

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

An Automated Approach to the Quantitation of Vocalizations and Vocal Learning in the Songbird.

David G. Mets^{1*} and Michael S. Brainard¹

¹Department of Physiology, University of California, San Francisco, CA 94158; Center for Integrative Neuroscience, University of California, San Francisco, CA 94158; Howard Hughes Medical Institute, University of California, San Francisco, CA 94158.

* Corresponding authors

E-mail: dmets@phy.ucsf.edu (DGM), msb@phy.ucsf.edu (MSB)

18 **Abstract**

19 Studies of learning mechanisms critically depend on the ability to accurately assess learning outcomes.
20 This assessment can be impeded by the often complex, multidimensional nature of behavior. We present a
21 novel, automated approach to evaluating imitative learning that is founded in information theory. Conceptually,
22 our approach estimates the amount of information present in a reference behavior that is absent from the learned
23 behavior. We validate our approach through examination of songbird vocalizations, complex learned behaviors
24 the study of which has provided many insights into sensory-motor learning in general and vocal learning in
25 particular. Historically, learning has been holistically assessed by human inspection or through comparison of
26 specific song features selected by experimenters (e.g. fundamental frequency, spectral entropy). In contrast, our
27 approach relies on statistical models that broadly capture the structure of each song, and then uses these models
28 to estimate the amount of information in the reference song but absent from the learned song. We show that our
29 information theoretic measure of song learning (contrast entropy) is well correlated with human evaluation of
30 song learning. We then expand the analysis beyond song learning and show that contrast entropy also detects the
31 typical song deterioration that occurs following deafening. More broadly, this approach potentially provides a
32 framework for assessing learning across a broad range of similarly structured behaviors.

33

34 **Author Summary**

35 Measuring learning outcomes is a critical objective of research into the neural, molecular, and behavioral
36 mechanisms that support learning. Demonstration that a given manipulation results in better or worse learning
37 outcomes requires an accurate and consistent measurement of learning quality. However, many behaviors (e.g.
38 speech, walking, and reading) are complex and multidimensional, confounding the assessment of learning. One
39 behavior subject to such confounds, vocal learning in Estrildid finches, has emerged as a vital model for sensory
40 motor learning broadly and human speech learning in particular. Here, we demonstrate a new approach, founded
41 in information theory, to the assessment of learning for complex high dimensional behaviors. Conceptually, we
42 determine the amount of information (across many dimensions) present in a reference behavior and then assess
43 how much of that information is present in the resultant learned behavior. We show that this measure provides

44 an accurate, holistic, and automated assessment of vocal learning in Estrildid finches. Potentially, this same
45 approach could be deployed to assess shared content in any multidimensional data, behavioral or otherwise.

46

47 **Introduction**

48 Songbird vocal learning shares many parallels with speech learning[1] and is a powerful, tractable model
49 system for elucidating neural and behavioral mechanisms underlying vocal control and vocal learning[2]. Birds,
50 like humans, learn vocalizations early in life through exposure to the vocalizations of an adult ‘tutor’ followed
51 by a period of practice that eventually results in typical adult vocalizations that require auditory feedback for
52 maintenance[1]. Song is composed of discrete units of sound (syllables) organized into higher order
53 sequences[3]. In the finch species examined here, a given bird’s song comprises about 5-10 categorically
54 distinct syllable types, with these distinct types defined by their unique spectro-temporal structure (Fig. 1A).
55 Hence, an individual bird’s song can be described as a specific set of categorically distinct syllable types (that
56 can be labeled ‘A’, ‘B’, ‘C’ and so on).

57 Qualitatively, learning (and failure to learn) can occur in different ways (Fig. 1B-C)[4]. For example,
58 juvenile ‘tutees’ could learn to produce all distinct syllable types present in an adult ‘tutor’ song, but the spectral
59 content of the syllables might be imperfect or noisy (Fig. 1Ci-Cii), while other tutees might completely fail to
60 learn some syllables (Fig. 1Ciii), and still others might improvise new syllables (Fig. 1Civ).

61 Because of these complexities, many studies have relied on human evaluation of song similarity and
62 learning [5–8]. Indeed, human scorers can provide a useful ‘holistic’ assessment of similarity between complex
63 behaviors, such as songs, which integrates across many stimulus dimensions. However, human scoring suffers
64 from several problems including 1) it requires scorers to be trained on species-specific vocalizations, and
65 analysis of different vocalization types often requires new training, 2) correspondingly, scores are potentially
66 inconsistent over time and across different evaluators, and 3) human scoring is labor intensive and does not
67 readily scale to the size of relevant datasets, which, in the case of bird song, can include many individuals and
68 thousands of vocalizations per day.

69 More recent attempts at quantification of song similarity have focused on assessing learning based on

70 specific, reliably and automatically measured features of song (e.g. the fundamental frequency of a specific
71 syllable, or song entropy)[9–12]. In these approaches, selected samples of songs are decomposed into sets of
72 feature values[13] and song similarity is then evaluated as the similarity between weighted feature sets
73 associated with those samples. These approaches can overcome some of the inter-evaluator variability
74 associated with human scoring, and additionally enable useful assessment of how specifically analyzed features
75 such as the ‘pitch’ or ‘noisiness’ of syllables differ across songs and conditions. However, due to the reductionist
76 nature of the extracted features, even measures incorporating many such features will potentially fail to capture
77 biologically important song information. Additionally, as with direct human assessment of song similarity, these
78 approaches often require significant human intervention for the selection of which samples of song material to
79 analyze and which features to weigh in assessing similarity.

80 For these reasons, we were interested in developing an approach to scoring the similarity between pairs
81 of songs (and other complex stimuli) that could reproduce some of the human capacity for holistically
82 integrating across complex stimulus dimensions but that was also automatic, reproducible and efficiently
83 deployed across large data sets. Our approach has four main components. First, we start with a representation of
84 song that is the set of power spectral densities (PSDs) associated with each discrete syllable. This representation
85 retains much of the complexity of the spectro-temporal information present in each song without specifically
86 extracting features such as the pitch or entropy of each syllable. Second, we represent the spectro-temporal
87 structure of each song by transforming each of a large number of individual syllables drawn from the song (as
88 represented by their corresponding PSDs) into a ‘syllable similarity’ space, in which each syllable (PSD) is
89 represented by its similarity to a large basis set of other syllables (PSDs). Intuitively, in this high dimensional
90 space, different iterations of a given syllable type (e.g. different iterations of the syllable ‘A’) will cluster near
91 each other, because they share a similar PSD. Each song is therefore associated with regions of high density
92 within the syllable similarity space, with each high density region corresponding approximately to a distinct
93 syllable type. Third, we then model each song by characterizing the high-density regions within the syllable
94 similarity space (corresponding to regions in which there is clustering of PSDs). We use a Gaussian Mixture
95 Model (GMM) to fit these regions of high density, in which the means and variances are fit to the data for each

96 song using a standard expectation-maximization algorithm [14,15], and the number of Gaussian mixture
97 components is determined using Bayesian Information Criteria (BIC) [16], a measure of model fit that is
98 penalized for increasing model complexity. Lastly, since the GMM describing the distribution of spectro-
99 temporal structure for syllables from each song is a statistical model, any differences in vocal structure (e.g.
100 between tutors and tutees, or before and after onset of manipulations) can be assessed according to a principled
101 information theoretic measure of the differences between the models fit to those songs. We describe the
102 difference in structure between songs as the difference between two cross entropies, here referred to as the
103 contrast entropy (CE, see Methods). CE quantifies the amount of information present in a reference song (e.g. a
104 tutor song, or baseline song before a manipulation such as deafening) that is not present in a comparison song
105 (e.g. the learned song of a tutee, or songs that are produced following deafening).

106 In this paper, we validate our approach by characterizing the similarities between pairs of songs for two
107 conditions. First, we assess song learning by comparing adult tutor songs and juvenile tutee songs in the
108 Bengalese finch (*Lonchura striata domestica*), a species with variability in both spectral content and syntax. We
109 show that CE provides a measure of the quality of song learning in Bengalese finches that is well correlated with
110 scores provided by human experts. Second, we assess song deterioration following deafening by comparing
111 baseline songs produced by adult Zebra finches (*Taeniopygia guttata guttata*) with songs of the same individuals
112 at varying times following deafening. We show that CE detects and quantifies the characteristic song
113 deterioration that follows deafening [7]. Together, these results demonstrate an unsupervised approach to the
114 quantification of both song learning and song deterioration. Such an approach allows holistic and reproducible,
115 high-resolution tracking of song similarity across both large populations of individuals and long periods of time,
116 and has potential relevance outside of birdsong as it could be applied to automatic analysis of human speech and
117 other complex multidimensional data.

