

1 **Extracting socio-spatial networks from photo-ID data using multilevel**

2 **multinomial models**

3

4 Tyler R. Bonnell <sup>a,b,c</sup>, Robert Michaud <sup>d</sup>, Angélique Dupuch <sup>a,b</sup>, Véronique Lesage<sup>e</sup>, Clément

5 Chion <sup>a,b</sup>

6

7

8

9

10 <sup>a</sup> Department of Natural Sciences, Université du Québec en Outaouais, Gatineau, Québec,

11 Canada

12 <sup>b</sup> Institut des Sciences de la Forêt Tempérée, Université du Québec en Outaouais, Ripon,

13 Québec, Canada

14 <sup>c</sup> Department of Psychology, University of Lethbridge, 4401 University Drive, Alberta T1K

15 3M4, Canada

16 <sup>d</sup> Groupe de Recherche et d'Éducation sur les Mammifères Marins (GREMM), 870 avenue

17 Salaberry, Bureau R24, Québec, QC G1R 2T9, Canada

18 <sup>e</sup> Maurice Lamontagne Institute, Fisheries and Oceans Canada, P.O. Box 1000, 850 route de

19 la Mer, Mont-Joli, QC G5H 3Z4, Canada

20

21

22 **Abstract:**

23 Photo identification of individuals within a population is a common data source that is  
24 becoming more common given technological advances and the use of computer vision and  
25 machine learning to re-identify individuals. These data are collected through hand-held  
26 cameras, drones, and camera traps, and often come with biases in terms of sampling effort  
27 and distribution. In spite of these biases, a common goal of collecting these datasets is to  
28 better understand the habitat use pattern of individuals and populations. Here, we examine the  
29 potential for multilevel multinomial models to generate socio-spatial networks that capture  
30 the similarities in individual users across the spatial distribution of a species. We use this  
31 approach with 18 years of photo-ID data to better understand population structuring of beluga  
32 whales in the St. Lawrence River. We show using permuted and simulated data that this  
33 approach can identify community network structures within populations in a way that  
34 accounts for biases in collections methods. Applying this method to the entire 18 years  
35 dataset for SLE beluga, we found three spatially distinct communities. These results suggest  
36 that within the population's summer range individuals are moving within restricted areas (i.e.,  
37 home ranges), and have implications for the estimated impacts of localized anthropogenic  
38 stressors, such as chemical pollution or acoustic disturbances on animal populations. We  
39 conclude that multilevel multinomial models can be effective at estimating socio-spatial  
40 networks that describe community structuring within wildlife populations.

41

42 **Keywords:** Multinomial Model, Beluga, Photo ID, Socio-Spatial Network, Bayesian

43 Network, Community Detection

## 44        **1. Introduction**

45            An understanding of the spatial and temporal distribution of a species of concern is of  
46        central importance to conservation and management (Evans & Hammond, 2004).

47        Increasingly, photo and video are being used to monitor individuals within populations  
48        (hereafter photo-ID data), providing a view of within population social mixing and habitat  
49        use (Koivuniemi et al., 2016). The increased use of machine learning to identify individuals  
50        from these data streams has greatly facilitated the use of these photo-ID data (Schneider et  
51        al., 2019). These individual identifications have facilitated the use of novel statistical and  
52        computational methods to quantify within population structures, such as social network  
53        analysis (Perryman et al., 2019; Schilds et al., 2019).

54            It is often the case, however, that efforts when collecting photo-ID are not evenly  
55        distributed. This differentiation in effort can heavily bias estimates of both habitat usage and  
56        population distribution estimates (Hupman et al., 2018). Here we propose the novel use of  
57        multilevel multinomial models to account for these biases and to estimate socio-spatial  
58        structures within populations.

59            The existence of social structuring within populations, such as communities, can have  
60        important ecological and management implications. If a population as a whole can be  
61        considered as highly mixed, i.e., with individuals showing no strong patterns of home range  
62        use or sub-structuring within the large population, then all individuals are equally likely to  
63        feel the impacts of local changes in the environment. In contrast, if the population cannot be  
64        considered to be highly mixed, and shows strong sub-structuring and site-fidelity patterns  
65        within the larger population, local stressors might have a disproportionate impact on  
66        subsections of the population. For example, if noise pollution increased in only one sector, in  
67        a highly mixed population all individuals would be lightly impacted, but in a structured  
68        population a subset of the population would be highly impacted. These differences in spatial

69 structuring of populations can lead to biased estimation of the likelihood and magnitude of  
70 impacts from local stressors both at the individual and population levels (DeFur et al., 2007).

71         The multilevel multinomial modeling approach does not have a large body of  
72 literature to draw on for use with photo-ID data. However, it does have some unique  
73 advantages (Koster & McElreath, 2017). For instance, if sampling effort is biased in different  
74 regions, the mean probability of being seen within highly sampled regions will be biased  
75 upwards. By taking advantage, however, of the multilevel structure of the model it is possible  
76 to extract individual deviations in the probability of being seen within a particular region.  
77 Decisively, these individual level deviations from the mean probability are not biased by  
78 changes in the sampling effort. That is the mean probability will increase with sampling  
79 effort, but the relative difference between individuals within the sector will not. High users of  
80 a particular region of a habitat will consistently be higher compared to low users of that  
81 habitat, and this difference between high and low users will not be biased by sampling effort.  
82 Furthermore, by comparing the individual differences between regions it is possible to see if  
83 the high/low users of one region are similarly the high/low users of another region. The  
84 similarity, or dissimilarity, between regions can then provide information about which  
85 regions share similar user profiles. We suggest that by using the correlations between these  
86 individual level deviations in high/low users between regions it is possible to generate socio-  
87 spatial networks and identify social structuring within the population. In particular it can help  
88 to identify spatial communities, i.e., a set of regions that share similar usage patterns and that  
89 differ from other regions.

90         To evaluate the use of multilevel multinomial models to identify socio-spatial  
91 structuring within a population, we make use of a long term photo-ID dataset of beluga  
92 whales in the St. Lawrence Estuary, Canada. This population has undergone a drastic  
93 decrease from around 10,000 in the late 1800s to less than 1,000 today, and is currently

94 considered as endangered in Canada according to the Species at Risk Act (COSEWIC 2014;  
95 Fisheries and Oceans Canada, 2012; Mosnier et al., 2015). The population is part of a larger  
96 study on the mitigation of noise pollution due to marine traffic (Chion et al., 2017; Lesage et  
97 al., 2014; McQuinn et al., 2011; Parrott et al., 2011).

98 In this paper we first evaluate the performance of the multilevel multinomial models  
99 using simulated data, testing if the method correctly estimates no community structuring  
100 when none is present, and identifies the correct structure when it is present. In both cases we  
101 use the 18 year beluga photo-ID dataset, randomly permuting uniquely identified individuals  
102 to generate unstructured datasets, and randomly placing individuals within pre-specified  
103 communities to generate structured datasets. We then apply the method to the observed data  
104 and quantify community structures within the population's summer range in the St. Lawrence  
105 Estuary. Finally, we discuss some potential extensions to the multilevel multinomial  
106 modeling approach.

