1 **Quantifying song categories in Adelaide's Warbler (*Setophaga adelaidae*)**

2

3 **Chinthaka D. Kaluthota[1,4], Orlando J. Medina[2], David M. Logue[1, 3]**

4 [1] Behaviour and Evolution Research Group, Department of Psychology, University of Lethbridge. 4401 University

5 Drive West, Lethbridge, AB T1K3M4, Canada. (ORCID: 0000-0001-8151-4289)

6 [2] Cabo Rojo National Wildlife Refuge. United States Fish and Wildlife Service. PO Box 510, Boqueron, Cabo Rojo,

7 PR 00622-0510, USA

8 [3] Departamento de Biología, Universidad de Puerto Rico. Mayagüez PR, 00682, USA

9 [4] Corresponding author email: kaluthota@uleth.ca, Phone: +1 (403) 715-8027

10

11 **Abstract**

12 Many migratory wood-warblers in the genus *Setophaga* divide their song repertoires into two categories. Category B

13 songs are usually sung before dawn, with immediate variety and short latencies between songs, whereas category A

14 songs are sung exclusively after dawn, with eventual variety and longer latencies between songs. Songs in different

15 categories may also differ with respect to their acoustic structure. We used an unsupervised clustering algorithm to

16 identify song categories in Adelaide's Warbler (*Setophaga adelaidae*), a year-round territorial species. We identified

17 two categories of song types, the characteristics of which are similar to song categories in other migratory wood-

18 warblers. Clusters were not well-separated, suggesting that song categories may not be discrete. Song structures in

19 the two categories were similar, but category B songs were shorter and had fewer notes than did category A songs.

20 On average, dyads of males shared more category B songs than category A songs, and were more likely to use

21 category B songs when song type matching other males. The most important song delivery variable for separating

22 clusters was residual average run length (residual values control for covariation with time of day), followed by

23 percent of songs delivered before dawn, residual latency, and percent of songs used as song-type matches. We

24 recommend a scheme based on the first three variables to classify novel song types.

25

26 Key words: Adelaide's Warbler, Neotropical birds, singing modes, song repertoires, song types, year-round

27 territoriality

28

41    In many songbird species, individuals can produce more than one song type. The set of song types that an individual

42    can sing is called a repertoire. Song type repertoires are hypothesised to function in both mate attraction and in

43    competition with rival conspecifics (Catchpole & Slater, 2008). Some species divide their song repertoires into

44    distinct categories (Illes, 2015; Molles & Vehrencamp, 1999; Spector, 1992). This pattern is especially common in

45    wood-warblers belonging to the genus *Setophaga,* many of which use two song categories (Family: Parulidae;

46    Bolsinger, 2000; Demko et al., 2013; Price & Crawford, 2013; Spector, 1992; Wiley et al., 1994). The two song

47    categories may be characterized by distinct song delivery patterns, acoustic structures, and / or behavioral contexts

48    (Spector, 1992). Throughout this report, we follow Spector's (1992) recommendation to refer to these categories as

49    "first category" and "second category" when referencing multiple species, and to use pre-existing terminology when

50    discussing individual species (e.g., "category A" and "category B," in the case of our focal species).  Many, but not

51    all, *Setophaga* wood-warblers deliver first category songs in a repetitive manner, with relatively long silent gaps

52    between songs, during the daytime, and in the context of male-female interactions (Spector, 1992). By comparison,

53    second category songs are typically delivered with immediate variety, at higher singing rates, mostly around dawn

54    and in the context of male-male interactions (Beebee, 2004b; Kroodsma et al., 1989; Price & Crawford, 2013;

55    Spector, 1992; Staicer, 1989, 1991, 1996a).

56        In addition to differences in delivery patterns and behavioral context, songs from the two categories can be

57    structurally distinct. In American Yellow Warblers (*Setophaga petechia*), for example, type I (first category) songs

58    have higher trill rates and wider frequency bandwidths than type II (second category) songs (Beebee, 2004b).

59    Similarly, Price and Crawford (2013) showed that Pine Warblers' (*Setophaga pinus*) first category songs have

60    significantly higher trill rates than their second category songs.

61        The existence of song categories in wood-warblers is a foundational example of functional variation in

62    song type repertoires (Catchpole & Slater, 2008; Macdougall-Shackleton, 1997). Most published research on song

63    categories in wood-warblers, however, has come from North-temperate, migratory species. These studies may not

64    represent the many wood-warbler species that breed at other latitudes because the vocal behavior of tropical and

65    southern-hemisphere songbirds is often different from that of North-temperate species (Curson et al., 1984;

66    Stutchbury & Morton, 2001). In the present study, we evaluate the evidence for two distinct song categories in a

67    year-round territorial, tropical wood-warbler.

68    Adelaide's Warbler (*Setophaga adelaidae*) is endemic to Puerto Rico and the nearby island of Vieques

69    (Toms, 2010). Socially monogamous pairs defend all-purpose territories throughout the year. Previous studies of this

70    species conclude that males have two distinct song categories (Staicer, 1991, 1996a, 1996b). Those studies indicate

71    that males deliver category A (first category) songs after dawn, throughout the year, with eventual variety, and

72    relatively long intervals between songs. Category B (second category) songs occur primarily before dawn, during

73    the breeding season, with immediate variety, and relatively short intervals between songs. Category B songs are also

74    reported to occur later in the day, when they are still delivered with immediate variety and short intervals (although

75    intervals are longer and variety is lower than during pre-dawn singing). Males tend to use Category A songs when

76    interacting with females and Category B songs when interacting with other males.

77    Although the categorization of songs into two distinct modes is easily achieved for some wood-warbler

78    species, Staicer (1996b) found that Adelaide's Warbler songs could not be easily categorized because the acoustic

79    properties of songs in the two categories overlap. Song categorization is further complicated by the fact that

80    different males can assign a given song type to different categories (throughout the manuscript we use the term

81    "assign" to indicate that a bird delivers a song in a way that is characteristic of a given category). For example, one

82    male could assign song X to category A, but a different male could assign song X to category B. This type of

83    individual-specific usage of song types has also been observed in some other wood-warbler species (Beebee, 2004b;

84    Lemon et al., 1985; Price & Crawford, 2013; Spector, 1991; Staicer, 1989).

85    Staicer (1991, 1996a) used song delivery patterns to categorize songs. In one study, songs sung before

86    dawn (hereafter, "dawn" after Staicer, 1991) and after dawn ("morning") were taken to represent category B and A

87    songs, respectively (Staicer, 1996a). In another study, Staicer (1991, p. 21) describes a slightly different

88    categorization scheme: "Songs that males sang in the dawn chorus were called B songs and those that males

89    switched to around sunrise were called A songs…. If I recorded a song type only after sunrise but at a rate typical of

90    B songs, and usually in the same bout (consecutive song sequence) with B songs, I classified it as a B song" (a

91    similar description can be found in Staicer, 1996b).