118

119 **Results**

120

121 We present an automated method for assessing song learning by calculating the amount of spectral

122 information which is present in the song of the tutor bird, but absent from the song of the tutee. To accomplish
123 this we construct statistical models from the song spectral content of both reference song (tutor) and comparison
124 song (tutee) then estimate the amount of reference song information accounted for by the tutor model but not the
125 tutee model. Below, we first describe how the statistical model for each song is constructed and how the
126 statistical models for two songs are compared to provide an estimate of song similarity, the contrast entropy
127 (CE). We then show that CE correlates well with human scores for data sets that characterize song learning and
128 demonstrate that CE also detects the typical deterioration of song following deafening. The descriptions below
129 elaborate a specific instantiation of our overall approach. We follow this with an examination of the robustness
130 of CE measures across a range of different possible instantiations, and a discussion of how our approach could
131 be modified or extended to address related issues of similarity in other domains.

132

133 **Overview of statistical model assembly**

134 The assembly of statistical models representing song spectral data is schematized in Figure 2. Starting
135 with song data from a given bird (Fig. 2A), we identify individual syllables as continuous amplitude traces above
136 a threshold. For each syllable, we calculate the power-spectral density (PSD; an estimate of acoustic power at a
137 set of specific frequency values) using Welch's method[17] (Fig. 2B). We then calculate the similarity between
138 each PSD and a “basis set” of PSDs (see Methods), where the basis set is a random draw from the set of PSDs
139 being modeled. These computed similarities create an $N \times M$ matrix of PSD similarities (Fig. 2C) where N is the
140 number of PSDs being analyzed ('target' PSDs) and M is the number of PSDs in the basis set ('basis' PSDs). For
141 the analyses presented below, we use a value for ‘ N ’ of 3000 PSDs drawn from the song to be modeled, and a
142 value for ‘ M ’ of 50 PSDs to form the basis set for construction of the syllable similarity matrix.

143 Translation of the raw syllable data into a matrix of syllable-similarities has the important consequence
144 that it naturally and automatically clusters syllables with similar PSDs into regions of high density (Fig. 2D, E).
145 Each of these regions corresponds approximately to groups of syllables that would be identified by a human
146 observer as belonging to a specific type. That is because each instance of a syllables type has a similar PSD, and
147 therefore each of these instances will have a similar pattern of distances from each of the elements of the basis

148 set. This transformation therefore results in a representation of song in syllable similarity space in which
149 thousands of exemplar syllables are clustered into a small number of high density regions, corresponding
150 approximately to the numbers and identities of categorically distinct syllable types in the original song.

151 We then fit a series of GMMs[14] (using expectation maximization[15]) to the distribution of similarity
152 values. Each successive model in the series has an incrementally higher number of mixture components. The
153 model with the best fit to the data based on the lowest Bayesian Information Criterion[16] is then used to
154 represent the song of that bird. Conceptually, the number of Gaussian mixture components in the song model
155 corresponds approximately to the number of discrete syllable types present in the song.

156 The distribution of similarity scores reveals GMMs to be appropriate and effective for modeling these
157 data. To illustrate this, Figure 2D depicts the similarities between three discrete target PSDs (Fig. 2D, blue
158 purple and yellow dots) and two basis PSDs (out of a total of 50 basis PSDs). Target PSD identity was assigned
159 as the maximum posterior probability over Gaussian mixture (see Methods) and is depicted by data color.
160 Examination of the joint distribution of similarities between target PSDs and basis PSDs (Fig. 2D) reveals three
161 clusters that are, even in just two dimensions, well separated. Importantly, the separation between clusters is
162 reliant not only on the mean and variance in any one dimension but also on the covariance structure between
163 values in the two exemplar dimensions; all three PSD categories are distributed elliptically. Examination of the
164 marginal distributions (Fig. 2D, top and right) in each single dimension exemplifies the Gaussian nature of the
165 data. Sorting the NxM matrix of similarity data by PSD identity (Fig. 2E) reveals the underlying ordered
166 structure of these data; the PSDs in each group (Fig. 2E, indicated by color) share stereotypical PSD similarity
167 structure. Spectrograms reveal good consistency of syllable identity (as indicated by corresponding PSD
168 identity) within GMM classified groups (Fig. 2E, right).

169

170 **Assessment of the quality of learning**

171 Because we capture song spectral content as statistical models, we can compare the spectral content of
172 two songs using a principled information theoretic measure. In Figure 3 we outline the calculation of this
173 measure, here called contrast entropy (CE). CE estimates the amount of information in the reference song that is

174 not present in the comparison song (Fig. 3A). In this example, some, but not all, syllables from the reference
175 song are well-represented in the comparison song. GMMs are fit independently to the two songs, as described
176 above, except a single basis set of PSDs, drawn from the reference song, is used for modeling both songs, so that
177 the syllables from both songs are transformed into the same syllable similarity space. In Figure 3B a single
178 dimension of these probability densities are shown for simplicity. We then calculate the mean-log likelihood of a
179 set of reference song data (Fig. 3C, light blue histogram). This specific data set was not used to estimate GMM
180 parameter values (held-out data). This provides an estimate of the cross entropy between the sampled
181 distribution of the reference song and the model for the reference song. We similarly estimate the cross entropy
182 between the sampled reference song distribution and the model for the comparison song. Here we define the
183 difference between these two cross entropies as the contrast entropy (CE).

184 An intuition for the relevance of CE to the content of song is provided by examination of Figure 3C.
185 Here, represented in a single dimension, the held-out data from the reference song fit the probability distribution
186 from the reference model well; very little of the data (histogram) fall outside the model (blue line). The same
187 data are less well matched to the probability distribution for the comparison song; some data fall outside the
188 model (Fig. 3C, purple arrow). Thus, the log-likelihood of these data given the reference song model is higher
189 than that of the same data given the comparison song model and the magnitude of this difference determines the
190 CE. However, the likelihood of the data is not influenced by portions of the distribution of the comparison song
191 model that are not occupied (red arrow) by data from the reference song. Hence, contrast entropy will be low if
192 the tutee (comparison) learned all elements of the tutor (reference) song well, but will be unaffected by tutee
193 song content that is not present in the tutor song (e.g. song content innovated by the tutee; we discuss later how
194 our approach can be extended to quantify innovation by the tutee).

195

196 **CE closely parallels human assessment of learning outcomes**

197 We evaluated CE as a holistic measure of song similarity through comparison of CE to human scoring
198 (Fig. 4). We used both CE and human scoring to assess learning in five cohorts of Bengalese finches. Each
199 cohort was tutored with a different song (cohort tutor song). Figure 4A shows an example of the cohort tutor

200 song for one group (top) and 5 tutee songs that illustrate a broad range in the quality of copying. For each
201 cohort, four expert human evaluators independently estimated the similarity between each tutee's learned song
202 and the tutor song on a scale of 0-4, with 0 being most similar.

203 Across all five groups, CE and human scores were well correlated. Figure 4A shows the average human
204 and CE scores assigned to 5 example tutee songs from one cohort. The rank ordering of similarities for these
205 five songs relative to the tutor song was the same for CE and human scores; the learned songs are displayed from
206 top to bottom in order of decreasing similarity to the tutor song by both measures. Figure 4B shows the
207 correlation between CE and average human scores for 65 birds from the 5 cohorts ($r=0.72$, $p<0.01$). When
208 calculated for each cohort individually, the median CE-human correlation was high (Fig. 4C, human-computer
209 correlation). We compared these CE-human correlations to Human-Human correlations. For each cohort of
210 birds, the scores of each evaluator were correlated with the average scores provided by the other evaluators (Fig.
211 4C, human-human correlation). The correlation between CE and human scores was comparable to the
212 correlation between different humans' scores. Together, these results indicate that CE provides a holistic and
213 automated assessment of song learning that closely parallels human evaluation.