107

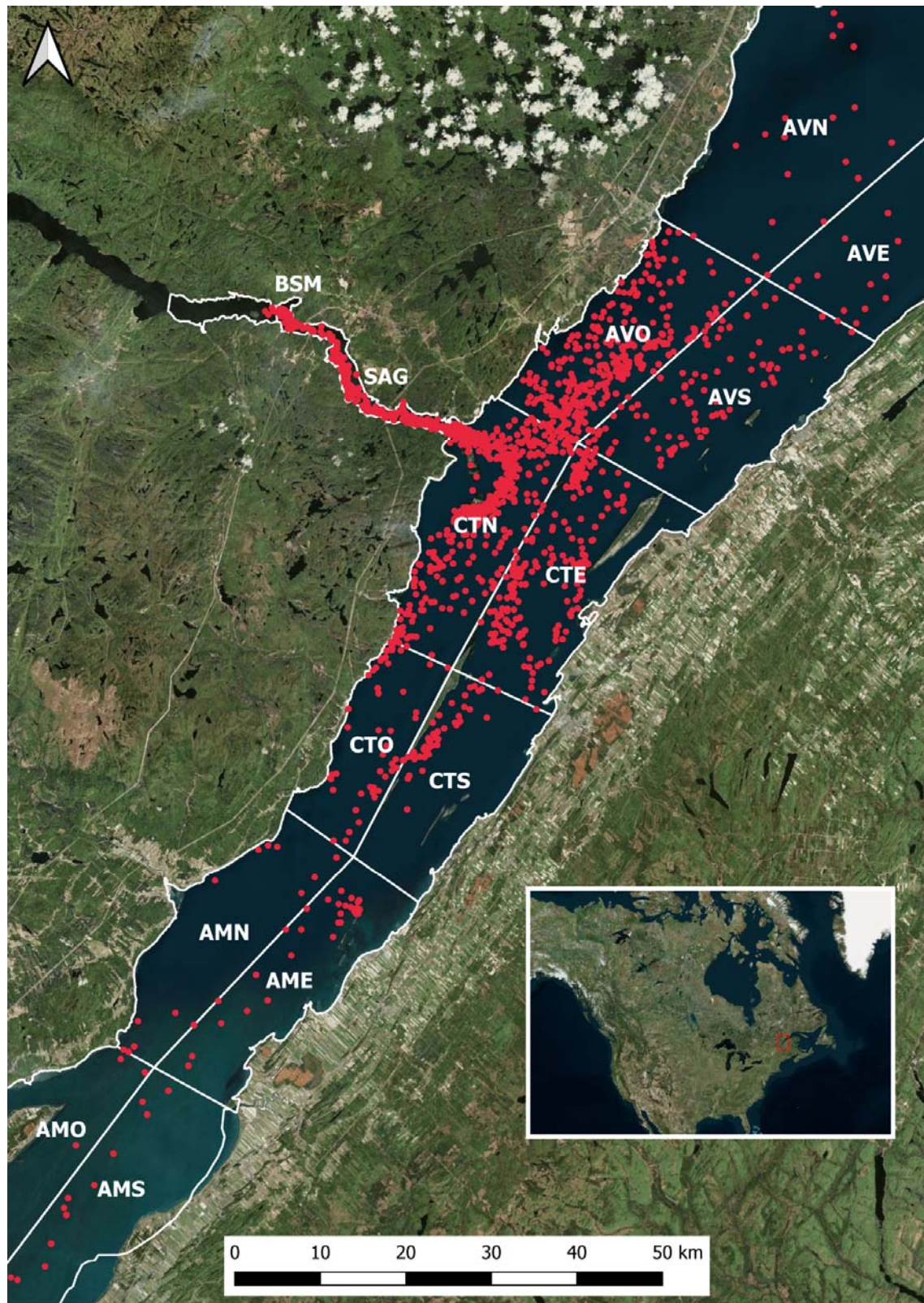
## 108 **2. Material and Methods**

### 109 *2.1 Data*

110 Photo-ID data were collected using a handheld camera onboard a boat that was able to  
111 navigate near to beluga, hereafter referred to as an encounter. Once near beluga a photo ID  
112 protocol was then followed to generate photographs used to attempt to identify individuals.  
113 Due to the logistical difficulty of covering a large body of water, the sampling effort was  
114 unequally distributed across 14 sectors within the St. Lawrence Estuary (Fig. 1). The photo-  
115 ID dataset used in this study was collected from 1989-2007 and are part of an ongoing  
116 project. The data is stored in a database that facilitates the identification and association  
117 between photos to help track the individual identification process. This resulted in a dataset  
118 of 7,525 individual encounters where the individual was successfully identified by photo

119 sampling (hereafter referred to as a photo-ID), and where a GPS point was taken and the  
120 sector recorded. This resulted in 821 unique individuals being successfully identified, with a  
121 mean number of photo-IDs per individual of 9 (min = 1, max = 90) (Fig. 1).

122



124 Figure 1: Spatial distribution of photo-ID data (red points) within the St. Lawrence Estuary,  
125 Quebec, Canada (red square in the inset map). The 14 sectors are outlined and labeled in  
126 white, and covers the summer habitat of this beluga population.

127

## 128 *2.2 Statistical Analysis*

129 Our aim was to use photo-ID data to estimate the probability of seeing an individual  
130 in each delineated sector of the St. Lawrence Estuary, and to use these individual  
131 probabilities to estimate socio-spatial structures within the beluga population (Fig. 1). To  
132 accomplish this aim we used a multilevel multinomial model, where the dependent variable  
133 was the number of times an individual was captured photographically (i.e., photo-identified)  
134 in each sector. The use of a multilevel model structure allows for the estimation of both the  
135 mean probability of photo-identifying an individual in each sector, and the individual level  
136 differences in this probability by using individual ID as a random intercept. If we take, as an  
137 example, a case where there is only two sectors, then the log-odds of finding individual  $i$  in a  
138 sector other than the reference sector can be modeled using a multilevel multinomial  
139 following (Koster & McElreath, 2017) as:

$$\log\left(\frac{p_{1,i}}{p_{r,i}}\right) = \mu_1 + v_{1i}$$
$$\log\left(\frac{p_{2,i}}{p_{r,i}}\right) = \mu_2 + v_{2i}$$

140 Where  $p_{1,i}$  is the probability of seeing beluga  $i$  in sector 1, whereas  $p_{r,i}$  is the  
141 probability of seeing beluga  $i$  in the reference sector. The  $\mu_1$  and  $\mu_2$  are the intercepts, i.e., the  
142 mean probability of seeing a beluga in sectors 1 and 2. This mean probability represents  
143 preference/avoidance of the specified sector, however, it is very likely to be biased due to  
144 variation in sampling effort. Finally, the  $v_{1i}$  and  $v_{2i}$  are the estimated individual differences  
145 (i.e., random intercepts) from the mean probability of capture in sectors 1 and 2, respectively.



146 These individual differences from the mean probability capture if an individual beluga is a  
147 high/low user of that sector, and is not biased by variation in sampling. It is then possible to  
148 model the covariance of the individual differences between two sectors using a multivariate  
149 normal distribution:

$$\begin{bmatrix} v_{1i} \\ v_{2i} \end{bmatrix} \sim \text{Multi Normal}(0, \Omega_v)$$

$$\Omega_v = \begin{bmatrix} \sigma_{1,1} & \sigma_{1,2} \\ \sigma_{2,1} & \sigma_{2,2} \end{bmatrix}$$

150 This multivariate normal distribution has a mean of 0 and a covariance matrix  $\Omega_v$ .  
151 Here the diagonal entries in the covariance matrix (e.g.,  $\sigma_{1,1}$ ) represent the magnitude of  
152 individual differences within a sector, i.e., are their high and low users in a sector or are all  
153 individuals equally likely to be seen? The off-diagonal entries (e.g.,  $\sigma_{2,1}$ ) are the covariance  
154 estimates between sectors, i.e., do sectors share similar high and low users? By estimating the  
155 correlation of individual differences between sectors, this multilevel modeling approach  
156 quantifies how much information individual differences in one sector can provide about  
157 another sector. Positive correlations suggest that the high/low users in one sector are similarly  
158 high/low users in another sector, while negative correlations suggest high/low users in one  
159 sector are the low/high users in another sector.

160 This model can be fit using brms in the R environment using a multivariate syntax:  
161 `bf(y | trials(n) ~ 1 + (1|q|ID)) + multinomial()`. Here, y is a set of column vectors where each  
162 column is a sector and each row is an individual. The values in this column vector indicate  
163 how many times each individual was seen in each sector. The n is the total number of times  
164 an individual was captured, and q is an arbitrary character choice that allows correlations  
165 between the estimates of random intercepts for each sector (Bürkner, 2017).