92    Staicer's categorization schemes succeeded in elucidating important patterns in Adelaide's Warblers' vocal

93    communication system (Staicer, 1996a, 1996b). Nevertheless, we see room for improvement. The first approach

94    described in the previous paragraph has the potential to miscategorise category B songs, which are reported to occur

95    both before and after dawn. The second approach improves on the first, but relies on the analyst to determine typical

96  rates of category B songs and whether a given song is sufficiently linked to other category B songs. We had several

97  motivations to develop a different approach to song categorization. First, we wanted to test whether Adelaide's

98  Warbler repertoires are best divided into two categories, as opposed to one category, three categories, etc. Second,

99  we were interested in which song delivery variables were most useful for categorizing songs. Third, we wanted to

100  develop an algorithm to categorize novel songs (see Chapter 3 of Hastie & Dawes 2010 for a discussion of the

101  benefits of automatic judgement algorithms). Finally, we wanted to compare our results to Staicer's (1991, 1996a, b)

102  results because replication with novel data and a novel analytic approach is a critical, but underutilized, component

103  of the scientific method (Baker, 2016).

104      The present study examines song categories in male Adelaide's Warbler. We applied an unsupervised

105  clustering algorithm to song type repertoires. Specifically, we ran cluster analyses for all fifteen combinations of

106  four song delivery variables believed to be important for distinguishing song categories. We chose the best number

107  of clusters (k) from each analysis, and then the best analyses based on a given number of variables. We evaluated

108  the four remaining clustering schemes with respect to the number of song categories, their characteristics, their

109  discreteness, and their distributions across individuals. Finally, we looked at patterns of song sharing and compared

110  song structures across categories.

111

112  **Methods**

113  **Field data collection**

114  We recorded nine mated male Adelaide's Warblers at the Cabo Rojo National Wildlife Refuge in south-western

115  Puerto Rico (17º 59' N, 67º 100' W) during the breeding season between March and June, 2012. Males were

116  captured using mist nets and marked with three colored leg bands and a numbered metal band prior to recording.

117      Individual males were recorded continuously for approximately 3.5 hours per day, for four days each.

118  Consecutive recordings of a given male were separated by at least four days, except on two occasions when

119  recordings were made on consecutive days because of logistical constraints. Observations started 45 minutes before

120  sunrise to ensure that we recorded the first song of the day. Recordists announced the singer's identity after each

121  song, as well as song type matches and fights with neighboring conspecifics. The identity of focal males was

122  confirmed by inspecting the colored leg bands prior to the end of each recording session. Recordings were made

123  with Marantz PMD 661 digital recorders and Sennheisser ME67 shotgun microphones (file format = wav, sampling

124  rate = 44.1 kHz, bit depth = 16 bits). We sampled a large number of songs from a small number of individuals. The

125  decision to emphasize recording effort per individual resulted in increased accuracy of singing variables at the song

126  type within individual level, but the small sample of individuals may have limited our ability to fully capture the

127  range of variation among individuals. This is the same set of recordings used in Schraft et al. (2017) and  Hedley et

128  al. (2018).

129

130  **Scoring and acoustic analysis**

131  We inspected sound spectrograms of all recordings in Syrinx PC v2.6f sound analysis software (John Burt,

132  http://www.Syrinxpc.com; Blackman window; window size = 1024 points). We assigned songs to song types based

133  on their appearance on a spectrogram. Several observers assigned songs to types as they scored the field recordings

134  for each male. Then, one person (D.M.L.) scored song types across males (i.e., decided which songs from different

135  males' repertoires belonged to the same type), and corrected scoring errors. We measured the inter-rater reliability

136  of song type scoring within an individual bird. We randomly selected 100 songs from one individual, and two

137  experienced bio-acousticians used a classification key to score them independently. The result was 100% agreement.

138  This analysis demonstrates that different observers agree on how to score song types within males. Most of the

139  analyses in the present study use "song type within male" as the independent sampling unit, so it is important to

140  verify that song type scoring within male is repeatable across observers. We then estimated the repeatability of

141  scoring song types among individual birds by having an experienced bio-acoustician (C.D.K.) re-classify 22-23

142  randomly selected songs from each of the nine males (total = 200 songs) using the population-level classification

143  key. In this analysis, 175 of 200 (87%) scores matched the originals. This number is lower than the within-individual

144  repeatability because song structure can vary among individuals. Among-individual repeatability is relevant to the

145  song type sharing analyses in the present study.

146  Song recordings with high signal-to-noise ratios, as determined by visual inspection of spectrograms, were

147  subjected to detailed acoustic analysis. We used Luscinia v.2.14 (Lachlan, 2007) sound analysis software to obtain

148  the minimum and maximum peak frequency, number of notes, and song duration for high signal-to-noise-ratio song

149  recordings from the breeding season dataset (settings: maximum frequency = 10 kHz, frame length [equivalent to

150  "bin" or "FFT" length] = 5 ms, time step = 1 ms, dynamic range = 35 dB, dynamic equalization = 100 ms, de-

151  reverberation = 100%, de-reverberation range = 100 ms, high-pass threshold = 1.0 kHz, noise removal = 10 dB). In

152  Luscinia, users identify focal signals by outlining their images on a sound spectrogram. Sounds in the outlined area

153  that exceed user-defined thresholds for amplitude and duration are labeled with a colored trace that users can

154  compare to the spectrogram to correct errors. The program stores acoustic information about the signal (as defined

155  by the trace) in a database, which users can query (e.g., for the minimum peak frequency for each note). Minimum

156  frequency was subtracted from maximum frequency to calculate the frequency bandwidth of each note. The

157  frequency bandwidth of each song was defined as the average frequency bandwidth of its notes. Trill rate was

158  calculated as note number / song duration (sec). We measured the frequency excursion (FEX) of the same songs

159  using the program FEX Calculator (Podos et al., 2016, J. McClure, https://github.com/BehaviorEnterprises/Fex).

160  FEX is a putative metric of vocal performance that measures changes in fundamental frequency, including changes

161  that are not voiced, over time.