214 As an additional reference for 'poor learning', we also computed CE relative to each cohort tutor song
215 for two additional groups of birds: 'isolate birds' that were raised without exposure to any tutor, and 'unrelated
216 birds' that were raised with a tutor different from any of the five cohort tutors. CE indicates information from
217 the reference song that is missing from comparison songs. Consistent with this, CE for isolate songs that contain
218 atypical vocalizations was higher than that for songs from normally tutored birds (Fig. 4A, example 'isolate
219 song' and Fig 4D, summary comparisons, $p<0.01$, Wilcoxon rank test). Similarly, CE for songs from birds that
220 copied an unrelated tutor was also higher than that for songs from birds that learned from cohort tutor (Fig. 4A,
221 example 'unrelated song' and Fig. 4D, summary comparisons; $p<0.01$, Wilcoxon rank test).

222

223 **Quantification of changes in the spectral content of song due to deafening**

224 Many studies assess changes in song structure following various manipulations. One such manipulation,
225 deafening, produces gradual deterioration of song structure and has been used extensively to provide insights

226 into the mechanisms of song learning[6,7,18]. However, as with learning, quantitative assessment of song
227 deterioration following deafening has often relied on human inspection[6,18]. To determine whether CE
228 captures this deterioration we examined the songs of seven Zebra finches before and after deafening. We used
229 CE to evaluate changes in song spectral content between baseline reference songs (before deafening) and
230 comparison songs from the same birds two, four, six, and eight weeks post deafening. Figure 5A illustrates
231 spectrograms from baseline songs and corresponding portions of post-deafening songs for 3 birds that
232 qualitatively exhibited small (green), medium (yellow) and large (blue) changes to syllable spectral content.
233 Figure 5B shows post deafening CE trajectories for 9 birds. These data reveal a gradual and continuous loss of
234 song information over time (as quantified by CE), and demonstrate our approach as a sensitive method for
235 evaluating changes in song following manipulations such as deafening.

236

237 **Robustness of similarity measures to parameter choices.**

238 Our approach is intended to provide a measure of song similarity that requires little in the way of user
239 intervention and selection of parameters. Correspondingly, all of the foregoing analyses were based on a specific
240 instantiation of our approach using fixed values for parameters that could in principle be set by a user. Here we
241 consider how different choices of parameter values affect similarity measures and demonstrate that CE measures
242 are indeed very robust across a brand range of values. The specific parameters that we consider are 1) the
243 number of ‘target’ syllables drawn from both the reference song and the comparison song, 2) the number of
244 ‘basis’ syllables used for the PSD similarity basis set, and 3) the number of mixture components in the GMM
245 used to model the structure of each song. In addition to these numerical choices, our approach as described
246 above uses the PSD for each syllable as a representation of the spectro-temporal complexity of song. We
247 therefore also consider and discuss below how different choices of song representation could affect similarity
248 measures or extend our approach to capture other aspects of song structure.

249 There are two numerical values important to CE that are necessarily under experimenter control: the
250 number of target syllables from the song to be modeled (N) and the number of syllables in the basis set (M). To
251 determine appropriate values, we conducted numerical titrations of both the number of target syllables in the

252 input data set (Fig. 6a and S1) and the number of syllables in the basis set (Fig. 6b and S2). For each of 44 birds,
253 CE (relative to the corresponding tutor song) was calculated using 100, 200, 500, 1000, 2000, and 3000 target
254 syllables from the songs that were being modeled (both reference and comparison songs). For each number of
255 target syllables, we computed CE values for the 44 birds and correlated these CE values with CE values
256 determined with 3000 syllables. Not surprisingly, CE values with little input data showed substantial deviation
257 from the 3000 syllable CE values (Fig. 6a and S1A). However, for target syllable numbers above ~500, CE-CE
258 correlations approached an asymptote, indicating little change to CE similarity measures above this value (Fig.
259 6a and S1B-D). We therefore used a fixed value of 3000 target syllables to model song structure throughout our
260 analysis. For the same set of birds, CE values were calculated using 5, 10, 20, 40, 80, and 160 basis syllables.
261 CE-CE correlations were calculated relative to CE values derived using the 160 PSD basis set. CE-CE
262 correlations approached an asymptote for basis set numbers above ~25 (Fig. 6b and S2). For computational
263 efficiency and to constrain the number of free parameters in our GMMs, we therefore used a basis set of 50
264 syllables throughout the study. These data indicate that CE measures of song similarity are robust to changes in
265 numbers of target and basis syllables above threshold minimum values and therefore that our approach can be
266 deployed effectively with fixed values of these parameters that do not require human tuning.

267 The number of Gaussian mixtures required to model the spectral complexity of a given song was
268 selected automatically using Bayesian Information Criteria. However, the number of Gaussian mixtures can in
269 principle be set to different values, for example in cases where the experimenter has an independent basis for
270 modeling a song with a specific number of discrete syllable types. We therefore evaluated the robustness of CE
271 similarity measures to variation in the number of Gaussian mixture components. Specifically, we calculated the
272 squared error between CE values when calculated using the number of mixture components determined by BIC
273 (nBIC) and CE values calculated using a series of other models in which the number of mixture components
274 ranged from nBIC-4 through nBIC+4 (Fig. 6C and S3). These values were used in both the reference model and
275 the comparison model. CE was very robust to variation in the number of mixture components; squared errors for
276 all CE comparisons were above 0.96 (Fig. 6C, n=44, $p < 0.001$ for all correlations).

277 Throughout the analyses described above, we used a single PSD to capture the spectro-temporal content

278 of a given syllable. At each frequency represented, this single PSD encodes acoustic power averaged across the
279 duration of the entire syllable (Fig. 2B) effectively capturing syllable spectral content collapsed across time.
280 Accordingly, time dependent spectral information is not captured by these representations and this missing
281 information may influence CE values. We therefore compared CE calculated using single PSDs per syllable
282 with a series of CE values calculated using multiple PSDs per syllable where each syllable was divided into
283 equal duration blocks (2, 5, or 10 blocks per syllable) and PSDs were calculated for each block (Fig. 6D and S4).
284 CE calculated with single PSD syllable representations was well correlated with CE calculated with 10 PSDs per
285 syllable (Fig. 6D and S4A; $n=44$, $r^2=0.78$). This strong correlation suggests that much of the spectro-temporal
286 information in a syllable is captured in single PSD representations. For purposes of computational efficiency, we
287 therefore use a single PSD to represent syllable spectral content for song modeling and calculation of CE.

288

289 **Syllable identity assignments provided by GMMs are well correlated with human assignments**

290 Our contrast entropy measure is intended primarily to provide a holistic measure of song similarity
291 between a reference song and comparison song. As part of this process we model the structure of each song
292 using a GMM in which an intuition is that each fit Gaussian corresponds approximately to what a human
293 observer would label as a single categorically defined syllable type. The assignment of syllables to specific
294 Gaussian mixtures is not required for the computation of CE, which relies solely on the models fit to the
295 distribution of (unlabeled) syllable similarity data. However, it is of potential interest to know how effectively
296 assignment of syllables to different Gaussian mixtures results in a categorization of syllables by type that
297 matches human labeling of syllables. Such automatic labeling of syllables has potential utility in objectively
298 determining the number of distinct syllable types in a bird's song repertoire and facilitating the assignment of
299 labels corresponding to these types to large amounts of song data.

300 Here, we explicitly examine the correspondence between syllable labels assigned using the GMM for
301 each song to labels assigned by expert human evaluators. For each of 90 birds, syllables were labeled by
302 determining the maximum posterior probability for assignment of syllable identity under the GMM fitted to
303 songs of that bird (see Methods). These classifications were then compared with classifications provided by

304 human inspection. Figure 7A illustrates labeled syllable categories for songs from two example birds. For some
305 birds (e.g. upper panel) there was perfect concordance between human assigned (black) and GMM assigned
306 (red) labels (categories). However, for most birds there were some discrepancies between human and GMM
307 assigned labels (e.g. lower panel, gray box). To determine the accuracy of GMM based classifications, for 90
308 birds, all assignments were inspected by an expert human observer and the classification of each was determined
309 to be accurate or inaccurate (Fig. 7B). Overall, 50% of birds had 96% or better correspondence between human
310 and GMM assignments, while 80% of birds had 93% or better correspondence (Fig. 7C). This correlation
311 between human and GMM based syllable classification indicates that much of the complex information
312 subjectively used by humans to classify song syllables is incorporated into the GMM models that were used to
313 model song structure.