166

167 *2.2.1 Addition of a common reference sector*

168 Tests using simulated data suggested that adding a preset reference category (i.e., a  
169 fifteenth sector used as a reference sector) to the data was required to estimate the correlation  
170 between sectors (Fig. S1). To create this reference category, we tally up the observations for  
171 each individual and add that number of observations to the new reference sector. This  
172 essentially sets the probability of a capture in the reference sector to 0.5 for all individuals,  
173 i.e., equal to the probability of being captured outside of this reference sector. This ensures  
174 that all individuals have the same baseline probability in the reference category, and as the  
175 parameters in the multinomial model measure deviations away from the reference sector, we  
176 gain better estimates of the relative deviations between individuals (Fig. S1).

177

#### 178 *2.2.2 Dealing with biases in photo-ID datasets*

179 This multilevel multinomial approach accounts for repeated sampling of individuals,  
180 and provides an estimate of whether some individuals are found more or less often than the  
181 mean probability of captures in each sector. We are particularly interested in the estimates of  
182 individual differences from the mean probability of capture (i.e., the random intercepts) as  
183 these estimates are not impacted by bias in sampling effort among sectors. This is not the  
184 case for estimates of the mean probability of capture for each sector, which are expected to  
185 increase in highly sampled sectors. For example, the Saguenay River is over-sampled  
186 compared to the other sectors (SAG in Fig. 1), increasing the mean probability of capturing  
187 individuals in that sector. However, over-sampling should not affect the individual  
188 differences in the probability of being captured, i.e., all individuals' chances of being captured  
189 go up or down equally.

190 Similarly, potential biases due to ease of recognition, e.g., some beluga or age classes  
191 might have more distinctive markings, are minimized using a multilevel multinomial  
192 approach where the differences in the probability of being seen between sectors is the main

193 focus. For example, if juveniles are less likely to be successfully identified by photo, then  
194 they might have a reduced number of photo-IDs compared to other age classes, but the  
195 difference in distribution of these fewer photo-IDs across sectors will not be impacted. For  
196 example, if both an adult and juvenile spend twice as much time in the SAG sector compared  
197 to all other sectors, you might expect a photo-ID distribution (Sag:not-sag) of 10:5 and 2:1,  
198 respectively. In both cases the probability of being captured in the SAG sector is twice that of  
199 the remaining sector. Due to the adaptive partial pooling properties of multilevel models,  
200 individuals with few photo-IDs will, however, be less likely to show differences to the mean  
201 probability, i.e., they contain less information. This means that if an age class has very little  
202 chance of being identified by photo-ID, they will likely contribute less to the estimated socio-  
203 spatial structures estimated by the multilevel multinomial approach.

204 By using a multilevel modeling approach we also reduce the chance of false positives  
205 when making comparisons between many different individuals in many different sectors (i.e.,  
206 problem of multiple comparisons). For example, if we were to estimate the differences in the  
207 probability of each sector separately for each individual, the risk of false positives would be  
208 increased. By using a multilevel approach to estimating the differences we can make effective  
209 use of partial pooling of information to reduce extreme values, especially where the number  
210 of photo IDs is not equal between individuals. Furthermore, by running this in a Bayesian  
211 framework we are able to place priors on the individual differences within sectors that start  
212 the model assuming that there are no differences between individuals in their use of each  
213 sector, e.g.,  $\text{student}_t(3,0,1)$ .

214

### 215 *2.3 Social Network Analysis*

216 Social networks are often used when visualizing and quantifying social structures  
217 within populations, with individuals often represented as nodes and their interactions as edges

218 between these nodes. In our case, we use sectors as nodes, and the similarities in user profiles  
219 between sectors as edges. The correlations between sectors estimated from the multilevel  
220 multinomial model can be used to create a network where the posterior predictions of each  
221 correlation parameter corresponds to an edge weight in the network. In this way, each edge  
222 has a posterior distribution and can be used to create many networks from which a  
223 distribution of network metrics can be generated, e.g., the distribution of node strength values  
224 can be calculated for each sector. The advantage of having distributions of network measures  
225 is that the measures can be readily compared, e.g., does one sector have a higher node  
226 strength than another? It is also possible to use the distribution of edge weights, and a chosen  
227 threshold (e.g., 95% credible interval), to highlight only the edges where the sign of the  
228 correlation is known with a particular range of certainty. In this paper, we used this latter  
229 approach to generate a signed network (i.e., a network with positive and negative edges) and  
230 use a simple signed-edge rule to define communities: where a distinct community is a set of  
231 nodes that share positive edges but no negative edges. We also made use of signed  
232 blockmodeling, an algorithm that can also be used to identify blocks of nodes that maximize  
233 within block positive edges and minimize within block negative edges (Doreian & Mrvar,  
234 2015). While the signed-edge rule generally provides relatively intuitive results with simple  
235 networks, using the signed blockmodeling is likely to be particularly advantageous when  
236 dealing with large networks.

237

#### 238 *2.4 Testing data*

239 To assess the accuracy of the multilevel multinomial modelling approach, we generated  
240 test datasets from the observed photo-ID data. We ensured that the test datasets contained the  
241 same number of unique individuals, distribution of sightings (i.e., some individuals are seen  
242 more than others), and overall number of photo-IDs as the observed dataset. We, however,

243 varied the spatial mixing of the test datasets. To test if the proposed method correctly  
244 detected no pattern when none existed, we created a completely random test dataset by  
245 permuting the sector associated with each photo-ID in the observed dataset. The expected  
246 result was to find no correlations between sectors given that the sectors for each photo-ID had  
247 been randomly permuted. To then test whether the proposed method could also correctly  
248 identify patterns when a known pattern existed, we generated a structured test dataset by  
249 randomly assigning each individual to four equally populated communities with the  
250 following and hypothetical home range of adjacent sectors: community 1-BSM, SAG, CTN,  
251 community 2- CTN, CTO, AMN, community 3- AVO, AVS, AVN, and community 4-AME,  
252 CTS, CTE. Following this, we altered the sector of where the individual photo-IDs were  
253 taken so as to fall within sectors associated with an individual's community, i.e., one of their  
254 home range sectors. We did this by choosing a sector for each photo-ID based on the  
255 individual's assigned community 80% of the time; a random sector was chosen for the other  
256 20% of the time, introducing noise in the assignment of sectors. We then tested whether the  
257 model correctly identified the correlations between sectors that defined the home range of  
258 each of the communities.

