

1

2

### 3 Supplementary Figure 1

4

5 **Pitch error corrections in two syllables that are misaligned with respect to each**

6 **other.** (a) Training models of two experimental birds in imitation task 2 (left, Audio

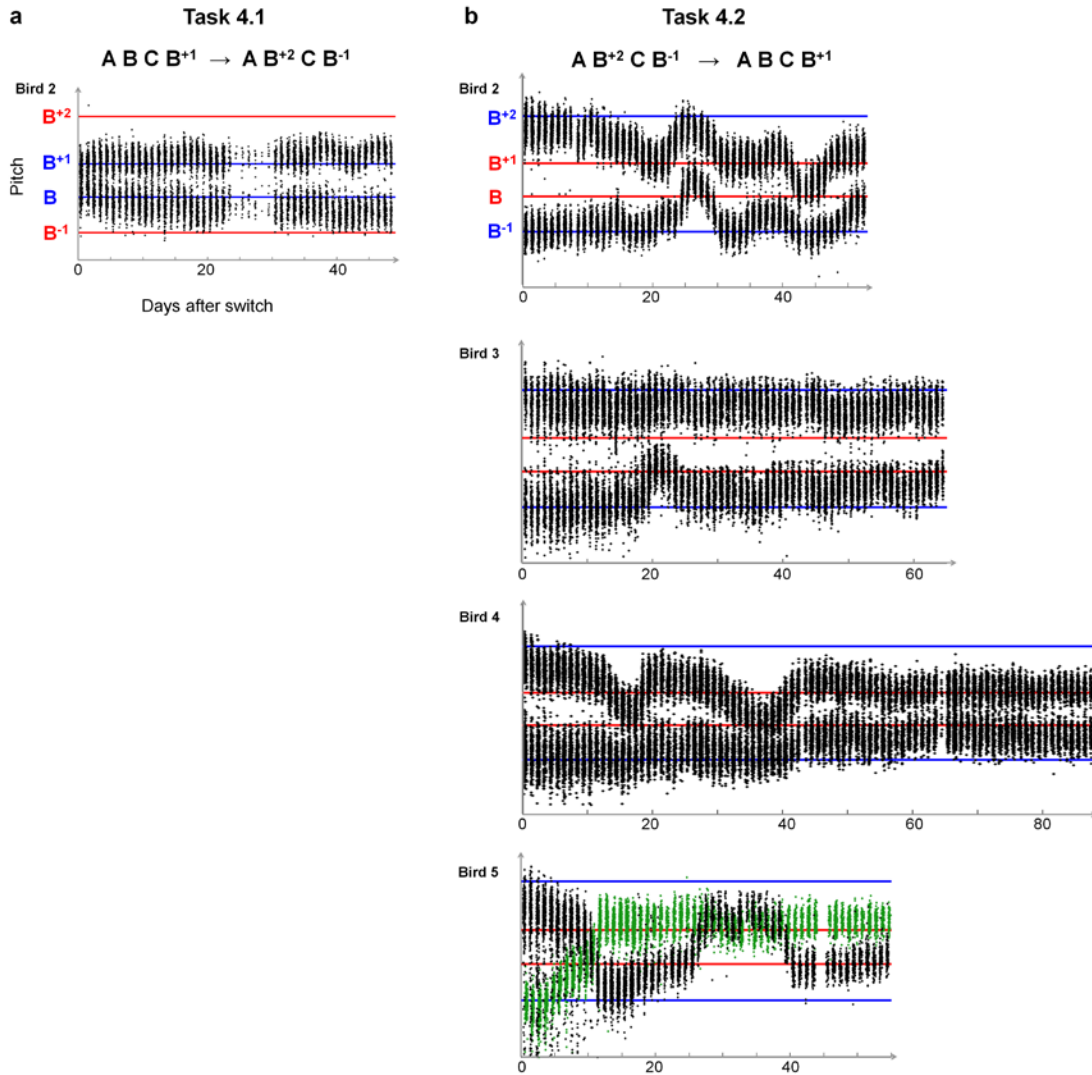
7 5, playbacks 1 and 2; right, Audio 6, playbacks 3 and 4; see Supplementary Table 2).

8 Scale bars for sonograms are 100ms (x axis) and 2kHz (y axis). (b) The median pitch

9 of consecutive renditions of syllable A (top) and C (bottom) in bird 1 and 2. Both

10 pitch errors were successfully corrected. (c) Stack plots showing consecutive motif

11 renditions in birds 1 and 2; colors, pitch of syllables A/A<sup>t</sup> and C/C<sup>t</sup> (t= target pitch);  
12 grayscale, Wiener entropy in neighboring syllables (as in Fig. 2b). Sonograms at  
13 bottom and top show song at start and end points. Birds corrected pitch errors in  
14 syllable A and C before changing syntax (bird 1 did not change syntax at all; bird 2  
15 matched the target syntax). **(d)** Fraction of pitch error correction (left) and time (days)  
16 to reach 50% pitch match (right) in syllables A and C across experimental birds.  
17 Black, individual birds (lines connect the two syllables in each bird); red, mean ±  
18 s.e.m; n=8.  
19

