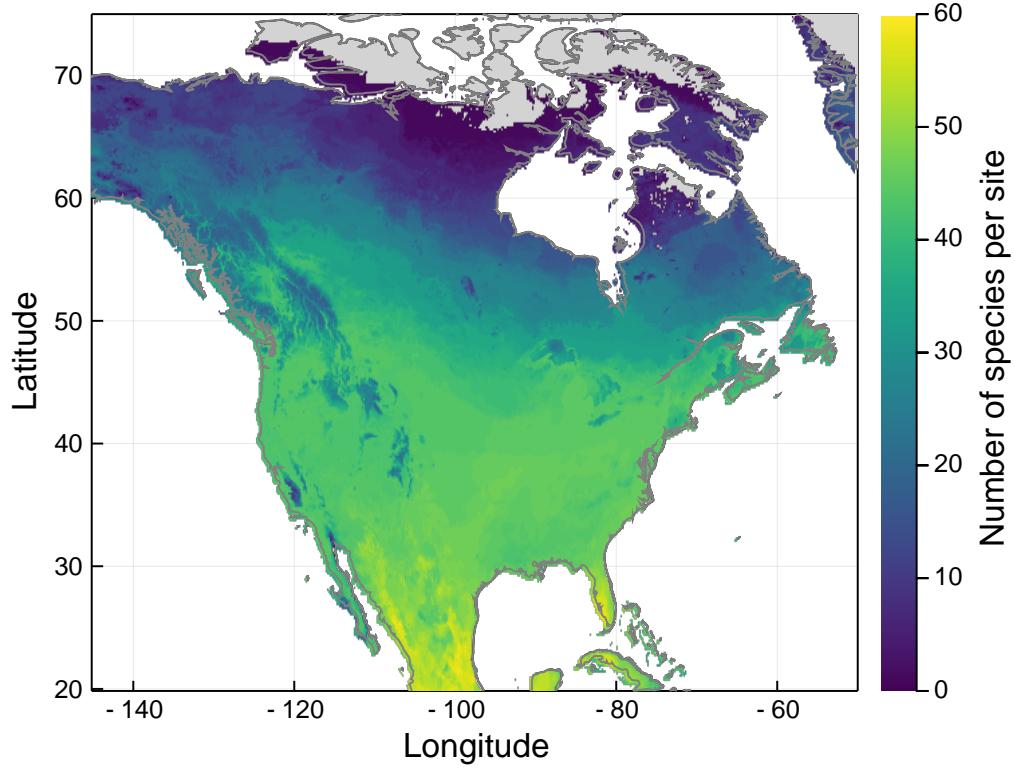
Country	Observations	Checklists	Species	Species per checklist (mean)	Species per checklist (median)	Species per checklist (maximum)
US	19 206 453	7 840 526	56	2.450	2.0	34
CA	3 360 650	1 115 625	45	3.012	2.0	31
MX	407 227	147 599	61	2.759	2.0	21
Total	22 974 330	9 103 750	63	2.523	2.0	34

Table 2: Description of the WorldClim 2 climate variables used in the analyses

Variable	Description
1	Annual Mean Temperature
2	Mean Diurnal Range (Mean of monthly (max temp - min temp))
3	Isothermality (BIO2/BIO7) (* 100)
4	Temperature Seasonality (standard deviation *100)
5	Max Temperature of Warmest Month
6	Min Temperature of Coldest Month
7	Temperature Annual Range (BIO5-BIO6)
8	Mean Temperature of Wettest Quarter
9	Mean Temperature of Driest Quarter
10	Mean Temperature of Warmest Quarter
11	Mean Temperature of Coldest Quarter
12	Annual Precipitation
13	Precipitation of Wettest Month
14	Precipitation of Driest Month
15	Precipitation Seasonality (Coefficient of Variation)
16	Precipitation of Wettest Quarter
17	Precipitation of Driest Quarter
18	Precipitation of Warmest Quarter
19	Precipitation of Coldest Quarter