162  We performed a series of unsupervised cluster analyses to categorize songs. We treated "song type within

163  individual" as the independent sampling unit, because different males may assign individual song types to different

164  categories (Staicer, 1991). There was a high risk of misclassifying rare songs types, so we omitted song types within

165  individual that were recorded 10 times or fewer (10 was chosen as a cut-off point because it is a round number).

166  Clustering was based on variables linked to four of the five attributes hypothesized to separate song categories

167  (Table 1): mean residual latency (a measure of the interval between songs), mean residual run length (a proxy for

168  delivery mode), percent matching (a proxy for social context), and percent of songs sung during the pre-dawn

169  period.

170  Latency was defined as the time since the focal male's prior song. Latency co-varies with time of day

171  (Staicer 1991). To better isolate latency from time of day, we regressed latency against time and used the residual

172  latency to generate means (logistic regression: $r^2 = 0.387$, $F_{1, 8998} = 5684$, $p < 0.001$, constant = 0.76, beta = 1.0). Run

173  length was the number of songs in a continuous run of a given song type (we only used a single run length value for

174  each run to avoid pseudoreplication). Like latency, run length was correlated with time of day, so we regressed run

175  length against time and used the residual latency to generate means (logistic regression: $r^2 = 0.159$, $F_{1, 5933} = 1124$, p

176  $< 0.001$, constant = 0.84, beta = 1.0).

177  We chose song type matches by the focal males as our measure of vocal interaction, because it was less

178  likely to occur by chance than other kinds of vocal interaction (e.g., song overlapping, unmatched counter-singing),

179  and so more likely to represent a deliberate interaction on the part of the focal male. Song type matching was scored

7

180  when a focal male sang the same song type as a neighbor had sung in the previous two seconds. Two seconds was

181  chosen as the cut-off for song type matching because it corresponds to the average duration of a male song (mean =

182  2.0 sec., n = 2776). Matching was scored if the recordist dictated "song type match" or if examination of sound

183  spectrograms revealed a match.

184       To calculate the percent of songs delivered before dawn, we first obtained sunrise time from the website

185  *www.timeanddate.com* and scored "time to sunrise" (in seconds) for each song. Negative values of this variable

186  correspond to times before sunrise, and positive values correspond to times after sunrise. We defined the end of the

187  dawn chorus as the time that song rates stabilized after the period of intense pre-dawn singing. Two of us (C.D. K.

188  and D.M.L.) visually inspected a histogram of song delivery times, and determined that 700 seconds after sunrise

189  was the optimal cut-off time. Following (Staicer, 1991), we refer to the period before this cut-off as "dawn," and the

190  period after the cut-off as "morning".  We controlled for sampling effort by dividing the proportion of dawn songs

191  from the focal male that belong to the focal song type by the proportion of all songs from the focal male that belong

192  to the focal song type, and multiplying by 100. All data were standardized by subtracting the mean and dividing by

193  the standard deviation prior to cluster analysis.

194       We ran the *NbClust* function in the R package *NbClust*  to cluster the data (Charrad et al., 2014). We chose

195  the k-means analysis with Euclidean distance because it is a simple and effective approach to unsupervised

196  clustering. We ran separate analyses for each of the 15 unique combinations of the four clustering variables (Table

197  2). For each analysis, we interpreted the number of clusters that produced the highest average silhouette index. The

198  average silhouette index measures the average similarity of objects to other objects in their own clusters relative to

199  objects in other clusters (Rousseeuw, 1987). The index ranges from -1 to +1, where higher values indicate better

200  clustering. We used the silhouette index to choose the best categorization schemes because it is a simple, widely-

201  used metric that reflects our intuitive concept of clustering. At this point, we had 15 clustering schemes (one for

202  each unique combination of variables). We then identified the clustering scheme with the highest average silhouette

203  index for each number of clustering variables, one through four. We were curious whether clustering strength was

204  influenced by sample size, so we calculated separate average silhouette indices for well-sampled (n > 40 songs) song

205  types and poorly-sampled (n ≤ 40 songs) song types for each of these four clustering schemes.

206       We compared the assignment of song types to categories across males to test the hypothesis that different

207  individuals assign the same song type to different categories. We then tested whether song type sharing differed

8

208  between song categories. First, we identified song types that were shared and assigned to the same category for all

209  possible pairs of males. We then conducted a paired t-test in Microsoft Excel 2013 (Microsoft Corp., Redmond,

210  WA). This test treated dyads of males as the independent sampling units. It tested the null hypothesis that the mean

211  difference between the number of shared category A songs and shared category B songs shared was zero.

212        To test whether there were structural differences between song categories, we developed mixed models

213  with the following structural properties as dependent variables: frequency excursion (FEX, unitless), song duration

214  (ms), number of notes, trill rate (notes / sec), minimum frequency (kHz), maximum frequency (kHz) and average

215  frequency bandwidth (kHz; averaged over the notes in the song). The sampling units were individual song

216  utterances, and the dataset was restricted to the high-quality song recordings that were suitable for fine-scale

217  structural analysis. For the structural models, song category was included as a fixed factor, and individual, recording

218  day nested within individual, and song type within individual (not a statistically nested variable) were included as

219  random factors. We chose to include song type in this model because structural variables (but not song delivery

220  variables) are inextricably linked to song type, and therefore different utterances of a given song type are not

221  independent with respect to acoustic structure. Mixed models were constructed with the *lmer* function of the

222  package *lme4* (Bates et al., 2015). We examined residual plots to evaluate model fit. All mixed models are in the

223  online resource.

224

225  **Results**

226  The breeding season recordings of nine focal males included 9,420 songs comprising 71 song types. Summing all

227  males' repertoires, there were 261 song types within individual (average ± SD = 29.0 ± 4.0 song types per male).

228  Removing uncommon song types (song types recorded 10 times or less within individual) resulted in 9032 song

229  utterances comprising 57 song types and 168 song types within individual (18.7 ± 2.4 song types per male). Of

230  those, 2776 song recordings were of sufficient quality for structural analysis.

231        The best clustering schemes for all fifteen combinations of song delivery variables are described in Table 2.

232  Most of these schemes included two (seven schemes) or three (six schemes) clusters, but one scheme used five, and

233  another used seven. The average silhouette index was negatively correlated to the number of variables, as is typical

234  of real world clustering problems (n = 15, r = -0.90; Anzanello & Fogliatto 2011).