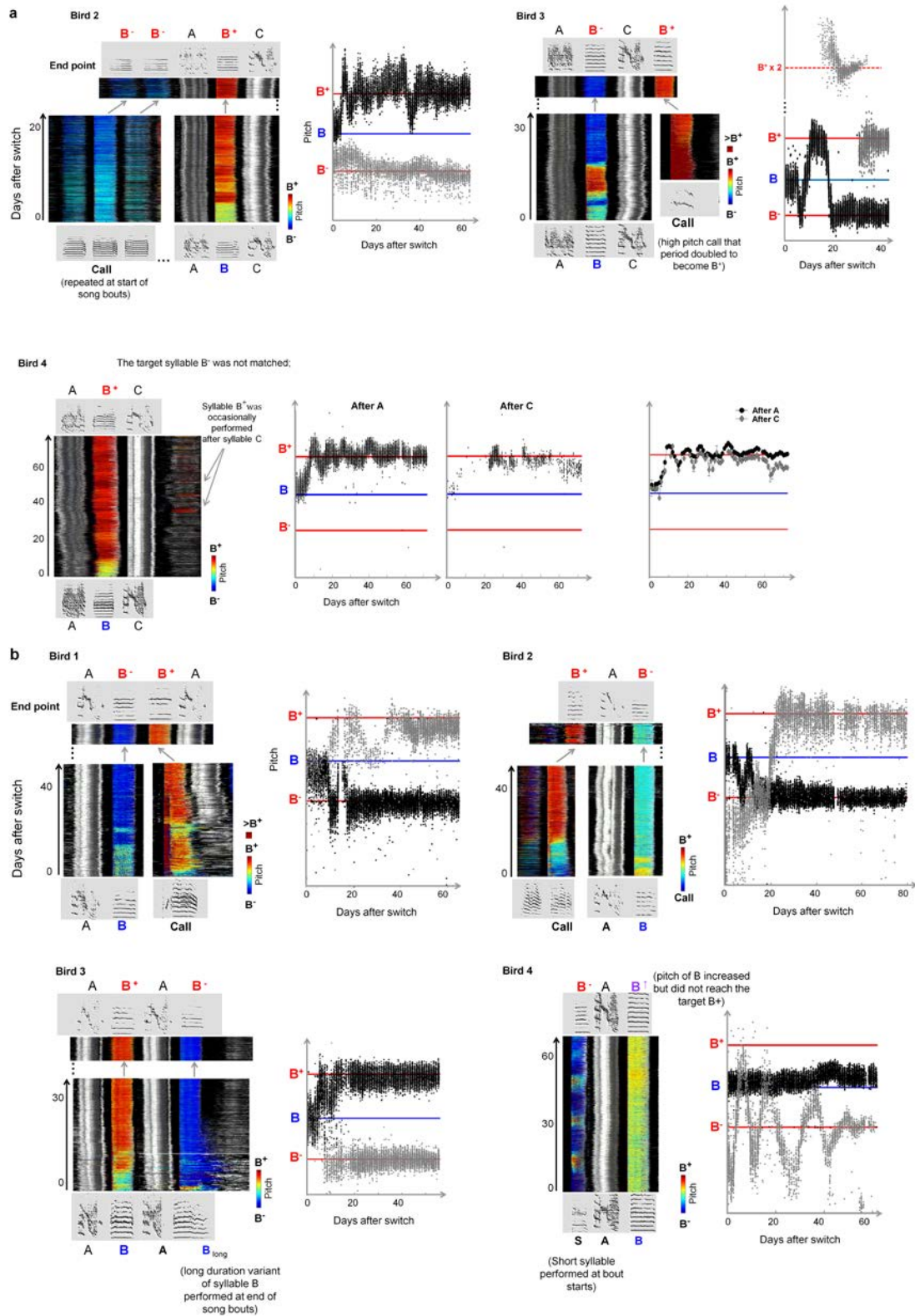


20

21 **Supplementary Figure 2**

22

23 **Pitch trajectories in birds trained with tasks 4.1 and 4.2.** Median pitch in  
 24 consecutive renditions of the two pitch shifted syllable types in birds trained with task  
 25 4.1 ( $ABCB^{+1} \rightarrow AB^{+2}CB^{-1}$ ) (**a**) and task 4.2 ( $AB^{+2}CB^{-1} \rightarrow ABCB^{+1}$ ) (**b**), not  
 26 including the two birds depicted in Fig 4c. Notation as in Fig 4c. In all birds except  
 27 one (bird 5 in (**b**)), the pitch of both syllable types shifted towards the spectrally  
 28 closer targets. In bird 5, the pitch of both syllable types (shown in black and green for  
 29 visual clarity) shifted towards the spectrally farther targets. Bird ages at the end of the  
 30 experimental period in task 4.2 were 121, 121, 128, 130 and 153 days post hatch. As  
 31 the sensitive period for song learning in zebra finches ends around day 90-100 post  
 32 hatch, it is unlikely that birds in this group that matched the 1-semitone targets were  
 33 on the way to matching the farther targets.