314 In this analysis, the spectral structure of each syllable was modeled using only a single PSD computed
315 from the entire syllable (Fig. 2B). To ask whether a richer spectro-temporal representation of each syllable
316 would increase concordance between human and GMM assignments, we built GMM models as above, but with
317 10 separate PSDs, evenly spaced over the duration of each syllable, used as an input representation for each
318 syllable. We previously found that this richer syllable representation had little effect on CE measures of song
319 similarity (Fig. 6D). In contrast, for the specific assignment of syllable labels, this richer representation resulted
320 in significantly improved correspondence between human assignments and GMM assignments of syllable labels
321 (Fig. 7D).

322

323 **Discussion**

324 We demonstrate an approach for analysis of song and song learning that is computationally efficient and
325 automated. We use syllable spectral content to assemble statistical models for song that can then be used to
326 estimate the amount of spectral information present in one song but absent in another. Our measure of song
327 similarity, the contrast entropy (CE), is well correlated with holistic song similarity scores provided by expert
328 human evaluators while providing several critical advantages. CE is automatically computed and thus consistent
329 given the same data, where human assessment is less reliable both across individuals and over time. Because CE

330 is automatically and efficiently computed, we can analyze large amounts of data (many thousands of songs)
331 facilitating dense analysis of learning across time and, for any given comparison, incorporating much more song
332 data in assessments of learning than can be accomplished by a human evaluator. Importantly, this also obviates
333 the need for selection of specific samples of song for comparison, and results in measurements of similarity that
334 neither require, nor are potentially biased by, human intervention in selection of representative samples of song
335 for analysis. Human evaluators of song learning also require species-specific training to ensure reliability and
336 increase consistency across individuals. We show that our approach can be applied, with no modification, to two
337 different songbird species vocalizations, indicating that this approach can be readily extend to analyze other
338 vocalizations and other complex but similarly structured data.

339 CE is asymmetrical in that it estimates the amount of spectral content in a reference song that is missing
340 from a comparison song. This asymmetry can be exploited to address distinct conceptual questions contingent
341 on the reference-comparison relationship. In the case of birdsong, when the reference is tutor song, and the
342 comparison is tutee song, the measure indicates how much information from the tutor song was not learned by
343 the tutee. Reciprocally, if the reference is tutee song, and the comparison is tutor song, CE indicates how much
344 information in the tutee song did not come from the tutor, providing an estimate of “innovation”.

345 We specifically focused on assessing the learning of song spectral content, but the general framework of
346 comparing the shared information between statistical models can be extended (or restricted) to different
347 categories of song information by changing the statistical descriptions of song. For example, an analysis might
348 focus on the means of syllables but not the rendition-to-rendition variation. In this case, syllable variances
349 would be excluded from cross entropy estimation. Alternatively, the model could be extended to include, not
350 only spectral content, but also syllable transition information using a Hidden Markov Model (HMM) with
351 Gaussian emissions. HMMs have been effectively used in the past to model song transition structure with
352 human assigned syllable identities[19]. If song structure was modeled using HMMs with Gaussian emissions,
353 CE would indicate discrepancies in spectral content as well as syllable ordering. Hence, our general approach to
354 the evaluation of learning can be applied to any aspect of song that can be incorporated into a statistical model.

355 Beyond song learning, our approach allows fitting GMMs to sparse, high dimensional data. GMMs

356 have been difficult to fit to high dimensional, sparse data, (like song PSDs) partially because standard
357 marginalization based dimensionality reduction approaches (e.g. principle components analysis) remove
358 covariance structure which is potentially critical to fitting accurate GMMs[20]. Here we reduce the
359 dimensionality and the sparseness of our data via calculation of syllable-syllable similarities. Our results
360 indicate that this intuitive approach allows modeling high dimensional data sets within a framework that
361 facilitates quantitative and meaningful comparisons. Similarity matrices are already used in the context of
362 spectral[21] and hierarchical[22] clustering, but neither approach provides a statistical description (provided by
363 GMMs) of data and, thus, cannot be easily used for information theoretic calculations. Our approach may have
364 broad application to high dimensional problems where statistical descriptions can be leveraged for more accurate
365 classification or where information theoretic calculations are desired.

366

367 **Methods**

368 **Song recordings**

369 For audio recording, animals were single housed in sound isolation chambers (Acoustic Systems).
370 Songs were recorded digitally at a sampling frequency of 32 kHz and a bit depth of 16 then stored uncompressed
371 using custom Python or LabView (National Instruments) software. Recording microphones were placed in a
372 fixed position at the top of the cage housing the bird. Prior to further analysis, all songs were digitally high pass
373 filtered at 500 Hz using a digitally implemented elliptical infinite impulse response filter with a passband edge
374 frequency of 0.04 radians.

375

376 **Syllable segmentation**

377 Discrete units of sound separated by silence (syllables) were identified based on amplitude. First an
378 “amplitude envelope” was created by rectifying the song waveform then smoothing the waveform through
379 convolution with an 8 ms square wave. This amplitude trace was then used, through thresholding, to identify
380 periods of vocalization. To automatically identify a threshold capable of separating vocalizations from silence,
381 we used Otsu's method[23]. Briefly, Otsu's method is an exhaustive search to identify a threshold that minimizes

382 the shared variance between data above threshold and data below threshold. Once the threshold is established,
383 “objects” are identified as contiguous regions of the amplitude envelope over threshold. To eliminate short and
384 spurious threshold crossing that can occur at the edge of syllables where syllable amplitude is low, any objects
385 separated by a gap of 5 ms or less are merged, and then any objects shorter than 10 ms are eliminated. The
386 onsets and offsets of each object are padded by 3 ms and then used to segment audio data from the original
387 filtered waveform. We refer to each of these segments of audio data as syllables.

388

389 **Power spectral density estimation**

390 To estimate frequency information of syllables while removing temporal information, we calculated the
391 power spectral density for each syllable at 2048 frequencies using Welch's method[17]. Briefly, the PSD was
392 computed via FFT for successive 4096 sample windows (at 32kHz), each overlapping by 256 samples. These
393 were then averaged over the duration of the syllable, and the power in the frequency range 600 Hz – 1600 Hz
394 (sampled at 1970 points) was used as the PSD for the syllable.

395

396 **Similarity matrix**

397 Instead of clustering syllables on their PSD values (sampled at 1970 points), we transformed each PSD
398 into a syllable similarity representation. For each of M syllables, the Euclidean squared distance between the
399 PSD of that syllable and the PSDs of a basis set of N reference syllables was calculated creating an MxN
400 distance matrix, D , in which $D_{ij} = \|p_i - q_j\|^2$ where p is the vector of M syllable PSDs and q is the vector of basis
401 set syllable PSDs. We then calculated A, a similarity matrix, where $A_{ij} = 1 / (D_{ij} / \max(D))$. For each bird in our
402 main analysis presented in results, an M=3000 data syllables and N=50 reference basis syllables were used.

403

404

405 **Gaussian mixture model and parameter estimation**

406 We model each song as a Gaussian mixture model fit to the distribution of syllables in the syllable
407 similarity space. These GMMs[14] are defined as:

$$p(x|\theta) = \sum_{k=1}^K \pi_k N(x|\mu_k, \Sigma_k) \quad (1)$$

408

409 Where K is the number of mixing components, π_k is the mixing weight, μ_k is the vector of means, and Σ_k is the
410 full covariance matrix for component k . The values of $\theta = \{\pi_k, \mu_k, \Sigma_k\}$ are the parameters. The values of μ_k and Σ_k
411 are initialized using the K-means algorithm and then estimated through standard expectation-maximization[15].
412 The value of π_k is $1/k$ for all k . The values of x are the observed data. For this work, the implementation of
413 expectation-maximization for GMMs in the scikit-learn software package[24] was used.

414

415 **Model estimation.**

416 We conducted model selection to identify the number of Gaussian mixtures needed to describe a song.
417 For each bird, we fit, as described, a series of GMMs to song data. The set of models had increasing numbers of
418 Gaussian components (K), ranging from 2-20. For each model we calculated a three-fold cross validated
419 Bayesian information criterion (BIC)[16], a measure of model fit that is penalized for increasing model
420 complexity. As the number of Gaussian components increase, the BIC decreases to a minimum value and then,
421 as the number of Gaussian components proceeds beyond optimal, the BIC increases again. The number of
422 Gaussian components in the model with the lowest BIC value was used to model each song's structure for
423 purposes of song similarity comparisons.