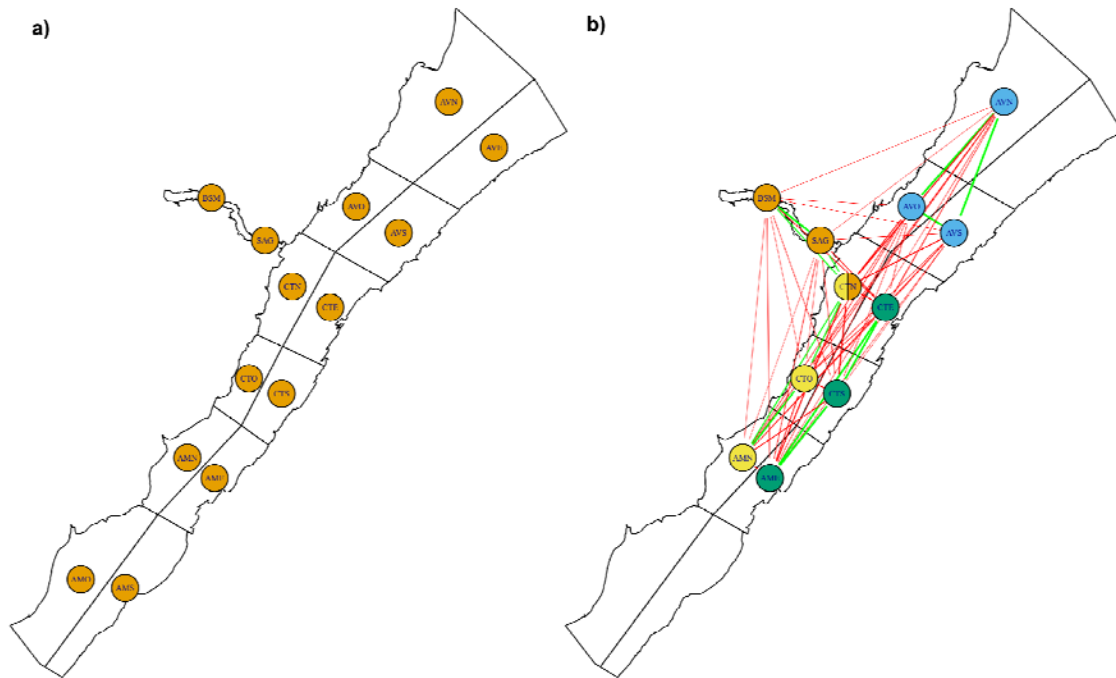
259

### 260 **3. Results**

#### 261 *3.1 Testing data*

262 When the multinomial multilevel model was fit to the data with sectors randomly  
263 permuted between all photo-IDs, the model found no evidence for positive/negative  
264 correlations between sectors (Fig. 2a). Similarly, when we simulated data with some known  
265 structure, i.e., when we artificially created spatially distinct communities, we found that the  
266 model accurately estimated the correlations between sectors that defined these artificial  
267 communities (Fig. 2b). The simple signed-edge rule and blockmodelling algorithm applied to

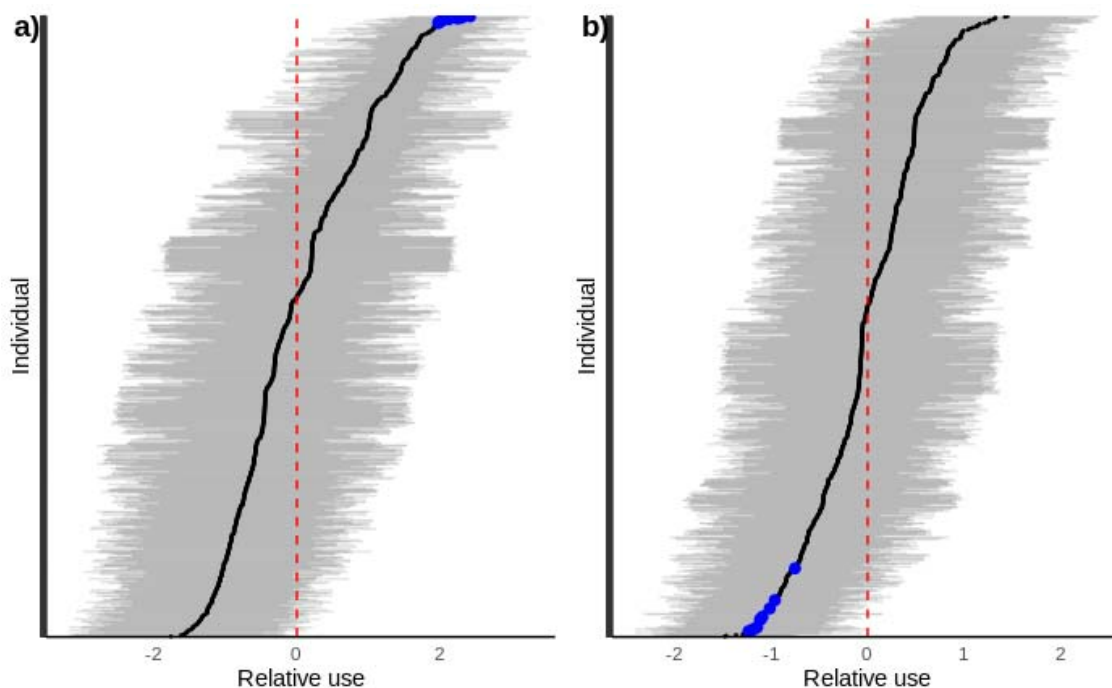
268 the simulated datasets both revealed the four artificially generated communities, though the  
269 blockmodeling algorithm had difficulty with the multi-membership node as it could not  
270 assign a node to two blocks (i.e., the CTN node that was shared between communities 1 and  
271 2).  
272



273  
274 Figure 2: Similarity and dissimilarity between sectors in the simulated datasets: a) randomly  
275 permuted data, where there are no spatial communities, and b) structured data, where there  
276 are four distinct communities. In b) the simulated communities are represented by color codes  
277 for each of their sectors (Note: CTN is part of the orange and yellow communities). The  
278 green edges between two sectors signify that the sectors share high/low users, while red  
279 edges signify that they have dissimilar high/low users. The lack of an edge signifies that the  
280 high/low users of one sector does not provide information about the high/low users of other  
281 sectors. Nodes represent sectors, and are coloured based on the communities imposed when  
282 simulating the data.  
283

284 3.2 Observed data

285 The results from our multilevel multinomial model found that within sectors there  
286 were consistent individual differences in how likely it was to see individual beluga (Table 1).  
287 That is, within sectors, there were some beluga that used the sector heavily, while others did  
288 not. The model also found that between sectors these individual differences were correlated  
289 (Table S1). These correlations quantify the magnitude of similarity/dissimilarity between  
290 sectors in terms of which beluga are using those sectors heavily or rarely. If we take two  
291 sectors as examples, e.g., the SAG and CTE sectors, representing, respectively, a tributary to  
292 the St. Lawrence Estuary and a sector on the opposite side close to the South shore of the  
293 Estuary, and we look at the top 10 estimated high users (i.e., relatively high probability of  
294 being found there) within the SAG, we find that they are found to be low users in the CTE  
295 sector (see blue dots in Fig 3 a) and b)).



296 Figure 3: Estimate of the relative use (i.e., deviation from mean use) for each individual  
297 within the SAG (a) and CTE sectors (b) of the St. Lawrence beluga summer habitat. The  
298 values are deviations from the mean probability of observing individuals within a sector and  
299

300 are on the logit scale. The red dashed line represents the mean use, black points represent the  
301 estimated deviation from the mean, while the horizontal grey lines represent the 95% credible  
302 interval. To highlight how correlations are estimated between sectors, the estimated top 10  
303 users of the SAG sector are represented by blue dots (panel a), and those same individuals are  
304 also highlighted in blue in the CTE sector (panel b).

305

306 The use of a multilevel model also allowed us to estimate the magnitude of individual  
307 differences in each sector, i.e., the extent to which there are high/low beluga users in a sector.  
308 Our model found that the CTN sector showed very little individual differences in use (Table  
309 1, i.e., low “sd” value) compared to other sectors, suggesting very little differences in high  
310 and low users of that sector. While the BSM sector showed large individual differences, with  
311 some very high/low users of that sector (Table 1).

312

313 Table 1: Parameter estimates from the multilevel multinomial model predicting the  
314 probability of capturing a photo-ID by sector. Estimated magnitudes of individual differences  
315 (sd) are presented for each sector. Higher ‘sd’ estimates indicate more individual differences  
316 in individual use of that sector, whereas lower estimates indicate individuals are using the  
317 sector at very similar levels. To facilitate interpretation we have ordered the table by lowest  
318 to highest estimates of individual differences. As the number of parameters in the model is  
319 large, the overall mean by sector, and estimated correlations between individual differences,  
320 are presented in the supplementary section (Table S1).