9

235    Of the four clustering schemes with the highest average silhouette indices for a given number of variables,

236    three were based on two clusters, and one (the best two-variable scheme) was based on three clusters. All four

237    schemes included one category characterized by high (> 50%) average percent predawn, low (< 50 s) average

238    latency, short (< 2.5 songs) average run lengths, and high (> 4.9%) average percent song matching (Table 3). We

239    refer to these clusters as "category B" because of their similarity to previous descriptions of category B (Table 1).

240    We refer to the other clusters as "category A" because of their similarity to previous descriptions of category A. The

241    best two-variable scheme includes two clusters that are like category A.  We refer to them as categories A1 and A2.

242    There was considerable variation among clustering schemes in the number of song types assigned to each

243    category for each male (Fig. 1). The best one-variable clustering scheme put most song types within individual

244    (86.3%) in category B. One male's (RDY's) songs were all assigned to category B. The best two-variable scheme

245    assigned 69.0% of songs to category B. Two males were not assigned any A1 songs, and one was not assigned any

246    A2 songs. The best three variable scheme, based on residual run length, residual latency, and percent pre-dawn,

247    assigned 64.3% of songs to category B. All males were assigned A and B songs. Similarly, the best four variable

248    scheme assigned 61.9% of songs to category B and assigned all males both A and B songs.

249    Cluster plots for the four best clustering schemes all included a relatively dense cluster that corresponds to

250    category B (Fig. 2). The points outside of that cluster were more diffuse. None of the scatterplots showed a clear gap

251    between the category B cluster and the other points. Average silhouette indices were significantly higher for well-

252    sampled song types than they were for poorly sampled song types in three of the four best clustering schemes (Table

253    SX). Separate plots for each male can be found in the electronic supplementary material (Fig. S1).

254    No song types were recorded from all nine subjects. We identified 41 (71.9%) song types that were shared

255    by at least two birds and 14 (24.6%) song types that were shared by at least five birds. Assignment of song types to

256    categories varied among individuals (Table S1). For example, in both the three-variable and four-variable schemes,

257    20 (48.8%) song types were used in different categories by different birds. Focusing on song types that are both

258    shared and assigned to the same category, dyads shared more category B songs than category A for all four schemes

259    (Table S1). For example, in both the three-variable and four variable schemes, dyads shared an average of 3.31 (±

260    2.33) category B songs, but only 1.25 (± 2.33) category A songs.

261    We studied the clustering results to rank variables by their importance for clustering. The best clustering

262    schemes for a given number of variables formed a nested hierarchy (Table 2). Residual run length was present in all

263    four schemes, residual latency was in three, percent pre-dawn was in two, and percent matching was in one. Among

264    the one-variable clustering schemes, the scheme based on residual run length produced the highest average

265    silhouette index (0.771), followed by schemes based on percent pre-dawn (0.738), residual latency (0.730), and

266    percent matching (0.688). In the four variable scheme, the mean difference in standardized values between clusters

267    was largest for percent pre-dawn (1.68), followed by residual latency (1.44), residual run length (1.04), and percent

268    matching (0.46).

269         Relative to category A songs, category B songs were shorter, with fewer notes and slightly lower trill rates

270    according to the three-variable clustering scheme (Tables 4, see Tables S2 & S3 for detailed results from the four

271    best clustering schemes). Frequency excursion, maximum frequency, minimum frequency, and frequency bandwidth

272    were not significantly different between song categories. Models based on the best one-, two-, and four-variable

273    clustering schemes produced qualitatively similar results, but the FEX effect was statistically significant in the four-

274    variable scheme, and the trill rate effect was not significant in the one-variable and two-variable schemes.

275

276

277    **Discussion**

278    The four clustering schemes with the highest average silhouette index for a given number of variables (hereafter, the

279    "best" schemes) all included one cluster that matched previous descriptions of category B song and one or two

280    clusters that matched previous descriptions of category A songs (Tables 1 & 3; Staicer 1991). Compared to category

281    A songs, category B songs were more likely to be sung before dawn. Category B songs were also delivered in

282    shorter runs and with shorter intervals between songs (controlling for time of day). Category B songs were more

283    likely than category A songs to be used as song type matches with neighbors, suggesting they may be especially

284    important for male-male vocal interactions. These song delivery patterns are similar to those of several migratory

285    species of wood-warblers (Demko et al., 2013; Price & Crawford, 2013; Spector, 1991, 1992). We found that male

286    Adelaide's Warbler repertoires included fewer A songs than B songs (Fig. 1; Table 3), which is consistent with the

287    results of previous studies on Adelaide's Warbler (Staicer, 1991, 1996a) and other wood-warbler species (Lemon et

288    al., 1985; Spector, 1991; Staicer, 1989). Overall, our results strongly support Staicer's (1991, 1996a)

289    characterization of song categories in Adelaide's Warbler.

290    In the four best clustering schemes, category B songs were clustered in song delivery space, whereas

291    category A songs were more dispersed (Fig. 2). There was no clear break between categories A and B in song

292    delivery space. The average silhouette index from the four-variable scheme (0.379) suggests only weak evidence of

293    clustering, but the indices from the other schemes suggest moderate-to-strong clustering (range = 0.495-0.771).

294    Importantly, the best sampled song types within male (n > 40 songs), were significantly more clustered than the

295    other song types in three of the four best clustering schemes (the sole exception was the one-variable scheme, Table

296    SX). This finding indicates that sampling error depressed the average clustering values in the full dataset, and

297    validates our decision to collect emphasize intensive sample within individuals (as opposed to collecting fewer

298    samples from more individuals). We conclude that there are song categories in this population, but they may not be

299    entirely discrete. Perhaps song types within male exist on a continuum from "B-like" to "A-like," with many songs

300    clustering near the B-like end of the continuum. Larger samples permitted more precise estimates of average song

301    delivery variables, which increased the separation between the song categories. Interestingly, a recent Pine Warbler

302    study concluded that species' song system lacked, "the clear distinction between song categories typical of most

303    other *Setophaga* wood-warblers" (Price & Crawford, 2013, p.559). Those results are based on recordings from the

304    later part of the breeding season, calling into question whether song delivery patterns may be more distinct earlier in

305    the year. Nevertheless, the hypothesis that song categories may not be entirely distinct merits further research.