424

425 **Information theoretic measurement of song similarity**

426 In order to quantify the similarity of two songs, we follow the procedure described above to derive two
427 song models, one from a reference song (i.e. the song of an adult tutor, or baseline song of a bird prior to a
428 manipulation such as deafening) and one from a comparison song (i.e. the song of a tutee, or song of a bird
429 following a manipulation such as deafening). We then calculate how much worse (in bits of information) the
430 model based on the comparison song is than the model based on the reference song at explaining held out data
431 from the reference song. Mathematically, this is a difference between two cross-entropies (or Kullback-Leibler
432 divergence). Here we refer to this quantity as contrast entropy, defined as follows:

433
$$H_c(R||P,Q) = H(R,P) - H(R,Q) \quad (4)$$

434 Where H_c is the contrast entropy, P is the density function for the model fit to the reference data set, Q is the
435 density function fit the comparison data set, R is the true (but unknown) density function of the reference data.
436 $H(R,P)$ is the cross entropy between P and R , and $H(R,Q)$ is the cross entropy between Q and R . H_c is calculated
437 as:

438
$$H_c(R||P,Q) = -\sum_{n=1}^N \frac{1}{N} \log_2 p(r_n) + \sum_{n=1}^N \frac{1}{N} \log_2 q(r_n) \quad (5)$$

439 where N is the size of a set of syllable similarity data (r) from the reference song, $\log_2 p(r_n)$ is the log likelihood
440 of that data under the GMM (p) fit to a different set of reference data, and $\log_2 q(r_n)$ is the log likelihood of that
441 data under the GMM (p) fit to data from the comparison song. Intuitively, H_c is an estimate of how much worse
442 the comparison (e.g. tutee) model is at explaining the true data from the reference song (e.g. tutor) than a model
443 based on different (held out) data from the reference song. Alternatively, H_c can be viewed as an estimate of the
444 information present in the reference song that is absent in the comparison song.

445

446 **Syllable classification**

447 For each syllable, the GMM based syllable classification was calculated as the maximum posterior
448 probability over the K syllable types defined by each of the Gaussian mixtures in the fit model. Calculated as:

449
$$z_n = \underset{k}{\operatorname{argmax}} p(k|x_n, \theta) \quad (2)$$

450
$$= \underset{k}{\operatorname{argmax}} \frac{\pi_k N(x_n|\mu_k, \Sigma_j)}{\sum_j \pi_j N(x_n|\mu_j, \Sigma_j)} \quad (3)$$

451 Where z_n is the assigned syllable type.

452

453 **Human scoring**

454 Contrast Entropy estimates of song similarity created by our method were compared with assessments

455 made by humans experienced in song analysis. We considered 5 groups of songs, corresponding to 5 nests in our
456 colony. Each group consisted of a tutor reference song (the adult male breeder in the nest) and multiple tutee
457 comparison songs (between 8 and 20 tutee songs per group, corresponding to juveniles that were hatched and
458 raised to independence in the nest of the adult male tutor). These groups were chosen such that each of the tutor
459 songs were qualitatively distinct, and each group of tutees expressed a broad range in the quality of tutor song
460 copying. Human judges that had extensive experience in song analysis were presented with spectrograms
461 (frequency range from 500-10000 Hz) representing four samples of tutor song and four sample of each tutee
462 song. Four seconds of each of the four sample songs were presented on the same time scale in a single
463 representation. Human judges assigned a score between zero (high similarity to tutor), and four (low similarity to
464 tutor), to each tutor-tutee song pair.

465

466 **Deafened birds**

467 Deafening data were presented previously in Kojima et. al. 2013[24]. For each of seven birds, song
468 was recorded before deafening and at two, four, six, and eight weeks post deafening. In each case the amount of
469 information loss at any time post deafening was taken with reference to the pre-deafening, baseline song.

470

471 **Bibliography**

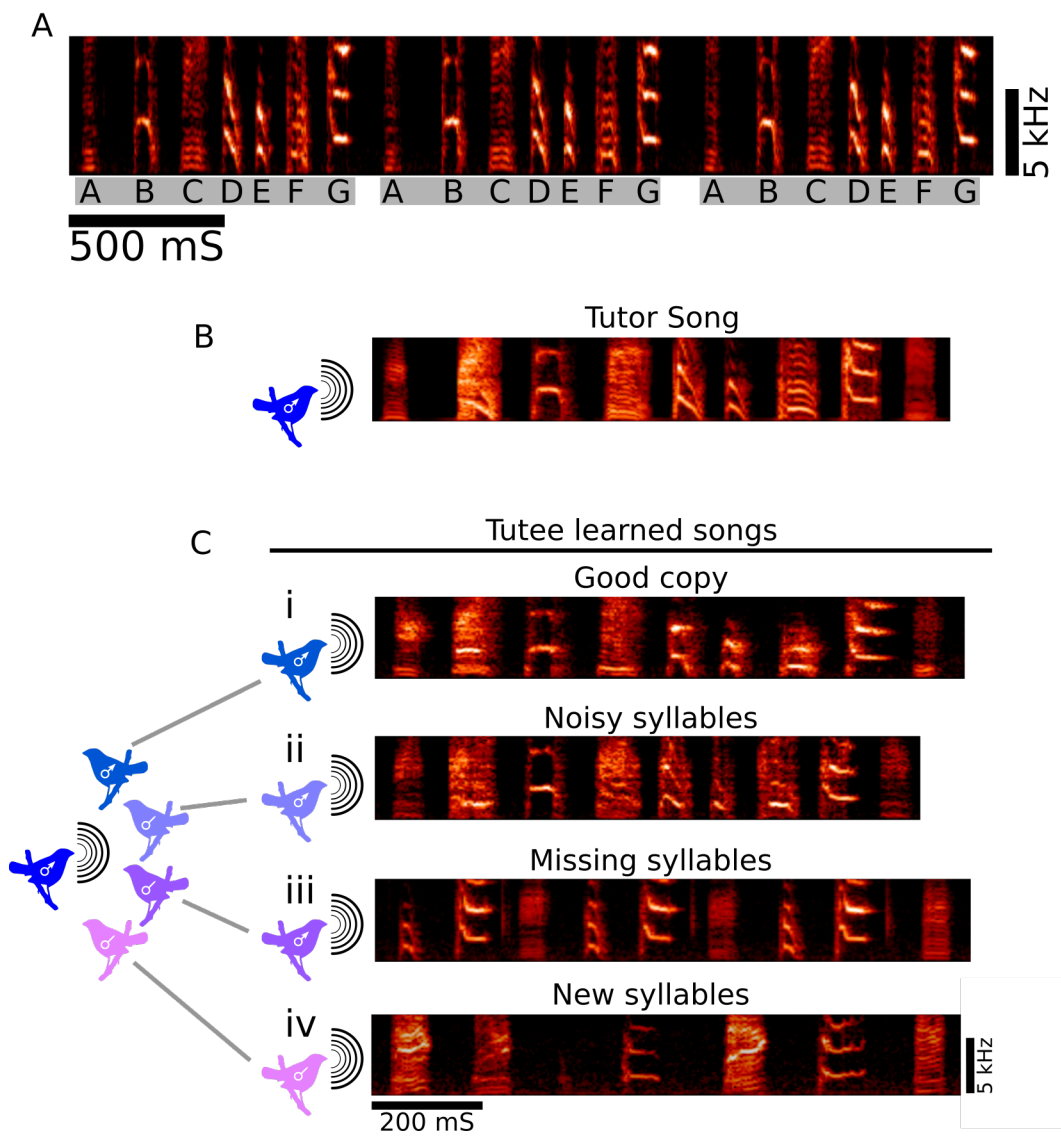
- 472 1. Doupe AJ, Kuhl PK. Birdsong and human speech: common themes and mechanisms. *Annu Rev*
473 *Neurosci.* 1999;22: 567–631. doi:10.1146/annurev.neuro.22.1.567
- 474 2. Brainard MS, Doupe AJ. Translating birdsong: songbirds as a model for basic and applied medical
475 research. *Annu Rev Neurosci.* 2013;36: 489–517. doi:10.1146/annurev-neuro-060909-152826
- 476 3. Catchpole C, Slater PJB. *Bird song : biological themes and variations.* Cambridge University Press; 2003.
- 477 4. Tchernichovski O, Lints T, Mitra PP, Nottebohm F. Vocal imitation in zebra finches is inversely related to
478 model abundance. *Proc Natl Acad Sci U S A.* 1999;96: 12901–4. Available:
479 <http://www.ncbi.nlm.nih.gov/pubmed/10536020>
- 480 5. Thorpe WH. The Process of Song-Learning in the Chaffinch as Studied by Means of the Sound
481 Spectrograph. *Nature.* 1954;173: 465–469. doi:10.1038/173465a0