<b>Parameter</b>	<b>Estimate</b>	<b>SD</b>	<b>l-95% CI</b>	<b>u-95% CI</b>
sd(mu_CTN)	0.33	0.03	0.27	0.39
sd(mu_AMO)	0.65	0.38	0.05	1.41
sd(mu_AMN)	0.76	0.29	0.14	1.31

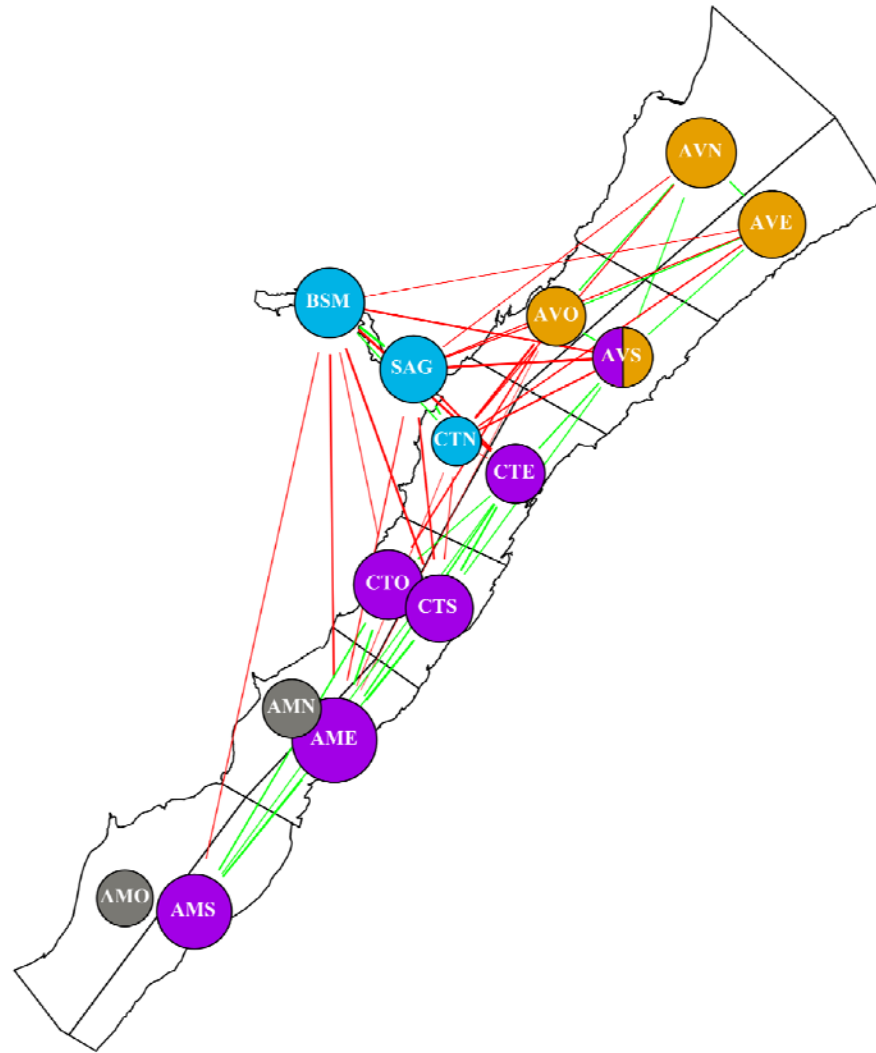


sd(mu_AVO)	0.78	0.05	0.68	0.88
sd(mu_CTE)	0.78	0.05	0.68	0.89
sd(mu_AVS)	0.83	0.07	0.70	0.97
sd(mu_AVE)	1.16	0.20	0.78	1.57
sd(mu_SAG)	1.16	0.08	1.01	1.31
sd(mu_CTS)	1.18	0.11	0.97	1.40
sd(mu_CTO)	1.27	0.14	1.00	1.55
sd(mu_AVN)	1.30	0.17	0.98	1.65
sd(mu_BSM)	1.30	0.11	1.09	1.52
sd(mu_AMS)	1.53	0.22	1.10	1.98
sd(mu_AME)	2.00	0.19	1.64	2.38

---

321

322           Using the between sector correlations to generate a signed network overlaid on top of  
323 the sectors in the St. Lawrence, suggests, for example, that individuals that are seen in the  
324 SAG sector a lot, are also seen in the BSM and CTN sectors a lot, but are seen very little in  
325 the CTE and CTS sectors (Fig. 4). Applying the simple-signed rule and the blockmodeling  
326 algorithm to delineate communities, both find that there are three distinct communities (Fig.  
327 4). Though, in the case of AVS the simple sign-rule suggested multi-membership, while the  
328 blockmodeling algorithm found AVS to be part of the cluster containing (AVO, AVN, AVE)  
329 or that the two clusters (orange and purple in fig. 4) merged into one depending on the choice  
330 of weighting parameter (i.e., emphasizing positive or negative edges).



331

332 Figure 4: Similarity and dissimilarity between sectors in the beluga whale population of the  
333 St. Lawrence. The green edges between two sectors signify that the sectors share high/low  
334 users, while red edges signify that they have dissimilar high/low users. The lack of an edge  
335 signifies that the high/low users of one sector does not provide information about the  
336 high/low users of other sectors. Nodes represent sectors, and are coloured based on shared  
337 communities: i.e., shared green edges, and no red edges. Node sizes represent the magnitudes  
338 of individual differences in use within the sector, i.e., larger nodes suggest larger differences  
339 between high and low users.

340

#### 341 **4. Discussion**

342 Here we've shown that using photo-ID data with a multilevel multinomial model it is  
343 possible to estimate socio-spatial networks, identifying spatial communities while controlling  
344 for sampling biases. Our results suggest that the beluga population shows a non-random  
345 social-spatial structuring within the summer range.

346 Within sectors of the St. Lawrence, our model suggests that subsets of belugas are  
347 heavily using some sectors, while other sectors show little evidence of differences in use. The  
348 magnitude of individual differences in each sector, i.e., how much individuals differ in their  
349 probability of being observed in a particular sector, shows that the CTN sector, in particular,  
350 has very little in the way of individual differences in the probability of being seen in that  
351 sector (Table 1). This result suggests that the CTN sector is used similarly by most  
352 individuals, and represents a potential high mixing zone for the population. In contrast, the  
353 AME sector shows the highest level of individual differences, suggesting that there are large  
354 differences in how beluga are using this sector. These results suggest that the population is  
355 not randomly mixing with the St. Lawrence, and that there are belugas that make use of some  
356 sectors more than other belugas.

357 Between sectors of the St. Lawrence, our results add to the evidence that the beluga  
358 population cannot be assumed to be randomly mixing within its summer habitat. Rather,  
359 comparing the individual differences in beluga usage patterns within sectors suggests  
360 similar/dissimilar user populations across sectors (Fig. 4). By using correlations between  
361 sector usage patterns to create a socio-spatial network, and running community detection  
362 algorithms, our results found that there are spatially distinct communities that make use of  
363 particular regions of the St. Lawrence and the Saguenay River (Fig. 4). We found that the  
364 beluga population in the St. Lawrence could be separated into three distinct communities: 1)

365 The lower St. Lawrence (AVO, AVS, AVE, AVN), the Saguenay River and mouth (BSM,  
366 SAG, CTN), and the upper and eastern portion of the St. Lawrence (CTE, CTS, CTO, AME,  
367 AMS) (Fig. 4).

368 Our findings have direct implications for estimating the impacts of anthropogenic  
369 disturbances on this population. As the population shows evidence of spatially restricted  
370 habitat use, disturbances to particular regions can have a disproportionate impact on  
371 particular segments of the larger population. In particular, the cumulative impacts over time  
372 are likely to be greatly increased in some segments while reduced in other segments of the  
373 population, altering estimations of the distribution of impacts. If cumulative impacts, such as  
374 noise, or environmental contaminants, have a threshold beyond which individual survival is  
375 greatly reduced, properly estimating the distribution of cumulative impacts can have large  
376 implications for conservation management. Our results add to the current understanding of  
377 socio-spatial structuring within this population (Michaud, 1993, 2005), and suggest that more  
378 empirical data, e.g., photo-ID data, movement data, aerial surveillance, should be collected to  
379 better refine socio-spatial mixing in this population.