34

35

36

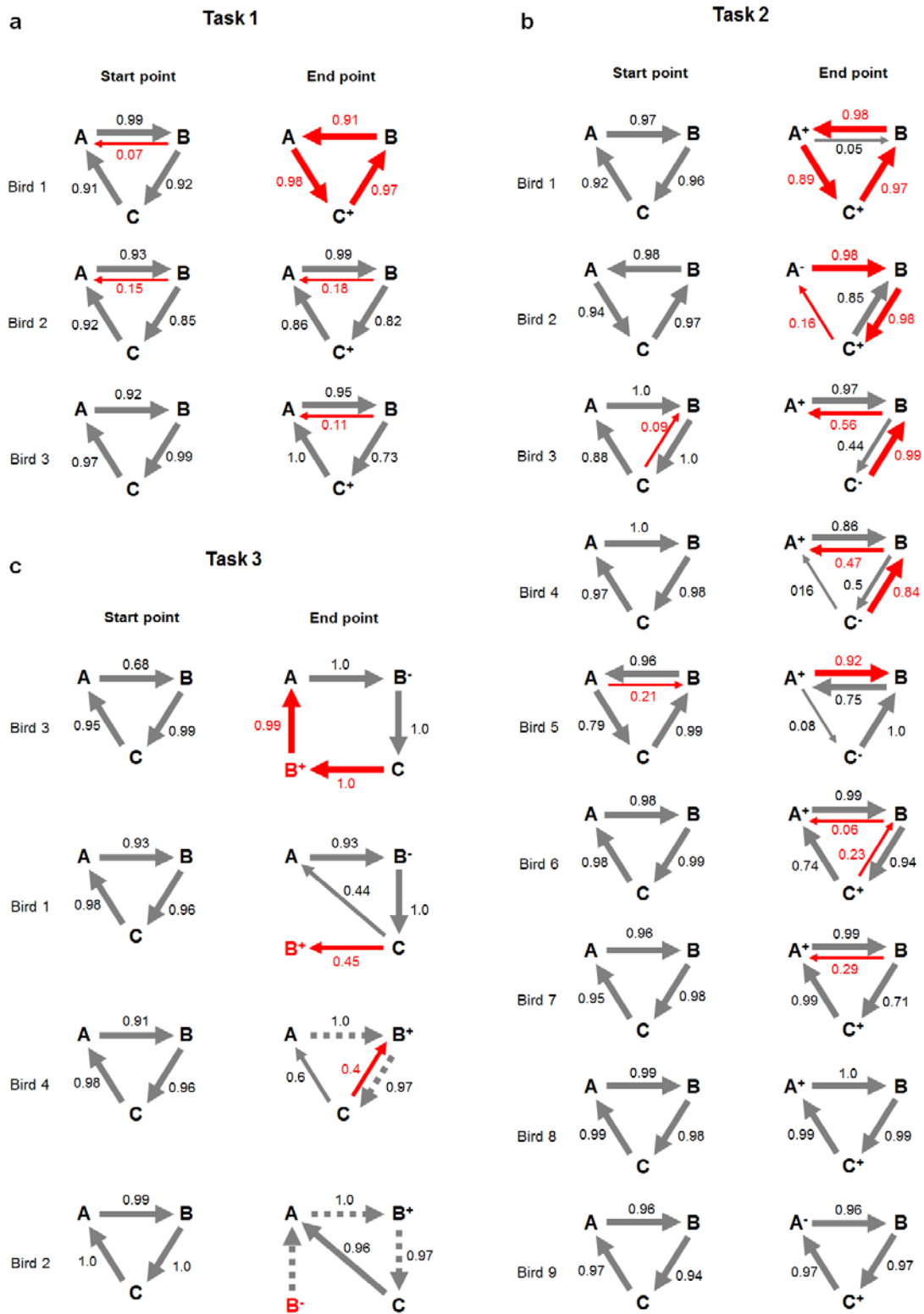
37 **Supplementary Figure 3**

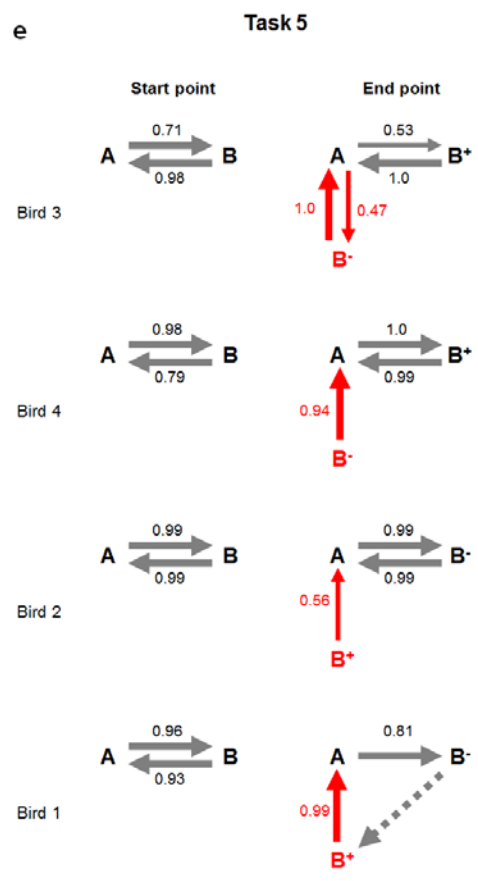
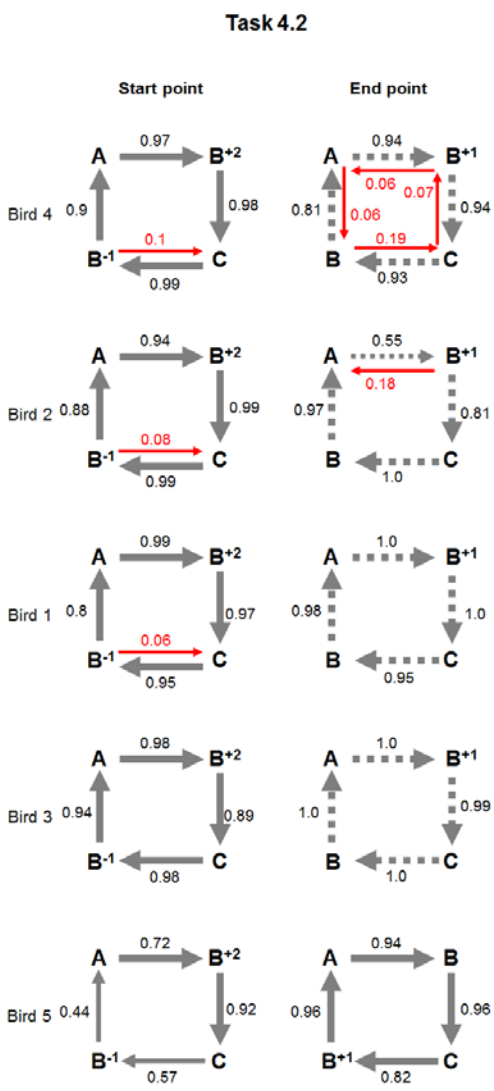
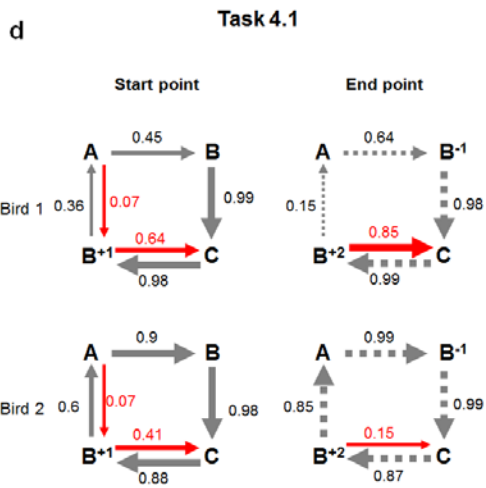
38

39 **Matching a vacant target with a vocalization initially performed outside of the**  
 40 **song motif.**

41 Developmental trajectories of experimental birds trained with tasks containing an  
 42 extra syllable type in the target versus the source: task 3 (ABC → AB<sup>-</sup>CB<sup>+</sup>) (a) and

43 task 5 ( $AB \rightarrow AB^+AB^-$ ) (**b**). Stack plots and pitch trajectories as in Fig 5 c and d  
44 (depicting Bird 1 of task 3). All birds except one (Bird 4 in (**a**)) matched the vacant  
45 target with a syllable type initially external to the song motif, usually a call. In bird 4  
46 (**a**, bottom), the target  $B^-$  was not matched. Syllable B shifted to  $B^+$  in the “wrong”  
47 context (namely, after syllable A), but was also performed sparsely after syllable C;  
48 pitch trajectories in this bird are shown separately for renditions after A and after C  
49 (middle plots); right-most plot shows daily pitch means  $\pm$  s.e.m. for renditions after A  
50 (black circles) and C (grey diamonds), showing a gradual divergence in pitch  
51 ( $756 \pm 1.6$  Hz after A versus  $731 \pm 1.9$  after C on last experimental day;  $p < 0.00001$ ,  
52 ttest). This could potentially result from incomplete “splitting” of  $B^+$  into two syllable  
53 types.  
54