- 482 6. Brainard MS, Doupe AJ. Interruption of a basal ganglia[ndash]forebrain circuit prevents plasticity of
483 learned vocalizations. *Nature*. Nature Publishing Group; 2000;404: 762–766. doi:10.1038/35008083
- 484 7. Konishi M. The Role of Auditory Feedback in the Control of Vocalization in the White-Crowned
485 Sparrow. *Z Tierpsychol*. Blackwell Publishing Ltd; 1965;22: 770–783. doi:10.1111/J.1439-
486 0310.1965.TB01688.X
- 487 8. Scharff C, Nottebohm F. A comparative study of the behavioral deficits following lesions of various parts
488 of the zebra finch song system: implications for vocal learning. *J Neurosci*. 1991;11.
- 489 9. Tchernichovski, Nottebohm, Ho, Pesaran, Mitra. A procedure for an automated measurement of song
490 similarity. *Anim Behav*. 2000;59: 1167–1176. doi:10.1006/anbe.1999.1416
- 491 10. Burkett ZD, Day NF, Peñagarikano O, Geschwind DH, White SA. VoICE: A semi-automated pipeline for
492 standardizing vocal analysis across models. *Sci Rep*. Nature Publishing Group; 2015;5: 10237.
493 doi:10.1038/srep10237
- 494 11. Wu W, Thompson J a, Bertram R, Johnson F. A statistical method for quantifying songbird phonology and
495 syntax. *J Neurosci Methods*. 2008;174: 147–54. doi:10.1016/j.jneumeth.2008.06.033
- 496 12. Mandelblat-Cerf Y, Fee MS, Nottebohm F, Williams H, Marler P, Tchernichovski O, et al. An Automated
497 Procedure for Evaluating Song Imitation. Bolhuis JJ, editor. *PLoS One*. Public Library of Science;
498 2014;9: e96484. doi:10.1371/journal.pone.0096484
- 499 13. Ho CE, Pesaran B, Fee MS, Mitra PP. Characterization of the structure and variability of zebra finch song
500 elements. *Proceedings of the joint symposium on neural computation*. 1998. pp. 76–83.
- 501 14. Pearson K, Pearson A. Contributions to the Mathematical Theory of Evolution. *Source Philos Trans R*
502 *Soc London A*. 1894;185: 71–110. Available: <http://www.jstor.org/stable/90667>
- 503 15. Dempster AP, Laird NM, Rubin DB. Maximum Likelihood from Incomplete Data via the EM Algorithm.
504 *Source J R Stat Soc Ser B J R Stat Soc Ser B. Methodological*; 1977;39: 1–38. Available:
505 <http://www.jstor.org/stable/2984875>
- 506 16. Schwarz G. Estimating the Dimension of a Model. *Ann Stat*. Institute of Mathematical Statistics; 1978;6:
507 461–464. doi:10.1214/aos/1176344136
- 508 17. Welch P. The use of fast Fourier transform for the estimation of power spectra: A method based on time
509 averaging over short, modified periodograms. *IEEE Trans Audio Electroacoust*. 1967;15: 70–73.
510 doi:10.1109/TAU.1967.1161901
- 511 18. Nordeen KW, Nordeen EJ. Auditory feedback is necessary for the maintenance of stereotyped song in
512 adult zebra finches. *Behav Neural Biol*. 1992;57: 58–66. Available:
513 <http://www.ncbi.nlm.nih.gov/pubmed/1567334>
- 514 19. Kogan JA, Margoliash D. Automated recognition of bird song elements from continuous recordings using
515 dynamic time warping and hidden Markov models: A comparative study.
516 <http://dx.doi.org/101121/1421364>. *Acoustical Society of America*; 1998; doi:10.1121/1.421364
- 517 20. Azizyan M, Singh A, Wasserman L. Efficient Sparse Clustering of High-Dimensional Non-spherical

- 518 Gaussian Mixtures. 2014; Available: <http://arxiv.org/abs/1406.2206>
- 519 21. von Luxburg U. A Tutorial on Spectral Clustering. 2007; Available: <http://arxiv.org/abs/0711.0189>
- 520 22. Johnson SC. Hierarchical clustering schemes. *Psychometrika*. Springer-Verlag; 1967;32: 241–254.
521 doi:10.1007/BF02289588
- 522 23. Otsu N. A Threshold Selection Method from Gray-Level Histograms. 1979;9: 62–66.
- 523 24. Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine
524 Learning in Python. *J Mach Learn Res*. 2011;12: 2825–2830.
525