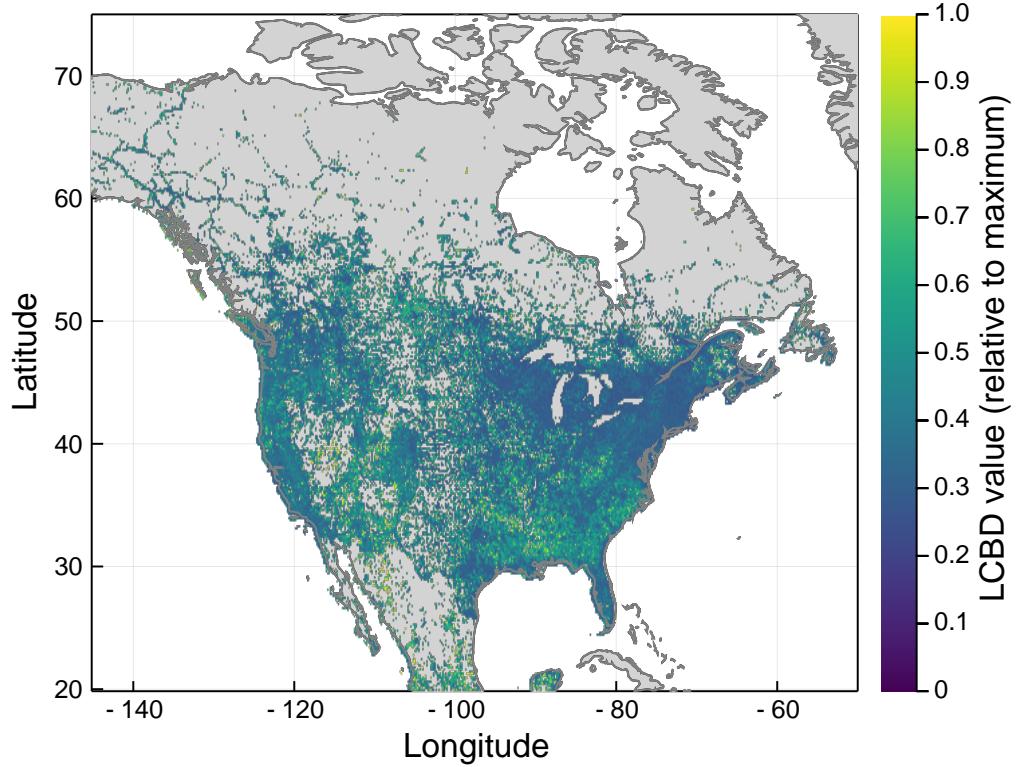


(a) Raw occurrence data

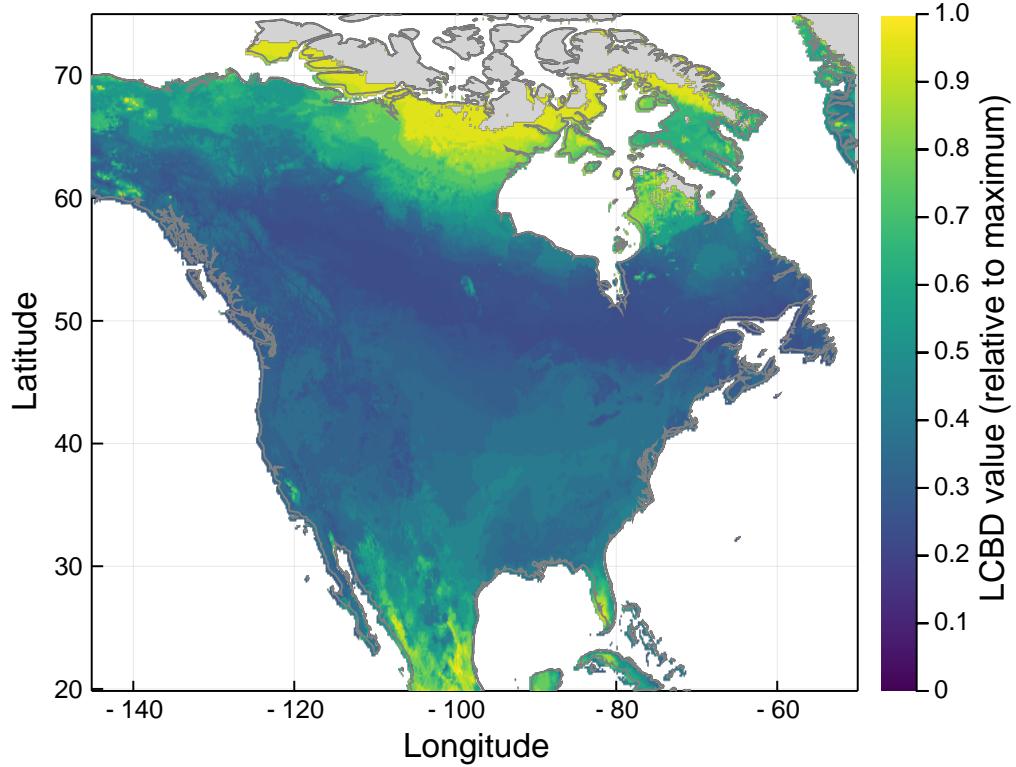


(b) SDM predictions

Figure 1: Distribution of species richness in North America, defined as the number of Warblers species per site. The raw occurrence observations from eBird (fig. 1a) and the SDM predictions from the BIOCLIM model (fig. 1b) were both transformed into presence-absence data per species before calculating richness.