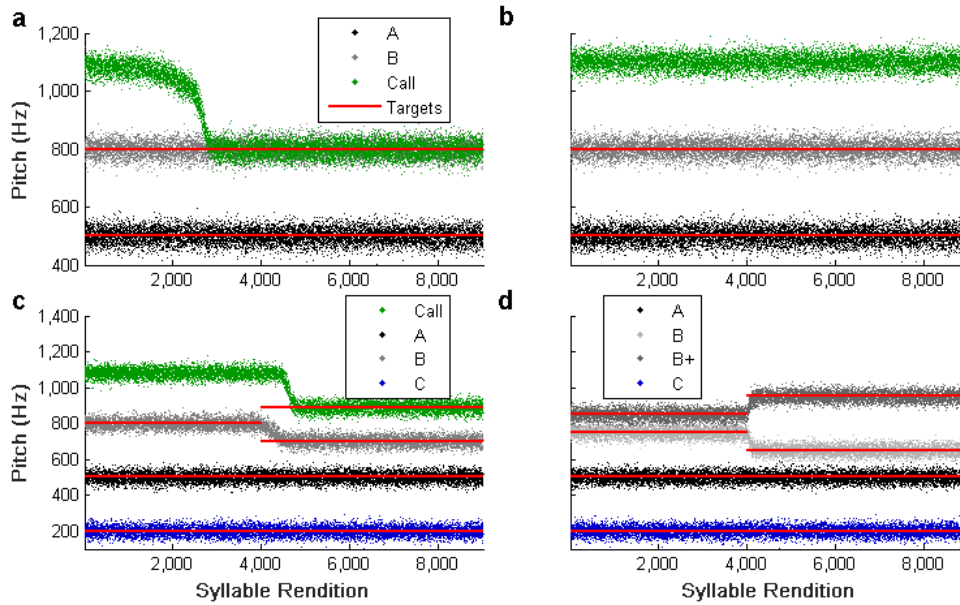


63 **Supplementary Figure 4**

64

65 **Syntax adjustments towards target across experimental groups. a-e**, Syntax  
66 diagrams of all experimental birds (tasks 1-5) at the start point of the experiment (the  
67 day of switching to target training, left) and at its endpoint (last experimental day,  
68 right). Arrows represent the fractions of performing the syllable transitions of the  
69 source song (gray) and the target song (red). The thickness of arrows corresponds to  
70 the fraction values, which are also indicated next to the arrows. Birds' names are the  
71 same as in Supplementary Figs. 2 and 3; birds in each group are vertically arranged  
72 according to the amount of syntax adjustments (maximum adjustment on top). **c-d**,  
73 dotted gray arrows represent incorrect transitions resulting from shifting syllable pitch  
74 to misaligned target syllables (note that bird 5 in **(d)** shifted pitch to the aligned  
75 targets, and therefore no syntax adjustments were necessary). Newly generated  
76 syllables in tasks 3 and 5 (**c** and **e**) are shown in red. In two cases (Bird 2 in **c** and Bird  
77 1 in **e**) the new syllable was placed in an incorrect context in the song motif, indicated  
78 by a dotted gray arrow.





79  
80  
81  
82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93  
94  
95  
96  
97  
98

### Supplementary Figure 5

#### Birds allocate vocalizations to models more sparsely than the EM algorithm. (a)

A bird singing two-syllable-song with syllable  $A$  (black dots) and  $B$  (grey dots) matching two targets (red lines). The unassigned call (green dots) converges to the occupied target  $B$  in case of the standard EM algorithm (without Step  $N3$ ), but (b) in case of the musical-chairs enhanced EM algorithm (including Step  $N3$ ) in which posterior probabilities can be large for at most one syllable, the call is not attracted towards an occupied target. (c) Simulation of the  $ABC \rightarrow AB^-CB^+$  task. After switching to the new targets (discontinuity in red lines) syllable  $B$ 's pitch (grey dots) moves almost instantly to  $B^-$  and the pitch of the (initially) unassigned call (green) moves gradually to  $B^+$  (top red line). Syllable  $A$  (black dots) and syllable  $C$  (blue dots) stay on their targets. (d) Simulation of the  $ABCB^{+1} \rightarrow AB^{+2}CB^{-1}$  task. After switching to the new targets the pitches of syllables  $B$  and  $B^+$  shift greedily to the closest targets, i.e.  $B \rightarrow B^{-1}$  (lighter grey) and  $B^{+1} \rightarrow B^{+2}$  (darker grey). Colored dots represent individual syllable renditions. The simulation parameters ( $\alpha = 0.02, \alpha_2 = 0.5, \sigma = 25, \sigma_b = 75$ ) are identical in all simulations.

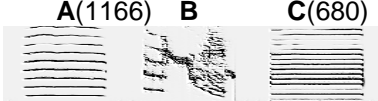
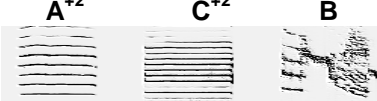
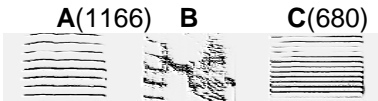
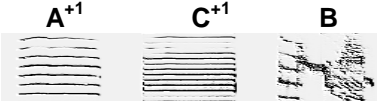
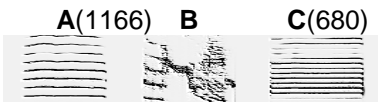

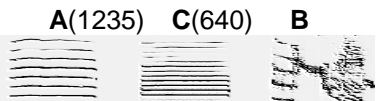
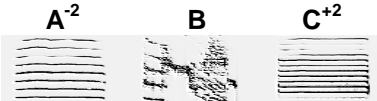

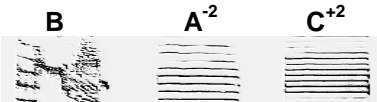
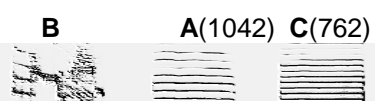
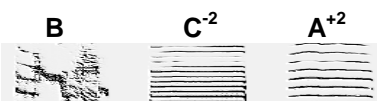
99 **Supplementary Table 1. Training models used in imitation task 1.**  
 100

n	Source model (baseline pitch of syllable C, Hz)	Target model (superscripts, pitch shift from baseline, semitones)	Sound file
1	<p>A B C(680)</p>	<p>A C<sup>+2</sup> B</p>	Audio 1 (playbacks 1 and 2)
2	<p>A B C(680)</p>	<p>A C<sup>+1</sup> B</p>	Audio 1 (playbacks 3 and 4 )

101 Each model playback contained two repetitions of the motif shown (Audio 1; same  
 102 for Supplementary Tables 2-5).

103  
104





**Supplementary Table 2. Training models used in imitation task 2.**

n	Source model (baseline pitch of syllables A and C, Hz)	Target model (superscripts, pitch shift from baseline, semitones)	Sound file
2			Audio 5 (playbacks 1 and 2)
2			Audio 5 (playbacks 3 and 4)
2			Audio 6 (playbacks 1 and 2)
1			Audio 6 (playbacks 3 and 4)
1			Audio 7 (playbacks 1 and 2)
1			Audio 7 (playbacks 3 and 4)

105  
106

107 **Supplementary Table 3. Training models used in imitation task 3.**

108

n	Source model (baseline pitch of syllable B, Hz)	Target model (superscripts, pitch shift from baseline, semitones)	Sound file
3			Audio 8 (playbacks 1 and 2)
1			Audio 8 (playbacks 3 and 4)

109

110

111 **Supplementary Table 4. Training models used in imitation task 4.**

112

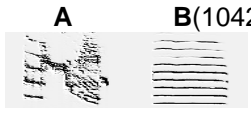
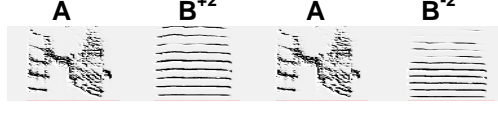
n	Source model	Target model	Sound file
2			Audio 9 (playbacks 1 and 2)
1			Audio 9 (playbacks 3 and 4)
4			Audio 9 (playbacks 5 and 6)

113 All indicated pitch shifts are with respect to the baseline pitch of syllable B (1166

114 Hz).

115 **Supplementary Table 5. Training models used in imitation task 5.**

116

n	Source model (baseline pitch of syllable B, Hz)	Target model (pitch shift from baseline, semitones)	Sound file
4			Audio 10

117 To avoid a large duration difference between source and target playbacks, in this task  
 118 the source playbacks included 4 motif repetitions (Audio 10).