380 The modeling approach presented in this paper relies on defined sectors within a  
381 particular spatial range, e.g., SAG sector, CTN sector... etc (Fig. 1). In some cases, these  
382 delineations can be justified as they identify management zones, but in other cases, the  
383 delineation and scale of these sectors can be delineated somewhat arbitrarily. Future work  
384 could assess the use of continuous random effects (as opposed to categorical) where  
385 individual differences in the probability of being seen could be on a continuous surface. Point  
386 estimates of individual differences could then be estimated at any location, and correlations  
387 between individual differences obtained between any two points in continuous space. This  
388 approach could avoid reliance on user-defined sectors and facilitate a means of looking at the  
389 results at different scales (e.g., grids of points at various scales could be used when estimating

390 correlations). Similarly, given that the method requires repeated sampling of individuals to  
391 obtain probability of being seen in any one sector, there is a reliance on longitudinal data. In  
392 well sampled populations, it could be feasible to estimate the change in time of individuals  
393 being seen in particular sectors. Here differences in being seen in any particular area could  
394 be explicitly modeled to capture any temporal changes in community substructures, or  
395 developmental trajectories related to habitat use.

396         In terms of implementing the multilevel multinomial model on other photo-ID  
397 datasets, the use of test datasets should hold a prominent role in the analysis. The use of  
398 permutation/randomization methods to both generate structured and unstructured datasets,  
399 while maintaining the sample size distribution of the original datasets, can be very valuable in  
400 helping to set model priors and to interpret the final model results. The use of permutation  
401 approaches is common in social network analysis (Croft et al., 2011; Farine, 2017), and is  
402 becoming more common in statistical workflows more generally (Gelman et al., 2013;  
403 McElreath, 2020).

404

## 405         **5. Conclusions**

406         We have introduced the use of multilevel multinomial modeling to estimate socio-  
407 spatial networks from photo-ID data. We've shown, using testing datasets, that the proposed  
408 method is effective at detecting socio-spatial structures. When applied to 18 years of photo-  
409 ID data from an endangered population of beluga whales in the St. Lawrence, our results  
410 suggest strong evidence that the population has three distinct spatial communities. We  
411 suggest that multilevel multinomial models can be effective in extracting socio-spatial  
412 structuring within animal populations monitored by photo-ID, and can have direct  
413 implications for conservation management.

414

415        **6. Acknowledgment**

416        We would like to thank all those who worked with and supported the GREMM in  
417        collecting these long term data.

418

419        **7. Funding**

420        Funding was provided by Ministère des Forêts, de la Faune et des Parcs du Québec and  
421        Secrétariat à la stratégie maritime du Québec.

422

423        **8. Data accessibility**

424        The permuted and simulated photo ID datasets are available on github  
425        ([github.com/tbonne/photoID\\_multinomial](https://github.com/tbonne/photoID_multinomial)), along with code used in the analysis.

426

427        **9. Author Contributions**

428        RM collected the data; TRB conceived the analytical methodology and performed the  
429        analysis; all authors contributed critically to the drafts and gave final approval for  
430        publication.

431

432

433

434 **References:**

- 435 Chion, C., Lagrois, D., Dupras, J., Turgeon, S., McQuinn, I. H., Michaud, R., Ménard, N., &  
436 Parrott, L. (2017). Underwater acoustic impacts of shipping management measures:  
437 Results from a social-ecological model of boat and whale movements in the St.  
438 Lawrence River Estuary (Canada). *Ecological Modelling*, 354, 72–87.  
439 <https://doi.org/10.1016/j.ecolmodel.2017.03.014>
- 440 Croft, D. P., Madden, J. R., Franks, D. W., & James, R. (2011). Hypothesis testing in animal  
441 social networks. *Trends in Ecology & Evolution*, 26(10), 502–507.
- 442 DeFur, P. L., Evans, G. W., Hubal, E. A. C., Kyle, A. D., Morello-Frosch, R. A., & Williams,  
443 D. R. (2007). Vulnerability as a function of individual and group resources in  
444 cumulative risk assessment. *Environmental Health Perspectives*, 115(5), 817–824.
- 445 Doreian, P., & Mrvar, A. (2015). Structural balance and signed international relations.  
446 *Journal of Social Structure*, 16, 1.
- 447 Evans, P. G., & Hammond, P. S. (2004). Monitoring cetaceans in European waters. *Mammal*  
448 *Review*, 34(1–2), 131–156.
- 449 Farine, D. R. (2017). A guide to null models for animal social network analysis. *Methods in*  
450 *Ecology and Evolution*, 8(10), 1309–1320.
- 451 Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., & Rubin, D. B. (2013).  
452 *Bayesian data analysis*. CRC press.
- 453 Hupman, K., Stockin, K. A., Pollock, K., Pawley, M. D. M., Dwyer, S. L., Lea, C., &  
454 Tezanos-Pinto, G. (2018). Challenges of implementing Mark-recapture studies on  
455 poorly marked gregarious delphinids. *PLOS ONE*, 13(7), e0198167.  
456 <https://doi.org/10.1371/journal.pone.0198167>

- 457 Koivuniemi, M., Auttila, M., Niemi, M., Levänen, R., & Kunnasranta, M. (2016). Photo-ID  
458 as a tool for studying and monitoring the endangered Saimaa ringed seal. *Endangered*  
459 *Species Research*, 30, 29–36.
- 460 Koster, J., & McElreath, R. (2017). Multinomial analysis of behavior: Statistical methods.  
461 *Behavioral Ecology and Sociobiology*, 71(9), 138.
- 462 Lesage, V., McQuinn, I. H., Carrier, D., Gosselin, J.-F., & Mosnier, A. (2014). Exposure of  
463 the beluga (*Delphinapterus leucas*) to marine traffic under various scenarios of transit  
464 route diversion in the St. Lawrence Estuary. *DFO Can. Sci. Advis. Sec., Res. Doc.*  
465 *2013/125. Iv*, 28.
- 466 McElreath, R. (2020). *Statistical rethinking: A Bayesian course with examples in R and Stan*.  
467 CRC press.
- 468 McQuinn, I. H., Lesage, V., Carrier, D., Larrivée, G., Samson, Y., Chartrand, S., Michaud,  
469 R., & Theriault, J. (2011). A threatened beluga (*Delphinapterus leucas*) population in  
470 the traffic lane: Vessel-generated noise characteristics of the Saguenay-St. Lawrence  
471 Marine Park, Canada. *The Journal of the Acoustical Society of America*, 130(6),  
472 3661–3673. <https://doi.org/10.1121/1.3658449>
- 473 Michaud, R. (1993). *Distribution estivale du béluga du Saint-Laurent: Synthèse 1986 à 1992*.  
474 Ministère des pêches et des océans, Direction de la gestion des pêches et de ....
- 475 Michaud, R. (2005). Sociality and ecology of the odontocetes. *Sexual Segregation in*  
476 *Vertebrates*, 303–326.
- 477 Parrott, L., Chion, C., Martins, C. C. A., Lamontagne, P., Turgeon, S., Landry, J. A., Zhens,  
478 B., Marceau, D. J., Michaud, R., Cantin, G., Ménard, N., & Dionne, S. (2011). A  
479 decision support system to assist the sustainable management of navigation activities  
480 in the St. Lawrence River Estuary, Canada. *Environmental Modelling & Software*,  
481 26(12), 1403–1418. <https://doi.org/10.1016/j.envsoft.2011.08.009>