306    The best two-variable clustering scheme included two "A" categories (Table 3, Figs. 1 & 2). Relative to

307    category A2, category A1 was characterized by lower percent pre-dawn, lower latency, and much longer runs. We

308    are reluctant to endorse a three-category hypothesis for several reasons. First, only one of the four best clustering

309    schemes suggested three categories, while the others all agreed that two categories were best. Second, we know of

310    no evidence of three-category systems in other *Setophaga* species. Third, samples from three of our nine subjects

311    were missing either category A1 songs or category A2 songs. We conclude that certain category A songs may be

312    used in particularly long runs, but the bulk of the evidence suggests these songs should be lumped in with category

313    A songs, rather than being placed in a separate category.

314    As previously described in this species, different individuals assigned specific song types to different

315    categories (Staicer, 1991, 1996b). Any hypothesis to explain the development of function of song categories in

316    Adelaide's warbler must account for this finding. Males shared more category B songs than category A songs, as has

317    previously been shown for Yellow Warblers (Beebee, 2002). This pattern was not found in Pine Warblers (Price &

12

318  Crawford, 2013), but again, that study was conducted much later in the breeding season than was ours. The findings

319  that males shared more category B songs, and that B songs were more likely to be used as song type matches, lend

320  credence to the hypothesis that category B songs are especially important in male-male vocal interactions.

321  The four best clustering schemes all indicated that relative to category A songs, category B songs were

322  shorter and contained fewer notes. Some clustering schemes indicated statistically significant effects of song

323  category on FEX and trill rate, but the effect sizes were small. Contrary to our results, Staicer (1996a) found no

324  difference in song duration between categories. Otherwise, Staicer's (1996a) song structure results agree with our

325  own: no differences in trill rate, minimum frequency, maximum frequency, or frequency bandwidth. Staicer (1996a)

326  did not compare note number or frequency excursion between song categories, but she did find differences in

327  syllable complexity (category B songs were more complex) and frequency spectra (category B songs had lower peak

328  frequencies) that we did not test for.

329  Structural differences between song categories vary among wood-warbler species. For example, category A

330  songs in Golden-cheeked Warblers (*Setophaga chrysoparia*) are shorter with fewer notes, the opposite of the pattern

331  identified in the current study (Bolsinger, 2000). In Pine Warblers (Price & Crawford, 2013), trill rates are

332  significantly higher in category A songs. Similarly, Yellow Warblers sing category A songs with lower vocal

333  deviation (corresponding to higher rates of frequency modulation) than category B songs (Beebee, 2004b). Taken

334  together with our results, these findings suggest that no single explanation can account for between-category

335  structural variation among wood-warblers with divided repertoires.

336  We used three lines of evidence to rank song delivery variables by their importance for clustering: the

337  hierarchical structure of the four best models, the average silhouette indices of the one-variable clustering schemes,

338  and the magnitude of between-group differences in the four-variable scheme. The ranks of residual run length,

339  residual latency and percent pre-dawn varied among analyses. Average ranks across analyses suggest that residual

340  run length (mean rank = 1.7) is the most important variable for clustering, followed by percent pre-dawn (2.0),

341  residual latency (2.3), and percent matching (4), which consistently ranked as the least important variable. This

342  ranking suggests a different approach from previous efforts to categorize song in Adelaide's Warbler, which relied

343  primarily on time of delivery and secondarily on latency (see Introduction). We propose three non-mutually-

344  exclusive explanations for the finding that percent matching was weakly predictive of category membership. First,

345  percent matching may not be very different between categories. Second, song type matching may be too rare to

346    accurately estimate percent matching with our sample size (sampling error). Third, we may have mismeasured

347    matching rates (measurement error). Specifically, it is likely we missed some instances of this behavior (e.g., when a

348    neighbor's song was audible to the focal bird, but did not show up on a spectrogram).

349        We favor the three-variable clustering scheme based on residual latency, residual run length, and percent

350    pre-dawn over the other clustering schemes considered in this study. When variables are not perfectly predictive of

351    category membership, the average silhouette index tends to correlate negatively with the number of variables

352    because each additional variable generates another dimension in which distance can be measured. Comparing

353    average silhouette indices across clustering schemes with different numbers of variables is therefore not a sound

354    method to identify natural clusters. Based on *a priori* expectations about song categories in wood-warblers, we

355    expected that our samples should contain examples of all categories from all males, which was not the case in the

356    one-variable and two-variable schemes. Prior research indicates that two clusters were more likely than three, further

357    discounting the validity of the two-variable scheme (Spector 1992). The three-variable and four-variable schemes

358    were very similar, but we prefer the three-variable scheme because (1) it is easier to measure three variables than

359    four, (2) the cluster corresponding to category B is more compact in the three-variable scheme, (3) and song type

360    matching (which contributes to the four-variable scheme but not the three-variable scheme) is harder to accurately

361    estimate than the other variables. Future studies could assign novel song types to categories by comparing the

362    distance in feature space between novel songs and the centroids of the two song categories using the three-variable

363    clustering scheme (Table 3).

364        The goal of this study was to test for and characterize song categories. In doing so, we set the stage to

365    address functional differences between song categories in future studies. Several studies of the behavioral context in

366    which birds sing songs from the two categories have concluded that first category songs are directed at females and

367    second category songs are directed at other males in *Setophaga* (Bolsinger, 2000; Kroodsma et al., 1989; Spector,

368    1991; Staicer, 1996b; Wiley et al., 1994). Other studies, however, have failed to find support for key prediction of

369    one or both of those hypotheses (Beebee, 2004a; MacNally & Lemon, 1985). For example, both females (Beebee

370    2004a) and males (Beebee 2004a, MacNally & Lemon 1985, Weary et al. 1992, Weary et al. 1994, but see Kelly &

371    Ward 2017) tend to respond similarly to playback of category A and category B songs. Our finding that song type

372    matching and song sharing were higher for category B songs are consistent with the hypothesis that category B

373    songs are important for male-male vocal interactions in this species. It is not clear, however, why category B songs

374    would be limited to the breeding season if they are only used for male-male communication. Nor is clear how

375    between-category variation in song structure might facilitate communication with different classes of receivers.

376    Information about when females prospect for mates, how female presence affects males' choice of song categories,

377    and whether females attend to male-male singing interactions like song type matching (Logue & Forstmeier, 2008),

378    would help to clarify the functional differences between song categories.

379          In summary, we applied an unsupervised clustering algorithm to Adelaide's Warbler song types, and

380    identified two song categories with different song delivery parameters. The delivery styles in the two categories are

381    similar to patterns described in several other *Setophaga* species, and are largely consistent with previous studies on

382    the focal species. Our findings do not allow us to exclude the possibility that song categories are not entirely

383    discrete. Whether the use of song categories, as defined in this study, varies seasonally remains an open question

384    (Table 1). We found evidence that category B songs are used more than category A songs during male-male vocal

385    interactions, but further studies are required to understand the functional significance of song categories in this

386    species.