119

120 **Supplementary Notes**

121

122 Mathematical Supplement:

123

124 Birds simplify the quadratic problem of computing performance error to a linear  
 125 assignment problem

126

127 Performance error

128 Associated with a song  $S$  we define a family of performance errors  $E(S, \Delta)$   
 129 parameterized by a set of unknown parameters grouped in the matrix  $\Delta = (\delta_{i,j})$ . The  
 130 errors are composed of an overall **phonology** (spectral) error and a **syntax** (sequence)  
 131 error

132 
$$E(S, \Delta) = \underbrace{\sum_{i,j} e(S_j, T_i) \delta_{i,j}}_{\text{phonology}} + c \underbrace{\|\Delta\|_{\#}}_{\text{syntax}} \quad (1)$$

133 The parameter  $c$  represents an unknown tradeoff between phonology and syntax  
 134 errors.

135 The overall **phonology error** between song  $S$  and target  $T$  is defined as a  
 136 weighted sum of the local phonological errors  $e(S_j, T_i)$  between song element  $S_j$   
 137 (e.g. syllable  $j = 1, \dots, n$ ) and target element  $T_i$  ( $i = 1, \dots, m$ ). The function  $e$   
 138 represents a distance metric, for example the Euclidean distance between specific  
 139 sound features such as pitch (e.g. pitch deviation). The unknown assignment matrix  
 140  $\Delta = (\delta_{i,j})$  specifies the weight  $\delta_{i,j}$  associated with the local phonology error, where  
 141 target assignments (which syllables are assigned to a specific target) correspond to  
 142 rows of  $\Delta$  and syllable assignments (which targets are assigned to a specific syllable)  
 143 to columns of  $\Delta$ . To illustrate this notation, a bird that does not assign a phonology  
 144 error to syllable  $S_2$  entails  $\delta_{i,2} = 0$  for all  $i$ ; and, a bird that compares  $S_2$  with the  
 145 first target (syllable)  $T_1$  entails  $\delta_{1,2} = 1$ . If there is local chaining of assignments then a  
 146 bird that compares  $S_1$  to  $T_3$  will also compare  $S_2$  to  $T_4$ ; in terms of  $\Delta$ , chaining of  
 147 assignments means that the condition  $\delta_{i,j} = 1$  implies  $\delta_{i+1,j+1} = 1$  with high probability.  
 148 By virtue of the assignment weights  $\delta_{i,j}$ , the phonology errors may parameterize any  
 149 imaginable comparison between song and target.

150

151 The **syntax error**  $\|\Delta\|_{\#}$  quantifies the amount of resequencing a bird must  
 152 perform in order to bring its song elements into global alignment with the template.  
 153 This error quantifies the new transitions to be created among existing song elements.  
 154 Because of stepwise acquisition of syntax in songbirds<sup>1</sup>, it makes sense to attribute to  
 155  $\|\Delta\|_{\#}$  a cost proportional to the number of new transitions to be generated.

156

157 In the case of binary and one-to-one syllable-target assignments ( $\delta_{i,j} = 0$  or  $1$ ;  
 158 up to one target per syllable,  $\sum_i \delta_{i,j} \leq 1, \forall j$ , and one syllable per target,

159  $\sum_j \delta_{i,j} = 1, \forall i$ ) we can write  $\|\Delta\|_{\#}$  as a sum of terms that are quadratic in the  
 160 assignment weights  $\delta_{i,j}$ :

$$161 \quad \|\Delta\|_{\#} = \sum_{i,j}^{m,n} \sum_{k \neq \lfloor j+1 \rfloor_n} \delta_{i,j} \delta_{\lfloor i+1 \rfloor_m, k} \quad (2)$$

162 where we use the short-hand notation  $\lfloor x \rfloor_l$  to denote  $1 + (x-1) \bmod l$ , which is  
 163 necessary to incorporate the circular boundary conditions arising from birds' tendency  
 164 to repeat motifs several times in a song bout. As can be seen, Equation (2) skips all  
 165 pairs of consecutive syllable assignments  $j$  and  $\lfloor j+1 \rfloor_n$  that do not need  
 166 resequencing, namely all those that are locally diagonally chained.

167  
 168 Equations (1) and (2) model the performance error attributed to any given song. The  
 169 assignment matrix is not specified therein; therefore, these equations can be thought  
 170 of as the space of all possible strategies for estimating phonology and syntax errors  
 171 between song and template (Fig. 1a). Our experiments were designed to resolve birds'  
 172 strategy in dealing with phonology and syntax errors and the on/off-diagonal structure  
 173 of the assignment matrix.

174  
 175 Next we describe the process of song learning. In terms of Equation (1), song  
 176 learning is the search of a song  $S$  with vanishing performance error.

177

178

### 179 Song learning

180

181 Song learning is the process of changing the current song  $S$  towards the **final**  
 182 **song**  $S^*$  that minimizes the performance error (ideally  $S^* = T$ ):

$$183 \quad S^* = \arg \min_S E(S, \Delta) \quad (3)$$

184 In the process of song learning,  $\Delta$  is either fixed or it evolves in time, possibly  
 185 giving rise to very complex learning trajectories. If birds want to perform song  
 186 learning optimally, they will try to compute the initially **optimal assignments**  $\Delta^*$ ,  
 187 which are the ones that achieve minimal initial performance error,

$$188 \quad \Delta^* = \arg \min_{\Delta} E(S, \Delta). \quad (4)$$

189 This optimal choice of assignment in Equation (4) is a quadratic assignment  
 190 problem (quadratic in the assignment weights)<sup>2</sup>. In the general case, that problem is  
 191 NP-hard, meaning that there is no known algorithm for solving this problem in  
 192 polynomial time. Moreover, it was proven that such problems do not even have an  
 193 approximation algorithm running in polynomial time<sup>3</sup>. Hence, almost certainly, birds  
 194 neither solve Equation (4) nor an approximation thereof. The question for us was how  
 195 birds actually assign performance errors.