(a) Raw occurrence data (Hellinger transformed)



(b) SDM predictions

Figure 2: Distribution of the LCBD values in North America, calculated from the variance of the community matrix  $Y$  and scaled to the maximum value observed. The Hellinger transformation was applied on the raw occurrence data (fig. 2a) before calculating the LCBD indices. SDM predictions (fig. 2b) were converted into presence-absence data, but no transformation was applied before calculating the LCBD indices.

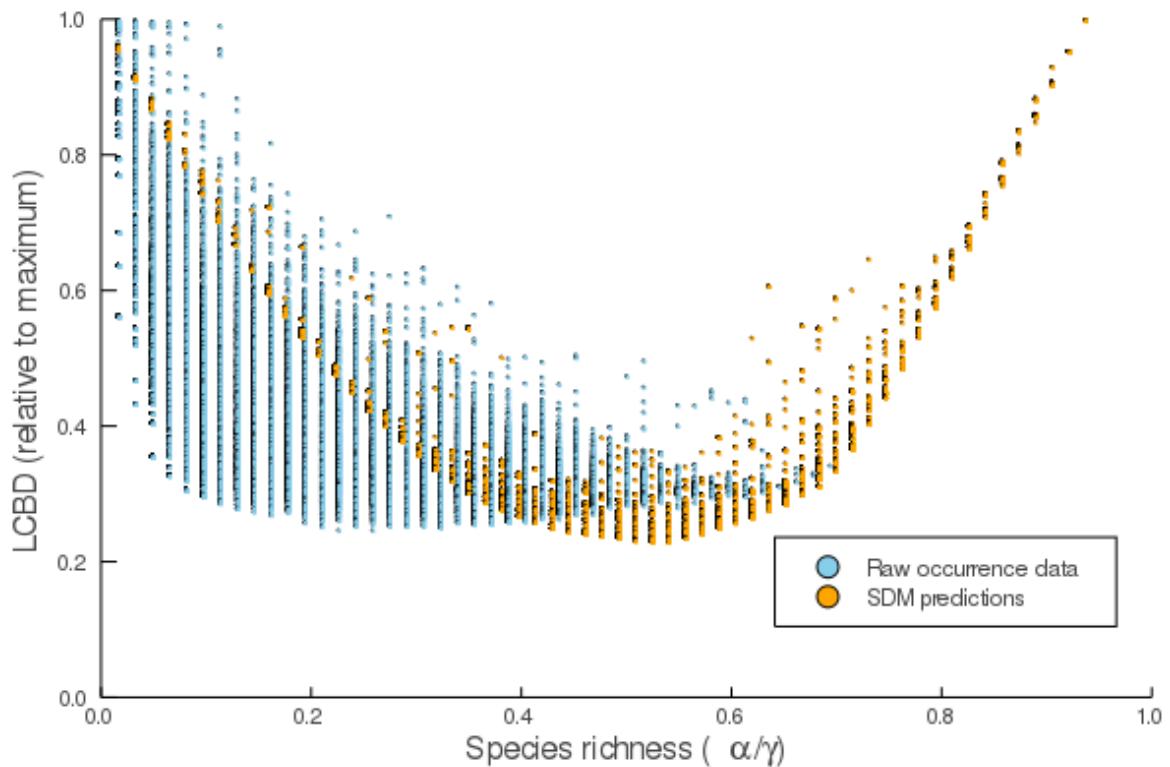


Figure 3: Relationship between the species richness and the LCBD value of the each site for raw occurrence data (blue) and SDM predictions (orange). Species richness was calculated as the number of species in a site ( $\alpha$ ), divided by the total number of species ( $\gamma$ ). LCBD values were scaled to the maximum value observed. Hellinger transformation was applied on the raw occurrence data before calculating LCBD indices.