387

388    **References**

389    Anzanello, M. J., & Fogliatto, F. S. (2011). Selecting the best clustering variables for grouping mass-customized

390          products involving workers' learning. *International Journal of Production Economics, 130*(2), 268-276.

391    Baker, M. (2016). 1,500 scientists lift the lid on reproducibility. *Nature, 533*(7604), 452-454. doi:10.1038/533452a

392    Bates, D., Machler, M., Bolker, B. M., & Walker, S. C. (2015). Fitting linear mixed-effects models using lme4.

393          *Journal of Statistical Software, 67*(1), 1-48.

394    Beebee, M. D. (2002). Song sharing by yellow warblers differs between two modes of singing: Implications for song

395          function. *Condor, 104*(1), 146-155. doi:10.1650/0010-5422(2002)104[0146:Ssbywd]2.0.Co;2

396    Beebee, M. D. (2004a). The functions of multiple singing modes: Experimental tests in yellow warblers, *dendroica*

397          *petechia*. *Animal Behaviour, 67*(6), 1089-1097. doi:10.1016/j.anbehav.2003.05.016

398    Beebee, M. D. (2004b). Variation in vocal performance in the songs of a wood-warbler: Evidence for the function of

399          distinct singing modes. *Ethology, 110*(531-542).

400    Bolsinger, J. S. (2000). Use of two song categories by golden-cheeked warblers. *The Condor, 102*, 539-552.

401  Catchpole, C. K., & Slater, P. J. B. (2008). *Bird song: Biological themes and variations* (2 ed.). Cambridge:

402     Cambridge University Press.

403  Charrad, M., Ghazzali, N., Boiteau, V., & Niknafs, A. (2014). Nbclust: An r package for determining the relevant

404     number of clusters in a data set. *Journal of Statistical Software, 61*(6), 1-36.

405  Curson, J., Quinn, D., & Beadle, D. (1984). *New world warblers*. London: A&C Black.

406  Demko, A. D., Reitsma, L. R., & Staicer, C. A. (2013). Two song categories in the canada warbler (*cardellina*

407     *canadensis*). *The Auk, 130*(4), 609-616. doi:10.1525/auk.2013.13059

408  Hastie, R., & Dawes, R. M. (2010). *Rational choice in an uncertain world: The psychology of judgment and*

409     *decision making*: Sage.

410  Hedley, R. W., Logue, D. M., Benedict, L., & Mennill, D. J. (2018). Assessing the similarity of song-type transitions

411     among birds: Evidence for interspecies variation. *Animal Behaviour, 140*, 161-170.

412     doi:10.1016/j.anbehav.2018.04.008

413  Hof, D., & Podos, J. (2013). Escalation of aggressive vocal signals: A sequential playback study. *Proceedings of the*

414     *Royal Society of London B: Biological Sciences 280.1768*, 20131553.

415  Illes, A. E. (2015). Context of female bias in song repertoire size, singing effort, and singing independence in a

416     cooperatively breeding songbird. *Behavioral Ecology and Sociobiology, 69*(1), 139-150.

417     doi:10.1007/s00265-014-1827-3

418  Kroodsma, D. E., Bereson, R. C., Byers, B. E., & Minear, E. (1989). Use of song types by the chestnut-sided

419     warbler: Evidence for both intra-sexual and inter-sexual functions. *Canadian Journal of Zoology, 67*(2),

420     447-456.

421  Lachlan, R. F. (2007). Luscinia: A bioacoustics analysis computer program. Version 1.0. *Computer program].*

422     *Retrieved from luscinia. sourceforge. net on October, 8*, 2012.

423  Lemon, R. E., Cotter, R., MacNally, R. C., & Monette, S. (1985). Song repertoires and song sharing by american

424     redstarts. *The Condor, 87*(4), 457-470. doi:10.2307/1367942

425  Logue, D. M., & Forstmeier, W. (2008). Constrained performance in a communication network: Implications for the

426     function of song-type matching and for the evolution of multiple ornaments. *Am Nat, 172*(1), 34-41.

427     doi:10.1086/587849

428    Macdougall-Shackleton, S. A. (1997). Sexual selection and the evolution of song repertoires. *Current ornithology*,

429        81-124.

430    MacNally, R. C., & Lemon, R. E. (1985). Repeat and serial singing modes in american redstarts (s*etophaga*

431        *ruticilla*): A test of functional hypotheses. *Zeitschrift fuˉr Tierpsychologie, 69*, 191-202.

432    Molles, L. E., & Vehrencamp, S. L. (1999). Repertoire size, repertoire overlap, and singing modes in the banded

433        wren (thryothorus pleurosticus). *The Auk, 116*(3), 677-689.

434    Podos, J., Moseley, D. L., Goodwin, S. E., McClure, J., Taft, B. N., Strauss, A. V. H., . . . Lahti, D. C. (2016). A

435        fine-scale, broadly applicable index of vocal performance: Frequency excursion. *Animal Behaviour, 116*,

436        203-212. doi:10.1016/j.anbehav.2016.03.036

437    Price, J. J., & Crawford, C. L. (2013). Use and characteristics of two singing modes in pine warblers. *The Wilson*

438        *Journal of Ornithology, 125*(3), 552-561. doi:10.1676/13-006.1

439    Rousseeuw, P. J. (1987). Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal*

440        *of Computational and Applied Mathematics, 20*, 53-65. doi:https://doi.org/10.1016/0377-

441        0427(87)90125-7

442    Schraft, H. A., Medina, O. J., McClure, J., Pereira, D. A., & Logue, D. M. (2017). Within-day improvement in a

443        behavioural display: Wild birds 'warm up'. *Animal Behaviour, 124*, 167-174.