526

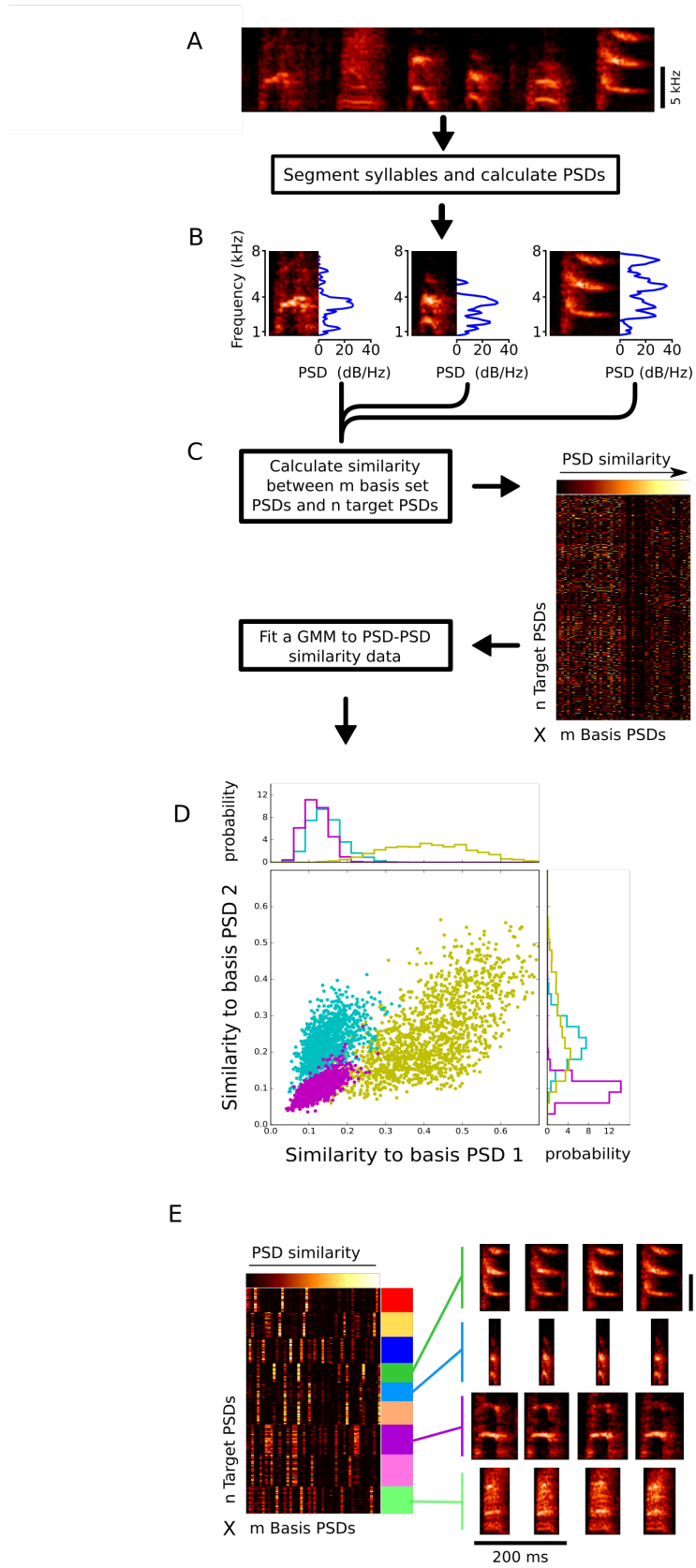
Fig 1



527 Fig 1. Quantification of song learning is complicated by variety in both learning and failure to learn.
528 (A) Typical sample of song from an adult Bengalese finch. Song is composed of a set of categorically
529 distinct syllable types (labeled ‘A’, ‘B’, ‘C’...) that are organized into larger, repeated, sequences (gray
530 bars). Both the spectral structure of syllables and their sequencing are learned features of song. Hence,
531 song is a complex, high dimensional behavior that differs across individuals. (B) Song of an adult male
532 ‘tutor’ and (C) songs of four juvenile ‘tutees’ that were all exposed to the same tutor song, illustrating
533 variation in the quality of song learning. (Ci) Song from a tutee that learned the spectral content of the
534 tutor song well, producing a song with accurate copies of all syllables. (Cii) Song from a tutee that
535 copied all syllables, but with noisier versions than those present in the tutor song. (Ciii) song from a
536 tutee that failed to copy some of the syllables from the tutor song. (Civ) song from a tutee that
537 included ‘new syllables’ that were not clearly present in the tutor song.
538

539

Fig 2



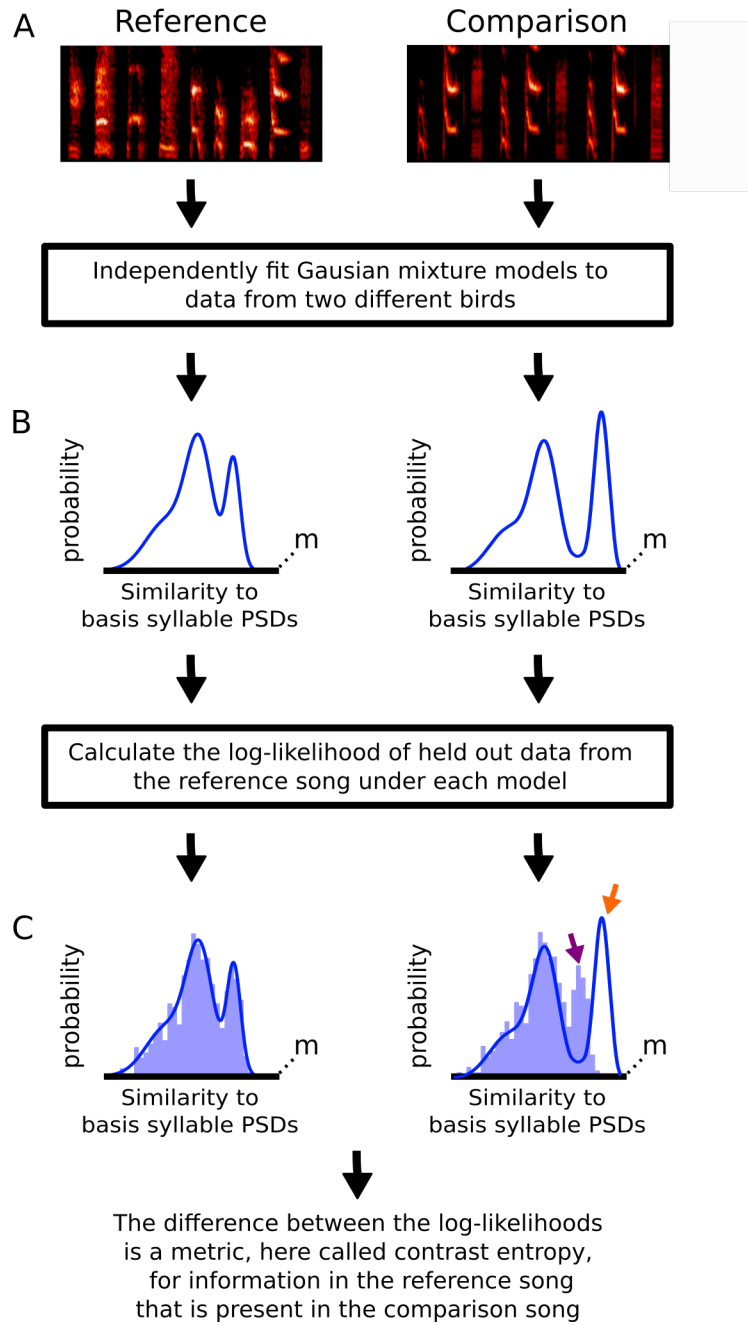
540

541 Fig 2. Assembly of statistical models for song. (A) To assemble a statistical model for the song of a
542 given bird, we first segment all syllables from a set of songs produced by that bird and compute their
543 corresponding PSDs. (B) Three examples of segmented syllables, each of a different type, and their
544 corresponding PSDs. (C) For each of 3000 ‘target’ syllables from the song to be modeled, similarity of
545 PSDs is calculated relative to a basis set of PSDs for 50 syllables randomly drawn from the same song.
546 This creates an M (number of basis syllables) by N (number of target syllables) similarity matrix. (D)
547 Visualization of how transformation of raw syllable data into the syllable similarity matrix results in a
548 clustering of syllables by type. Each point in the plot indicates the similarity between the PSD for one
549 target syllable and two basis PSDs (‘basis PSD1’ and ‘basis PSD2’) from the set of 50 basis PSDs. For
550 clarity of exposition, only data that fall into one of three regions of high density are plotted here. Each
551 of these regions corresponds approximately to multiple instances of one syllable type (which cluster
552 near each other because of the similarity in their PSDs). In practice, there were more than three regions
553 of syllable clustering (corresponding approximately to the number distinct syllable types in the bird’s
554 song), and these regions were represented in the 50 dimensional space defined by the basis set of PSDs
555 (only two of which are illustrated here). The regions of high density in this similarity space were fit
556 with a Gaussian mixture model, in which the optimal number of Gaussian mixtures was determined by
557 Bayesian Information Criteria. Individual data points here are color-coded by their assignment to one
558 of three Gaussian mixtures, with data points corresponding to 6 additional Gaussian mixtures not
559 shown. In any single dimension (top and right) data points assigned to each Gaussian mixture were
560 approximately normally distributed. (E) Similarity matrix shown in C, reordered so that data are
561 grouped by assignment to each of 9 Gaussian mixtures fit to the data (represented by colored blocks at
562 the right of the similarity matrix). In this reordered representation, it is apparent that syllables
563 assigned to each Gaussian mixture have a shared ‘bar code’ reflecting a shared pattern of PSD
564 similarity values relative to the basis PSDs. The spectrograms at the right illustrate that syllables
565 assigned to a given Gaussian mixture tend to be of the same type.