196

197 Whatever birds do, we imagined they must be facing a tradeoff between  
 198 phonology and syntax errors, illustrated by the following example: Consider two birds  
 199 that need to change their songs from syllable sequence  $ABC \ ABC$  to  $ACB \ ACB$ .  
 200 The first bird forms 3 new bigrams ( $AC$ ,  $CB$ , and  $BA$ ) among the existing



201 syllables, which would imply that its syntax error initially is  $c \begin{vmatrix} 1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 1 & 0 \end{vmatrix}_{\#} = 3c$  and the

202 phonology error is zero. However, the second bird achieves the same target song not  
203 by permuting syllables but by making the local transformations  $B \rightarrow C$  and  $C \rightarrow B$ ,

204 implying that this latter bird produces no syntax errors  $\begin{vmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{vmatrix}_{\#} = 0$  but a

205 phonology error of  $2e(B, C)$ . This example illustrates that there is a tradeoff between  
206 phonology and syntax errors and that the latter are avoidable in principle. The reason  
207 is that the second bird learns the song sequence intrinsically by globally aligning the  
208 song to the template ( $\Delta$  is the identity matrix), that latter bird needs only to correct  
209 phonology errors to also automatically learn the correct syntax. However, global  
210 alignment may not be an ideal strategy because it can entail a high phonology cost,  
211 which is absent in the first bird. As a tradeoff we imagine that birds may chain  
212 alignments locally rather than globally. Such chaining is for example suggested by the  
213 sequence requirement for correct acoustic models in white-crowned sparrows<sup>4</sup>.

214

#### 215 Experimental characterization of $\Delta$

216

217 Our experiments provide the following constraints  $C1$  to  $C4$  on error  
218 assignments in zebra finches:

219  $C1$ )  $ABC \rightarrow AC^+B$  and  $ABC \rightarrow A^+C^+B$  (imitation tasks 1 and 2): Birds  
220 perform the local change  $C \rightarrow C^+$ ; therefore, off-diagonal elements of  $\Delta$  can be  
221 nonzero (**no global alignment**).

222  $C2$ )  $ABC \rightarrow AB^-CB^+$  (imitation task 3): Birds do not chain locally: as target  
223 for  $B$ , they select either  $B^-$  (in context) or  $B^+$  (out of context). Because birds  
224 do not choose an interpolation between  $B^-$  and  $B^+$  as target, it follows that in  
225 each column of  $\Delta$  at most one assignment weight is nonzero, i.e.,  $\delta_{i,j} = \{0,1\}$

226 (**selection**); and  $\sum_i \delta_{i,j} \leq 1, \forall j$  (**winner-takes-all**).

227  $C3$ )  $ABC B^{+1} \rightarrow AB^{+2}CB^{-1}$  (imitation task 4): Birds make the changes  $B \rightarrow B^{-1}$   
228 and  $B^{+1} \rightarrow B^{+2}$  implying that assignments are **greedy**:  $\delta_{i,j} = 1$  in syllable-target  
229 pairs for which the local phonology error  $e(S_j, T_i)$  is minimal, and  $\delta_{i,j} = 0$   
230 otherwise.

231  $C4$ ) In all experiments, no two syllables or calls converge on the same target  
232 (**musical chairs**); therefore, in each row of  $\Delta$  exactly one assignment weight is  
233 nonzero, i.e.  $\sum_j \delta_{i,j} = 1, \forall i$ .

234

235

236

237

238

239 Linear assignment problem

240

241 In combination, our observations show that birds greedily choose a (binary)  
242 assignment matrix  $\Delta^*$  associated with minimal phonology error:

243 
$$\Delta^* = \arg \min_{\Delta} \sum_{i,j} \|S_j - T_i\| \delta_{i,j} \quad (5)$$

244 where  $e(S_j, T_i) = \|S_j - T_i\|$  is the absolute pitch difference between syllable  $j$   
245 and target  $i$ , and where  $\Delta$  is a  $m \times n$  permutation matrix with at most one nonzero  
246 entry per column and exactly one nonzero entry per row ( $\delta_{i,j} \in \{0,1\}$ ,  $\sum_i \delta_{i,j} \leq 1, \forall j$ ,  
247 and  $\sum_j \delta_{i,j} = 1, \forall i$ ). Equation (5) fully specifies the assignment matrix; what is  
248 particularly interesting is that the optimization in Equation (5) does not depend on the  
249 tradeoff constant  $c$ , implying absence of a tradeoff. The optimization in Equation (5)  
250 is known as the linear assignment problem which can be conveniently solved using  
251 for example the Hungarian method<sup>5</sup>.

252

253 In the context of natural language processing, the solution to Equation (5) (the  
254 minimum in Equation (5) rather than its argument) is also known as the word mover's  
255 distance<sup>6</sup> that represents the distance between two text documents. In that analogy,  
256  $\|S_j - T_i\|$  represents the distance between an individual word  $S_j$  in a source document  
257 and a word  $T_i$  in a target document. The assignments  $\delta_{i,j}$  (which do not have to be  
258 binary but can take arbitrary nonnegative values), represent the flows between words.  
259 These flows have to sum up to match the bag-of-words (vocabulary) representations  
260  $d_i$  and  $d'_j$  of the source and target documents,  $\sum_j \delta_{i,j} = d_i$  and  $\sum_i \delta_{i,j} = d'_j$ , in  
261 analogy to the musical chairs competition we find. The word mover's distance  
262 outperforms other approaches on many benchmark document categorization tasks<sup>6,7</sup>.

263

264 The fact that birds choose the assignment of minimal overall phonology error,  
265 irrespective of syntax, demonstrates a radical way of dealing with the intractability of  
266 the general assignment problem. Namely, rather than getting entangled with high  
267 complexity and large cognitive demand, birds decide to solve a much simpler  
268 tractable problem and do this remarkably well.

269

270 The surprising implications are that birds do not consider the cost of  
271 resequencing at all when correcting phonology errors. Phonology errors seem to be  
272 associated with a high cost, perhaps reflecting the amount of effort required to change  
273 syllable pitch. Counterintuitively, birds behave in this process as if there were no  
274 resequencing cost at all, despite the fact that this cost is seemingly very high, given  
275 that most birds try to re-sequence their syllable strings but only few succeed. Namely,  
276 we found that many birds do not reach the global performance error minimum in  
277 Equation (3) but get stuck somewhere on the way where some syntax errors but  
278 usually no phonology errors remain.

279

280 In summary, song learning is a modular, two-fold process. In a first process,  
281 birds choose assignments  $\Delta^*$  by solving a linear problem based on their vocal  
282 repertoire but not on their song sequence. In the second process, birds reduce

283 phonology errors defined by these correspondences and independently and more  
284 slowly, also reduce the resulting syntax error.

285

286 Sub-syllabic notes

287

288       What is the smallest song unit to which our formalism applies? We deliberately  
289 called  $S_j$  a song element and  $T_i$  a target element, implying these elements do not  
290 necessarily have to represent entire song syllables but could also represent sub-  
291 syllabic notes. In the following we discuss this possibility.