444        doi:10.1016/j.anbehav.2016.12.026

445    Spector, D. A. (1991). The singing behaviour of yellow warblers. *Behaviour, 117*(1/2), 29-52.

446    Spector, D. A. (1992). Wood-warbler song systems. In D. M. Power (Ed.), *Current ornithology* (pp. 199-238).

447        Boston, MA: Springer US.

448    Staicer, C. A. (1989). Characteristics, use, and significance of two singing behaviors in grace's warbler (*dendroica*

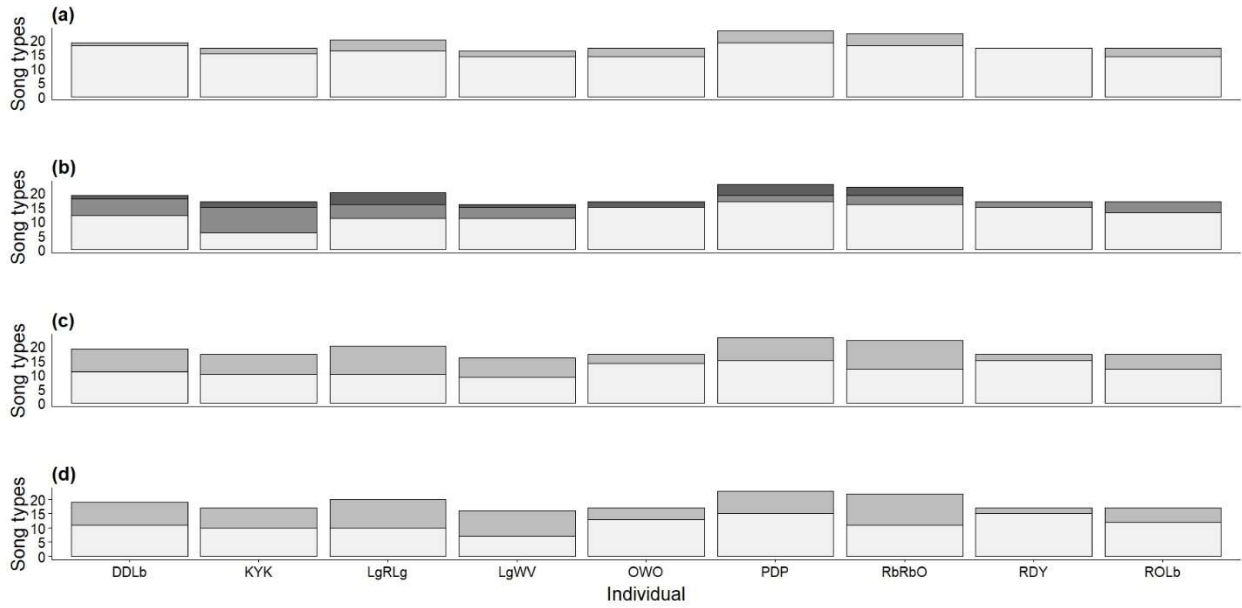449        *graciae*). *The Auk, 106*(1), 49-63. doi:10.2307/4087756

450    Staicer, C. A. (1991). *The role of male song in the scioecology of the tropical resident adelaide's warbler*

451        *(dendroica adelaidae).* (PhD), University of Massechussetts.

452    Staicer, C. A. (1996a). Acoustical features of song categories of the adelaide's warbler (*dendroica adelaidae*). *The*

453        *Auk, 113*(4), 771-783.

454    Staicer, C. A. (1996b). Honest advertisement of pairing status: Evidence from a tropical resident wood-warbler.

455        *Anim Behav, 51*, 375-390.

456    Stutchbury, B. J., & Morton, E. S. (2001). *Behavioral ecology of tropical birds*: Academic Press.

457    Toms, J. D. (2010). Adelaide's warbler (*setophaga adelaidae*). *Neotropical Birds Online.* version 1.0. Retrieved

458        from https://doi.org/10.2173/nb.adewar1.01

459    D.M Weary, R.E Lemon, S Perreault. (1992). Song repertoires do not hinder neighbor–stranger discrimination.

460        *Behavioral Ecology and Sociobiology*, 31 , 441-447.

461    D.M Weary, R.E Lemon, S Perreault. (1994). Different responses to different song types in American redstarts

462        Auk, 111, 730-734.

463    Wiley, H. R., Godard, R., & Thompson, A. D. (1994). Use of two singing modes by hooded warblers as adaptations

464        for signalling. *Behaviour, 129*(3), 243-278. doi:http://dx.doi.org/10.1163/156853994X00631
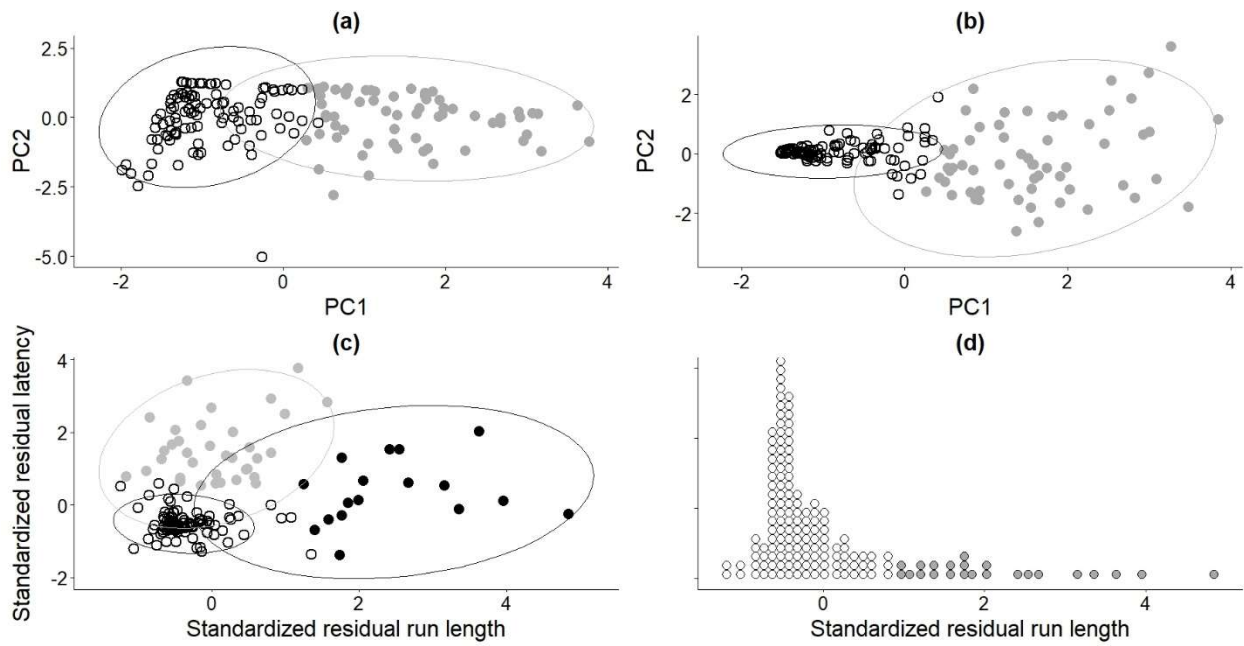
465

**Figure 1.** Assignment of song types to song categories by individual male Adelaide's Warbler according to the best (a) four-variable, (b) three-variable, (c) two-variable, and (d) one-variable clustering schemes. White bars represent category B songs, and gray bars represent category A songs. For the two-variable clustering scheme, dark gray represents category A1 and light gray represents category A2.