482 Perryman, R. J. Y., Venables, S. K., Tapilatu, R. F., Marshall, A. D., Brown, C., & Franks, D.  
483 W. (2019). Social preferences and network structure in a population of reef manta  
484 rays. *Behavioral Ecology and Sociobiology*, 73(8), 114.  
485 <https://doi.org/10.1007/s00265-019-2720-x>

486 Schilds, A., Mourier, J., Huveneers, C., Nazimi, L., Fox, A., & Leu, S. T. (2019). Evidence  
487 for non-random co-occurrences in a white shark aggregation. *Behavioral Ecology and*  
488 *Sociobiology*, 73(10), 138. <https://doi.org/10.1007/s00265-019-2745-1>

489 Schneider, S., Taylor, G. W., Linquist, S., & Kremer, S. C. (2019). Past, present and future  
490 approaches using computer vision for animal re-identification from camera trap data.  
491 *Methods in Ecology and Evolution*, 10(4), 461–470. [https://doi.org/10.1111/2041-](https://doi.org/10.1111/2041-210X.13133)  
492 [210X.13133](https://doi.org/10.1111/2041-210X.13133)

493  
494  
495  
496

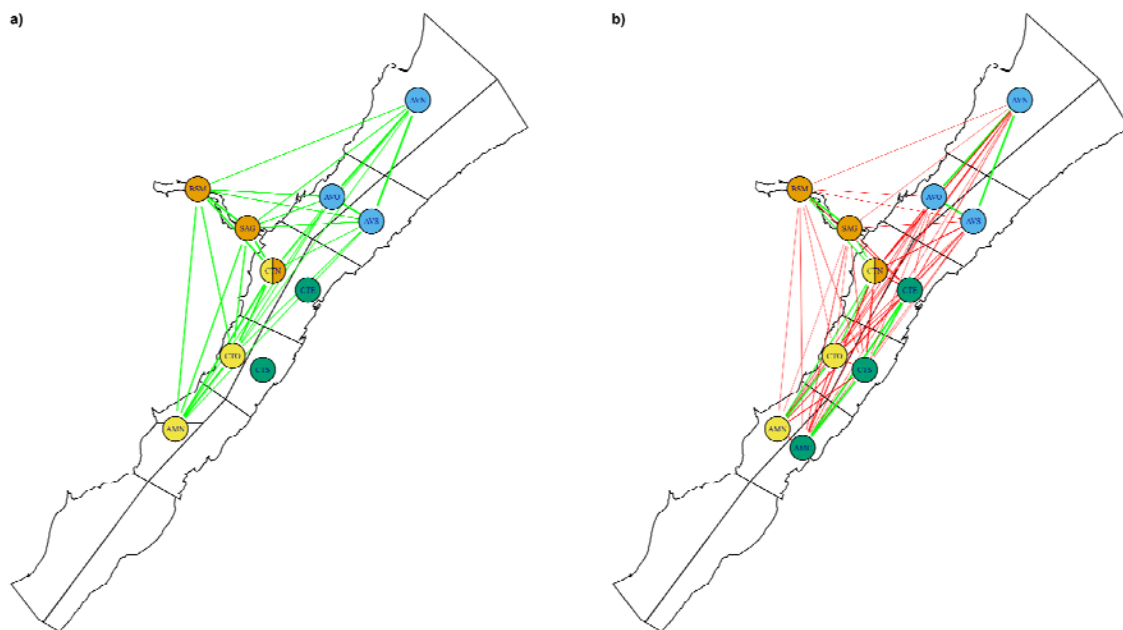
497 **Supplementary material:**

498

499 *Use of a preset reference sector:*

500 We use the simulated dataset with structure, presented in the main text, to fit a  
501 multilevel multinomial model with and without a preset reference sector. That is, in the  
502 model with the preset reference sector we duplicated each individual's photo-IDs and placed  
503 them in the reference sector. This results in a probability of 0.5 for being seen in the reference  
504 sector for all individuals. The results suggest that using a preset reference sector where all  
505 individuals have the same probability of being seen is required to estimate correlations  
506 between sectors and produces appropriate socio-spatial networks (Fig S1).

507



508

509 Fig S1: The estimate socio-spatial networks for a multilevel multinomial model a) without a  
510 preset reference sector, i.e., AME is used as the reference, and b) with a preset reference  
511 sector. The green edges between two sectors signify that the sectors share high/low users,  
512 while red edges signify that they have dissimilar high/low users. The lack of an edge signifies  
513 that the high/low users of one sector does not provide information about the high/low users of

514 other sectors. Nodes represent sectors, and are coloured based on the communities imposed  
515 when simulating the data.

516

517 *Full model table:*

518

519 Table S1: Parameter estimates from the multilevel multinomial model predicting probability  
520 of capturing a photo-ID by sector. Estimated mean probability (logit scale), individual  
521 differences in mean probability ('sd'), and correlation of individual differences between the  
522 different sectors are presented along with an estimate of their 95% credible intervals. Positive  
523 correlations suggest that the high/low users in one sector are similarly high/low users in  
524 another sector, while negative correlations suggest high/low users in one sector are the  
525 low/high users in another sector.

Type	Parameter	Estimate	SD	l-95% CI	u-95% CI
<b>Mean probability (logit scale)</b>					
	mu_AME	-5.85	0.28	-6.43	-5.33
	mu_AMN	-5.60	0.24	-6.15	-5.17
	mu_AMO	-6.68	0.34	-7.41	-6.08
	mu_AMS	-6.33	0.33	-7.00	-5.73
	mu_AVE	-5.47	0.24	-5.99	-5.02
	mu_AVN	-5.06	0.22	-5.52	-4.65
	mu_AVO	-1.82	0.05	-1.93	-1.72
	mu_AVS	-2.90	0.07	-3.04	-2.76
	mu_BSM	-3.78	0.13	-4.04	-3.52
	mu_CTE	-2.20	0.05	-2.30	-2.10
	mu_CTN	-1.24	0.03	-1.30	-1.18
	mu_CTO	-4.58	0.16	-4.92	-4.28

mu_CTS	-4.09	0.13	-4.35	-3.85
mu_SAG	-2.78	0.09	-2.96	-2.61

### Individual Differences

sd(mu_AME)	2.00	0.19	1.64	2.38
sd(mu_AMN)	0.76	0.29	0.14	1.31
sd(mu_AMO)	0.65	0.38	0.05	1.41
sd(mu_AMS)	1.53	0.22	1.10	1.98
sd(mu_AVE)	1.16	0.20	0.78	1.57
sd(mu_AVN)	1.30	0.17	0.98	1.65
sd(mu_AVO)	0.78	0.05	0.68	0.88
sd(mu_AVS)	0.83	0.07	0.70	0.97
sd(mu_BSM)	1.30	0.11	1.09	1.52
sd(mu_CTE)	0.78	0.05	0.68	0.89
sd(mu_CTN)	0.33	0.03	0.27	0.39
sd(mu_CTO)	1.27	0.14	1.00	1.55
sd(mu_CTS)	1.18	0.11	0.97	1.40
sd(mu_SAG)	1.16	0.08	1.01	1.31