292

293       In our treatment of the song learning problem, we implicitly assumed that birds  
294 compute phonological error of a syllable by integrating over the errors in its  
295 constituent notes. Essentially, we assumed that birds compute the error of a syllable  
296 by globally aligning its notes with that of a template syllable. However, we have no  
297 evidence for this mini-version of global alignment. Thus, it remains to be explored  
298 whether birds can assign one of its syllable notes either to a note in a different syllable  
299 of the template or to a note in a different position within the same template syllable.

300

301       Although it will not be possible to resolve this issue without further  
302 experimenting, we imagine that our discovered assignment strategy cannot apply to  
303 ever smaller song units. Namely, at some point, there must be an overload to short-  
304 term memory arising from all these pairwise comparisons between song and template  
305 elements. It is therefore likely that assignment capabilities of zebra finches are limited  
306 to the syllable level and do not generalize to smaller song units below that level.

307

308

309 How to match syllable vocabulary using the expectation maximization (EM)  
 310 algorithm and Gaussian mixture models

311  
 312 Song learning can be considered a density estimation problem in which the  
 313 unknown parameters of the developing song syllables must be identified such that  
 314 good matches with the sensory targets are achieved. In a Gaussian mixture model, the  
 315 observable data points  $T_i$  (renditions of the  $i=1, \dots, m$  target syllables) are modeled  
 316 as a superposition of  $n$  Gaussian probability densities, where  $n$  is the number of  
 317 distinct song syllables and calls in the juvenile's repertoire. In the one dimensional  
 318 case of fitting the pitch parameter alone, the distribution of pitch  $X_j$  of syllable  $j$

319 ( $j=1, \dots, n$ ) is given by  $P(X_j) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x_j-s_j)^2}{2\sigma^2}}$ , where  $S_j$  is the mean pitch of  
 320 that syllable. It follows that the likelihood density that syllable  $j$  will produce the  
 321 pitch  $T_i$  of template  $i$  is given by

$$322 \quad P_{ij} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(T_i-S_j)^2}{2\sigma^2}}. \quad (6)$$

323 Here we assume that  $\sigma^2$  is the constant pitch variance of syllable  $j$ . The closer  
 324 the mean pitch  $S_j$  is to the target pitch  $T_i$ , the more likely the production of syllable  
 325  $j$  will match the target  $i$ .

326  
 327 The goal of the EM algorithm is to identify the set of mean syllable pitches  $S_j$   
 328 that maximize the total probability  $P$  (or its logarithm) of reproducing all the target  
 329 pitches  $T_i$ . The following function  $L$  is usually maximized by the EM algorithm:

330  $L = \sum_i \ln(P_i) = \sum_i \ln\left(\sum_j P_{ij} p_j\right)$ , where  $P_i = \sum_j P_{ij} p_j$  is the probability that target  
 331  $i$  is produced by any of the  $n$  syllables,  $p_j$  being the prior probability of singing  
 332 syllable  $j$ .

333 In the following, we assume that all syllables have identical prior probability,  
 334  $p_j = p$  (zebra finches sing motifs with linear syllable arrangement) and drop any  
 335 further mention of  $p$  because it is an irrelevant constant. Note that the EM algorithm  
 336 is an iterative algorithm aimed not at directly maximizing  $L$ , but rather its lower  
 337 bound  $L_L < L$ :

$$338 \quad L_L = \sum_{i,j} P_{j|i} \ln(P_{ij}), \quad (7)$$

339 where  $P_{j|i}$  is referred to as the posterior probability (that syllable  $j$  is assigned  
 340 to target  $i$ ). To avoid confusing the posterior probability  $P_{j|i}$  with the Gaussian  
 341 likelihood  $P_{ij}$ , it is always assumed that the index  $j$  refers to a syllable and index  $i$   
 342 refers to a target. The maximization of  $L_L$  in Equation (7) is usually done in two  
 343 steps, an E step in which the posterior probabilities are updated:

344 
$$P_{ji} = \frac{P_{ij}}{\sum_j P_{ij}}, \quad (8)$$

345 and an M step in which the mean pitches are updated according to

346 
$$S_j = \frac{\sum_i T_i P_{ji}}{\sum_i P_{ji}}. \quad (9)$$

347 Back-and-forth iteration of Equations (8) and (9) usually leads to convergence of  
348 the set of mean pitches  $S_j$  towards the set of targets  $T_i$ .

349

350

351 Birds' strategy compared with the EM algorithm

352

353 The EM algorithm operating on Gaussian mixture models just outlined exhibits  
354 several similarities with birds' strategy of minimizing performance error. Namely, if  
355 we place the Gaussian model in Equation (6) into the function  $L_L$  to be maximized in  
356 Equation (7), then we obtain the following expression:

357 
$$L_L = K - \frac{1}{2\sigma^2} \sum_{i,j} P_{ji} (T_i - S_j)^2, \quad (10)$$

358 with  $K = -m \ln(\sigma\sqrt{2\pi})$  being a constant without relevance for the

359 maximization (because  $\sigma$  is assumed to be constant). The maximization in Equation  
360 (7) is identical with the minimization in Equation (5), provided we interpret the  
361 posterior probabilities  $P_{ji}$  as assignment weights  $\delta_{i,j}$ .

362 In the EM algorithm, the posterior probabilities  $P_{ji}$  are not constrained to be  
363 binary variables that take values either zero or one. Nevertheless, the E and M steps in  
364 Equations (8) and (9) achieve to a good approximation the musical chairs competition  
365 we found.

366 To see this, consider the case in which for a given target there is only a single  
367 syllable with similar pitch ( $P_{ij}$  is large only for a single syllable  $j$ ). According to  
368 Equation (8),  $P_{ji}$  is close to 1 for that best matching syllable  $j$  and close to 0 for the  
369 other, non-matching syllables. This means that Equation (8) implements a soft  
370 competition among syllables ( $\sum_j \delta_{i,j} = 1$ ), which is an approximation of the musical  
371 chairs interactions among syllables we found.

372 In a similar way, the normalization in Equation (9) implements a soft winner-  
373 takes-all mechanism. Namely, if for a given syllable  $j$  one of the posterior  
374 probabilities (assignment weights)  $P_{ji}$  is large and the other very small, then by the  
375 weighted sum in Equation (9), that syllable's pitch is drawn towards the pitch of the  
376 assigned target.

377 To model the slow and gradual song development seen in birds, we simulated a  
378 finely discretized version of the EM algorithm in Equations (8) and (9). Because birds  
379 change pitch continuously and slowly unlike in Equations (8) and (9), we  
380 implemented a slow dynamical system in which we replaced the possibly large and  
381 discontinuous posterior probability and mean pitch changes in Equations (8) and (9)  
382 by gradual iterative processes (iterating over renditions  $t$ ):

383

384 1) We replaced the mean pitch  $S_j$  defined in Equation (9) by the iterative  
 385 variables

$$386 \quad S_j^{t+1} = S_j^t + \alpha \left( \sum_i T_i P_{ji} - S_j^t \sum_i P_{ji} \right), \quad (11)$$

387 where the index  $t$  labels the rendition number of the syllable, and  $\alpha$  is a small  
 388 integration rate. At a steady state  $S_j^t$  equals  $S_j$ .