475
476



**Figure 2.** Cluster plots showing the best (a) four-variable, (b) three-variable, (c) two-variable, and (d) one-variable

clustering schemes. Open circles represent category B songs, and filled circles represent category A songs. For the

two-variable clustering scheme, black circles represent category A1 songs and gray circles represent category A2

songs.

482
483
484
485
486

487    **Table 1.** Hypothesized song delivery patterns of category A and category B songs in Adelaide's warbler.

| Attribute | Category A | Category B |
|---|---|---|
| Time of year* | All year | Breeding season only |
| Time of day | After dawn | Before dawn and sporadically after dawn |
| Social context | Male-female interactions | Male-male interactions |
| Delivery mode | Eventual variety | Immediate variety |
| Latency between songs | Longer | Shorter |

488    After Staicer (1991)

489    * Time of year was not examined in the present study.

490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527

528
529
530
531
532
533

**Table 2.** Fifteen clustering schemes for Adelaide's Warbler song types.  Variables included in each model are marked with an "X."

| % predawn | residual latency | residual run length | % matching | # of variables | best k | ASI of best k |
|:---:|:---:|:---:|:---:|:---:|:---:|:---:|
| **X** | **X** | **X** | **X** | **4** | **2** | **0.379** |
| **X** | **X** | **X** | | **3** | **2** | **0.495** |
| | X | X | X | 3 | 2 | 0.433 |
| X | | X | X | 3 | 5 | 0.404 |
| X | X | | X | 3 | 2 | 0.402 |
| | **X** | **X** | | **2** | **3** | **0.617** |
| X | | X | | 2 | 3 | 0.583 |
| X | X | | | 2 | 3 | 0.552 |
| X | | | X | 2 | 3 | 0.491 |
| | | X | X | 2 | 3 | 0.490 |
| | X | | X | 2 | 3 | 0.450 |
| | | **X** | | **1** | **2** | **0.771** |
| X | | | | 1 | 2 | 0.738 |
| | X | | | 1 | 2 | 0.730 |
| | | | X | 1 | 7 | 0.688 |

534
535 ASI = Average silhouette index.
536 The models with the highest ASI for each number of variables are in bold font.

**Table 3.** Means ± standard deviations of song delivery variables from the best (highest average silhouette index) categorization schemes based on one, two, three, and four song delivery variables.

| # of variables | cluster | % predawn | latency (s) | residual latency | run length (songs) | residual run length | % match | category |
|---|---|---|---|---|---|---|---|---|
| 4 | 1 | 71.87 ± 26.64 | 21.92 ± 16.48 | 1.64 ± 17.59 | 2.24 ± 2.71 | 0.09 ± 0.41 | 5.62 ± 4.98 | B |
|  | 2 | 7.06 ± 13.43 | 112.43 ± 48.25 | 70.87 ± 50.59 | 5.85 ± 4.60 | 1.26 ± 1.49 | 3.40 ± 4.23 | A |
| 3 | 1 | 69.89 ± 28.38 | 23.79 ± 18.54 | 2.74 ± 18.16 | 2.15 ± 2.23 | 0.07 ± 0.35 | 5.17 ± 4.54 | B |
|  | 2 | 6.31 ± 12.24 | 115.09 ± 48.86 | 73.52 ± 51.16 | 6.25 ± 4.91 | 1.37 ± 1.51 | 4.05 ± 5.24 | A |
| 2 | 1 | 63.21 ± 33.81 | 27.41 ± 22.69 | 3.16 ± 15.68 | 2.35 ± 2.38 | 0.1 ± 0.43 | 5.02 ± 4.48 | B |
|  | 2 | 3.22 ± 6.47 | 90.72 ± 41.89 | 45.1 ± 42.49 | 11.65 ± 5.32 | 3.31 ± 1.13 | 4.16 ± 6.70 | A1 |
|  | 3 | 15.40 ± 24.69 | 135.80 ± 46.55 | 102.09 ± 40.99 | 3.93 ± 2.48 | 0.6 ± 0.7 | 4.25 ± 4.92 | A2 |
| 1 | 1 | 4.07 ± 6.93 | 102.95 ± 61.51 | 55.34 ± 62.87 | 11.45 ± 4.97 | 2.93 ± 1.17 | 3.67 ± 5.86 | A |
|  | 2 | 54.02 ± 37.30 | 49.02 ± 49.87 | 23.68 ± 43.79 | 2.37 ± 1.71 | 0.15 ± 0.44 | 4.95 ± 4.63 | B |

541 **Table 4:** Effects of song category on song structure, based on the three-variable clustering scheme. Estimated

542 effects are relative to category A songs. See Table S2 for equivalent data from the other clustering schemes.

| | duration (ms) | notes | trill rate (notes / s) | FEX | $F_{min}$ (kHz) | $F_{max}$ (kHz) | bandwidth (kHz) |
|---|---|---|---|---|---|---|---|
| *fixed effects* | | | | | | | |
| Intercept | 2098.16 | 25.00 | 11.91 | 66.11 | 3178.41 | 6137.28 | 2951.89 |
| category B | -123.67*** | -2.12*** | -0.31* | 1.44 | -46.02 | -22.22 | 32.10 |
| | | | | | | | |
| *random effects* | | | | | | | |
| song type | 5479 | 2.72 | 0.28 | 16.15 | 16.805 | 11.78 | 17.82 |
| bird:day | 3998 | 1.01 | 0.15 | 6.14 | 5.767 | 11.32 | 7.52 |
| Bird | 691 | 0.52 | 0.06 | 2.25 | 0.283 | 3.81 | 1.29 |
| | | | | | | | |
| *Average ± SD* | | | | | | | |
| A | 2119.66 ± 291.14 | 25.70 ± 5.24 | 12.09 ± 1.56 | 65.87 ± 9.70 | 3190.03 ± 402.29 | 6115.27 ± 475.84 | 2925.25 ± 451.16 |
| B | 2119.66 ± 291.14 | 25.70 ± 5.24 | 12.09 ± 1.56 | 65.87 ± 9.70 | 3190.03 ± 402.29 | 6115.27 ± 475.84 | 2925.25 ± 451.16 |

543 $* P < 0.05, ** P < 0.01, *** P < 0.0001$

544

545

546

547

548