### Correlations between individual differences

cor(mu_AME,mu_AMN)	0.39	0.19	-0.03	0.71
cor(mu_AME,mu_AMO)	0.27	0.24	-0.26	0.67
cor(mu_AMN,mu_AMO)	0.14	0.24	-0.36	0.57
cor(mu_AME,mu_AMS)	0.63	0.11	0.39	0.82
cor(mu_AMN,mu_AMS)	0.25	0.21	-0.19	0.62
cor(mu_AMO,mu_AMS)	0.27	0.24	-0.26	0.67
cor(mu_AME,mu_AVE)	0.18	0.16	-0.14	0.46
cor(mu_AMN,mu_AVE)	0.15	0.21	-0.28	0.56

cor(mu_AMO,mu_AVE)	0.10	0.23	-0.36	0.52
cor(mu_AMS,mu_AVE)	0.05	0.19	-0.32	0.41
cor(mu_AME,mu_AVN)	-0.04	0.15	-0.33	0.25
cor(mu_AMN,mu_AVN)	-0.02	0.21	-0.42	0.39
cor(mu_AMO,mu_AVN)	0.00	0.22	-0.42	0.44
cor(mu_AMS,mu_AVN)	-0.07	0.18	-0.41	0.28
cor(mu_AVE,mu_AVN)	0.50	0.14	0.19	0.75
cor(mu_AME,mu_AVO)	-0.21	0.10	-0.40	0.00
cor(mu_AMN,mu_AVO)	0.03	0.19	-0.35	0.40
cor(mu_AMO,mu_AVO)	-0.06	0.21	-0.45	0.36
cor(mu_AMS,mu_AVO)	-0.28	0.14	-0.54	0.01
cor(mu_AVE,mu_AVO)	0.48	0.13	0.22	0.71
cor(mu_AVN,mu_AVO)	0.50	0.10	0.28	0.69
cor(mu_AME,mu_AVS)	0.13	0.12	-0.10	0.36
cor(mu_AMN,mu_AVS)	0.21	0.18	-0.16	0.53
cor(mu_AMO,mu_AVS)	0.13	0.21	-0.31	0.50
cor(mu_AMS,mu_AVS)	0.03	0.15	-0.26	0.32
cor(mu_AVE,mu_AVS)	0.40	0.14	0.12	0.66
cor(mu_AVN,mu_AVS)	0.38	0.12	0.14	0.60
cor(mu_AVO,mu_AVS)	0.54	0.08	0.38	0.69
cor(mu_AME,mu_BSM)	-0.57	0.09	-0.75	-0.38
cor(mu_AMN,mu_BSM)	-0.37	0.17	-0.66	0.00
cor(mu_AMO,mu_BSM)	-0.22	0.23	-0.62	0.26
cor(mu_AMS,mu_BSM)	-0.37	0.13	-0.62	-0.10
cor(mu_AVE,mu_BSM)	-0.34	0.13	-0.59	-0.06
cor(mu_AVN,mu_BSM)	-0.22	0.12	-0.44	0.02

cor(mu_AVO,mu_BSM)	-0.16	0.08	-0.32	0.01
cor(mu_AVS,mu_BSM)	-0.66	0.07	-0.80	-0.51
cor(mu_AME,mu_CTE)	0.54	0.10	0.34	0.72
cor(mu_AMN,mu_CTE)	0.34	0.18	-0.04	0.65
cor(mu_AMO,mu_CTE)	0.21	0.23	-0.29	0.60
cor(mu_AMS,mu_CTE)	0.38	0.14	0.10	0.62
cor(mu_AVE,mu_CTE)	0.07	0.14	-0.21	0.36
cor(mu_AVN,mu_CTE)	-0.11	0.12	-0.35	0.14
cor(mu_AVO,mu_CTE)	-0.10	0.08	-0.26	0.07
cor(mu_AVS,mu_CTE)	0.52	0.08	0.34	0.67
cor(mu_BSM,mu_CTE)	-0.80	0.06	-0.89	-0.68
cor(mu_AME,mu_CTN)	-0.26	0.12	-0.49	-0.03
cor(mu_AMN,mu_CTN)	-0.24	0.19	-0.59	0.16
cor(mu_AMO,mu_CTN)	-0.12	0.21	-0.50	0.32
cor(mu_AMS,mu_CTN)	-0.09	0.15	-0.39	0.21
cor(mu_AVE,mu_CTN)	-0.53	0.13	-0.77	-0.26
cor(mu_AVN,mu_CTN)	-0.44	0.12	-0.67	-0.18
cor(mu_AVO,mu_CTN)	-0.75	0.07	-0.86	-0.60
cor(mu_AVS,mu_CTN)	-0.63	0.09	-0.79	-0.44
cor(mu_BSM,mu_CTN)	0.47	0.09	0.28	0.64
cor(mu_CTE,mu_CTN)	-0.25	0.10	-0.46	-0.05
cor(mu_AME,mu_CTO)	0.63	0.09	0.43	0.79
cor(mu_AMN,mu_CTO)	0.28	0.20	-0.15	0.63
cor(mu_AMO,mu_CTO)	0.16	0.22	-0.30	0.56
cor(mu_AMS,mu_CTO)	0.55	0.13	0.27	0.78
cor(mu_AVE,mu_CTO)	-0.12	0.17	-0.45	0.22

cor(mu_AVN,mu_CTO)	-0.24	0.15	-0.54	0.07
cor(mu_AVO,mu_CTO)	-0.55	0.09	-0.72	-0.36
cor(mu_AVS,mu_CTO)	-0.16	0.12	-0.39	0.07
cor(mu_BSM,mu_CTO)	-0.33	0.11	-0.53	-0.12
cor(mu_CTE,mu_CTO)	0.39	0.11	0.17	0.60
cor(mu_CTN,mu_CTO)	0.19	0.12	-0.06	0.42
cor(mu_AME,mu_CTS)	0.70	0.08	0.53	0.85
cor(mu_AMN,mu_CTS)	0.42	0.19	0.00	0.73
cor(mu_AMO,mu_CTS)	0.21	0.23	-0.27	0.61
cor(mu_AMS,mu_CTS)	0.42	0.14	0.14	0.67
cor(mu_AVE,mu_CTS)	0.16	0.16	-0.15	0.45
cor(mu_AVN,mu_CTS)	0.06	0.14	-0.23	0.33
cor(mu_AVO,mu_CTS)	-0.04	0.10	-0.23	0.15
cor(mu_AVS,mu_CTS)	0.42	0.11	0.20	0.61
cor(mu_BSM,mu_CTS)	-0.72	0.08	-0.85	-0.55
cor(mu_CTE,mu_CTS)	0.64	0.09	0.46	0.80
cor(mu_CTN,mu_CTS)	-0.37	0.11	-0.57	-0.14
cor(mu_CTO,mu_CTS)	0.48	0.11	0.25	0.68
cor(mu_AME,mu_SAG)	-0.43	0.09	-0.61	-0.24
cor(mu_AMN,mu_SAG)	-0.36	0.17	-0.66	0.00
cor(mu_AMO,mu_SAG)	-0.17	0.22	-0.56	0.29
cor(mu_AMS,mu_SAG)	-0.22	0.13	-0.47	0.03
cor(mu_AVE,mu_SAG)	-0.44	0.12	-0.68	-0.19
cor(mu_AVN,mu_SAG)	-0.37	0.10	-0.57	-0.16
cor(mu_AVO,mu_SAG)	-0.48	0.06	-0.60	-0.35
cor(mu_AVS,mu_SAG)	-0.80	0.06	-0.90	-0.68

cor(mu_BSM,mu_SAG)	0.83	0.04	0.74	0.91
cor(mu_CTE,mu_SAG)	-0.68	0.06	-0.79	-0.55
cor(mu_CTN,mu_SAG)	0.67	0.07	0.52	0.80
cor(mu_CTO,mu_SAG)	-0.12	0.11	-0.32	0.09
cor(mu_CTS,mu_SAG)	-0.62	0.08	-0.77	-0.45

---

526

527

528

529 *Note on interpreting low sd within sectors:*

530 Similar to the CTN sector, the sector AMO is also estimated to have a low magnitude of  
531 individual differences. However, it has a large uncertainty in this estimate. This highlights  
532 that it is possible to have little individual differences in a sector due to either: 1) limited data,  
533 resulting in all individuals being pooled to the mean value, and 2) limited data is not a factor,  
534 but individuals are using this sector relatively equally. Care should therefore be taken when  
535 interpreting the magnitude of individual differences, nevertheless, the estimated uncertainty  
536 around magnitude estimates is one way to help identify sectors with limited data.