389 2) We sampled the pitch  $X_j^t$  of the rendition  $t$  of syllable  $j$  according to the  
 390 Gaussian model

$$391 \quad P(X_j^t) = \frac{1}{\sigma \sqrt{2\pi}} e^{-\frac{(x_j^t - S_j^t)^2}{2\sigma^2}}. \quad (12)$$

392 3) We assume that birds can compute random samples and estimate their density,  
 393 but they cannot explicitly compute probabilities. We thus sampled the likelihoods  $P_{ij}$   
 394 in Equation (6) via random variables  $b_{ij}^t$  and their running averages  $n_{ij}^t$ . First, we  
 395 computed the instantaneous likelihood as a binary random variable  $b_{ij}^t$ :

$$396 \quad P(b_{ij}^t = 1) = \frac{1}{\sigma_b \sqrt{2\pi}} e^{-\frac{(r_j - x_j^t)^2}{2\sigma_b^2}} \quad (13)$$

397

398 and  $b_{ij}^t = 0$  otherwise with parameter  $\sigma_b$ . Second, we estimated the posterior

399 probabilities in Equation (8) according to  $P_{ji}^t = \frac{n_{ij}^t}{\varepsilon + \sum_i n_{ji}^t}$  with  $\varepsilon$  a small

400 regularization constant and  $n_{ij}^t$  a running average of  $b_{ij}^t$ :

$$401 \quad n_{ij}^{t+1} = (1 - \alpha_2) n_{ij}^t + \alpha_2 b_{ij}^t \quad (14)$$

402 with integration rate  $\alpha_2$ .

403 We simulated birds that produced three syllables ( $n = 3$ ) and had to match three  
 404 targets ( $m = 3$ ). Therefore, we iterated back and forth the above expressions for  $S_j^t$   
 405 and  $P_{ji}^t$ .

406

407 In simulations, we realized that the musical chairs competition is not well  
 408 captured by these equations: A call that did not match any target tended to converge  
 409 to a nearby target regardless whether the latter was occupied or not (Supplementary  
 410 Fig. 5a). To remedy this discrepancy, we hardened the musical chairs competition by  
 411 adding the constraint that at each sampled likelihood,  $b_{ij}$  could be 1 for at most one  
 412 vocalization, implying that matched targets could not attract any unassigned syllables  
 413 or calls (Supplementary Fig. 5b). We achieved this constraint by setting  $b_{ij}$  to zero for  
 414 all  $j$  whenever for a given target  $i$  two or more  $b_{ij}$ 's were sampled to be one (i.e.,  
 415 when  $\sum_j b_{ij} > 1$  we set  $b_{ik} = 0, \forall k$ ).

416

417

418

419

420 In summary, we iteratively simulated the system in 5 steps ( $N1$  to  $N5$ ):

421  $N1$ : sample pitch  $X_j^t$  according to Equation (12)

422  $N2$ : sample the instantaneous likelihood  $b_{ij}^t$  according to Equation (13)

423  $N3$ : enforce musical chairs:  $\forall i$ , if  $\sum_j b_{ij}^t > 1$  set  $b_{ik}^t = 0 \forall k$

424  $N4$ : update running average according to Equation (14)

425  $N5$ : update mean pitch according to Equation (11)

426 The syllable trajectories  $X_j^t$  resulting from running this system are shown in

427 Supplementary Fig. 5b-d.

428 The interesting property of these equations is that they in essence capture the  
429 observations without requiring any parameter fitting other than  $\sigma_b$ . The latter was set  
430 to exceed the pitch standard deviation  $\sigma$ . This allowed for medium size pitch shifts  
431 of two semitones. The integration rates  $\alpha$  and  $\alpha_2$  dictates the speed of pitch shifting,  
432 these parameters are set to yield smooth looking transitions.

433

434 In summary, the competition we find in birds is harder than that in the standard EM  
435 algorithm, in the sense that the EM algorithm brings all Gaussian models (syllables)  
436 to observables (targets), even if there are just 2 targets and 3 models, very unlike birds  
437 that bring only one syllable or call to each target, in a presumed attempt to limit used  
438 syllable resources. In a sense, birds are more efficient than the traditional EM  
439 algorithm, similarly to ongoing machine-learning approaches for restricting the  
440 effective number of model parameters to prevent overfitting, such as sparse priors<sup>8,9</sup>,  
441 Bayesian learning, and Dirichlet processes<sup>10,11</sup>. We believe that greedy and  
442 competitive error assignment during vocal learning illustrates the importance of  
443 minimizing used resources.

444 **Supplementary References**

445

- 446 1. Lipkind, D. *et al.* Stepwise acquisition of vocal combinatorial capacity in  
447 songbirds and human infants. *Nature* **498**, 104–8 (2013).
- 448 2. Koopmans, T. C. & Beckmann, M. J. Assignment Problems and the Location  
449 of Economic Activities Author ( s ): Tjalling C . Koopmans and Martin  
450 Beckmann. *Econometrica* **25**, 53–76 (1957).
- 451 3. Sahni, S. & Gonzalez, T. P-Complete Approximation Problems. *J. ACM* **23**,  
452 555–565 (1976).
- 453 4. Rose, G. J. *et al.* Species-typical songs in white-crowned sparrows tutored with  
454 only phrase pairs. *Nature* **432**, 753–8 (2004).
- 455 5. Kuhn, H. W. The Hungarian method for the assignment problem. *Nav. Res.*  
456 *Logist. Q.* **2**, 83–97 (1955).
- 457 6. Kusner, M. J., Sun, Y., Kolkin, N. I. & Weinberger, K. Q. From Word  
458 Embeddings To Document Distances. *Proc. 32nd Int. Conf. Mach. Learn.*  
459 (2015).
- 460 7. Huang, G., Guo, C., Kusner, M., Sun, Y. & Sha, F. Supervised Word Mover’s  
461 Distance. *Adv. Neural* (2016).
- 462 8. Olshausen, B. A. & Field, D. J. Sparse coding with an overcomplete basis set:  
463 A strategy employed by V1? *Vision Res.* **37**, 3311–3325 (1997).
- 464 9. Tipping, M. Sparse Bayesian learning and the relevance vector machine. *J.*  
465 *Mach. Learn. Res.* **1**, 211–244 (2001).
- 466 10. Neal, R. M. *Bayesian Learning for Neural Networks*. (Springer Science &  
467 Business Media, 2012).
- 468 11. Ferguson, T. A Bayesian analysis of some nonparametric problems. *Ann. Stat.*  
469 (1973).

470