566

567

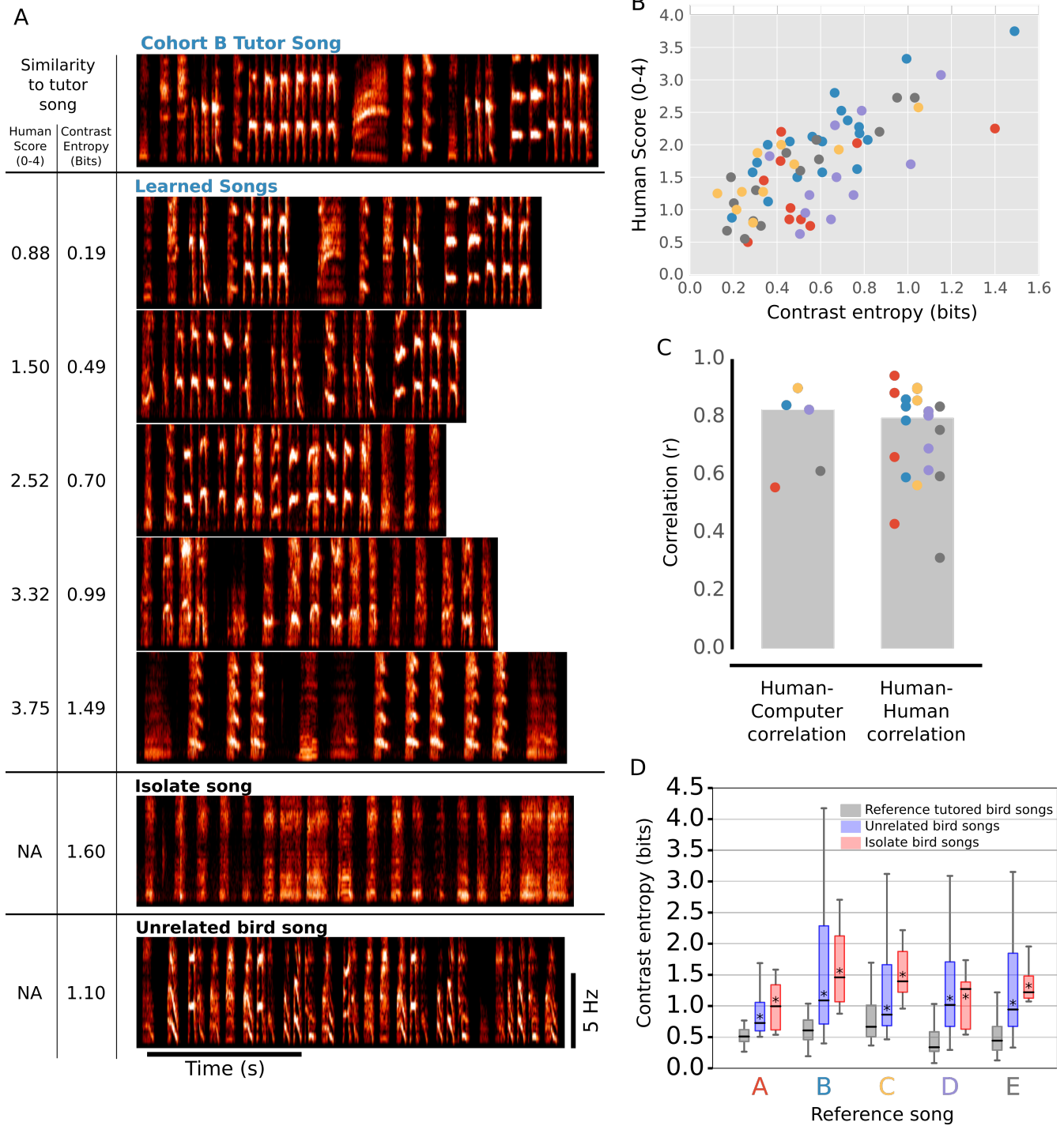
Fig 3



568

569 Fig 3. Estimation of the amount of spectral information present in the reference (tutor) song that is
570 absent from the comparison (tutee) song. (A) Example reference and comparison songs. To compute
571 contrast entropy for these songs, we first fit GMMs to the data from each song. (B) Representation in
572 one dimension of the GMMs fit to song spectral content for both the reference song (left) and the
573 comparison song (right). (C) We then calculate the log-likelihood of reference song syllables withheld
574 from model fitting (light blue histogram) under both the model for the reference song and the model for
575 the comparison song. Spectral content present in the reference song and not well represented in the
576 comparison song model (purple arrow) will result in lower likelihood. Consequently, the difference
577 (CE) between the two (reference song and comparison song) log-likelihoods will increase if there is
578 information present in the reference song that is absent from the comparison song. However,
579 information present in the comparison song, but not in the reference song (orange arrow) will not
580 impact the CE.
581

Fig 4



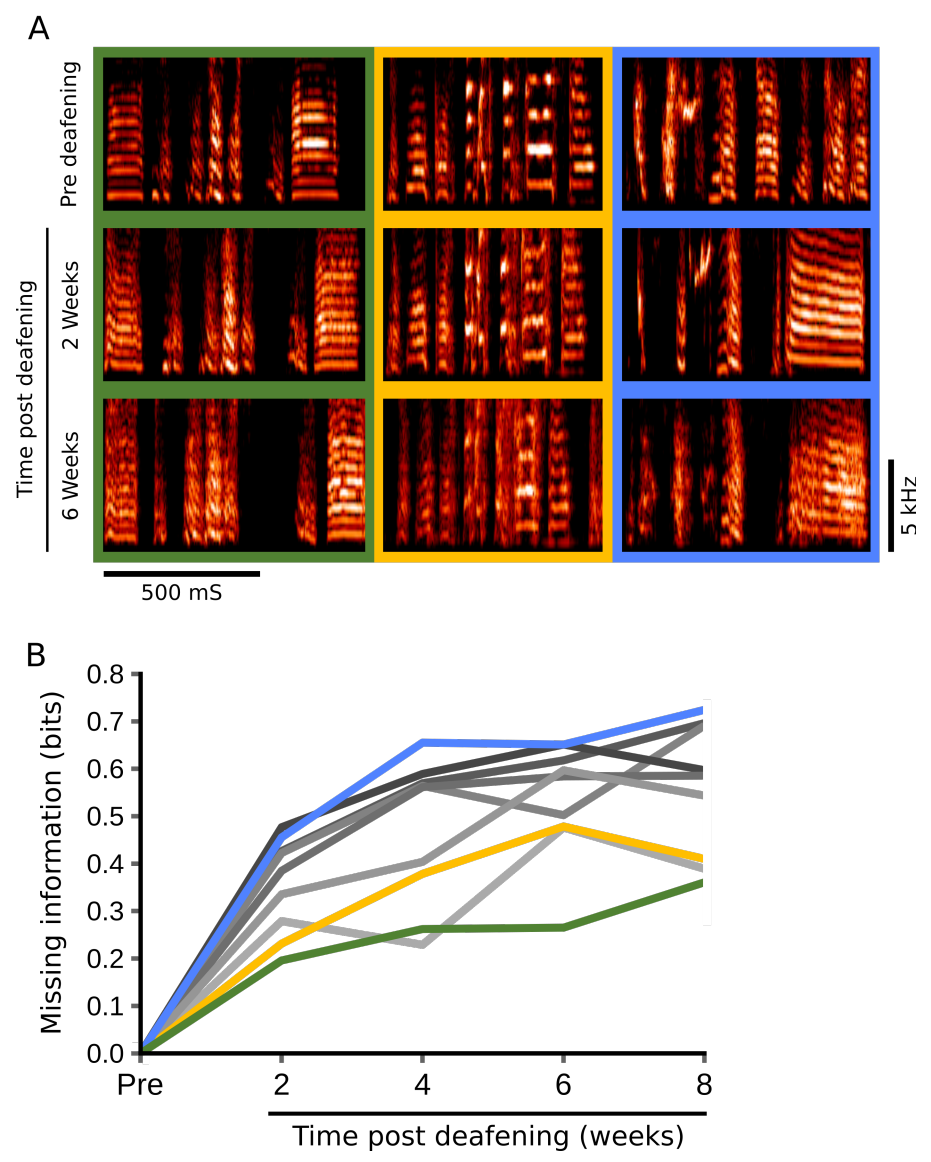
583

584 Fig 4. Contrast entropy closely parallels human assessment of learning outcomes. The quality of
585 learning for individuals from five cohorts, each with a distinct tutor song, were evaluated by contrast
586 entropy (CE) and human inspection. (A) Example spectrograms of the tutor song from one cohort and
587 the songs of 5 tutees from the same cohort (cohort B). Also shown for comparison is the song of one
588 isolate bird, raised without tutor song exposure (isolate song), and the song from one bird raised with a
589 different tutor (unrelated bird song). Numbers at left indicate the CE and human similarity scores for
590 each song relative to the tutor song from cohort B. (B) There was a good correspondence between CE
591 and human evaluations of learning across a broad range. Here, human scores are the average of four
592 human judges. Across all five cohorts, CE and human scores were well correlated ($p < 0.01$, $r = 0.722$,
593 OLS). (C) Comparison of CE and human scores for each of the five cohorts. Human-computer
594 correlation (left) shows the correlation between CE values and average human scores for each of the
595 five cohorts. Human-human correlation (right) indicates the correlation between the scores of each of 4
596 individual humans and the average of the remaining human scores for each cohort. Medians are
597 indicated as gray bars. (D) Summary of CE scores for the five cohorts (gray) were significantly lower
598 than scores from a cohort of unrelated birds (purple, $p < 0.01$, Wilcoxon rank test) and from a cohort of
599 'isolate birds' raised without a tutor (red, $p < 0.01$, Wilcoxon rank test). Across all panels, bird cohort
600 identity is indicated by color.

601

602

Fig 5

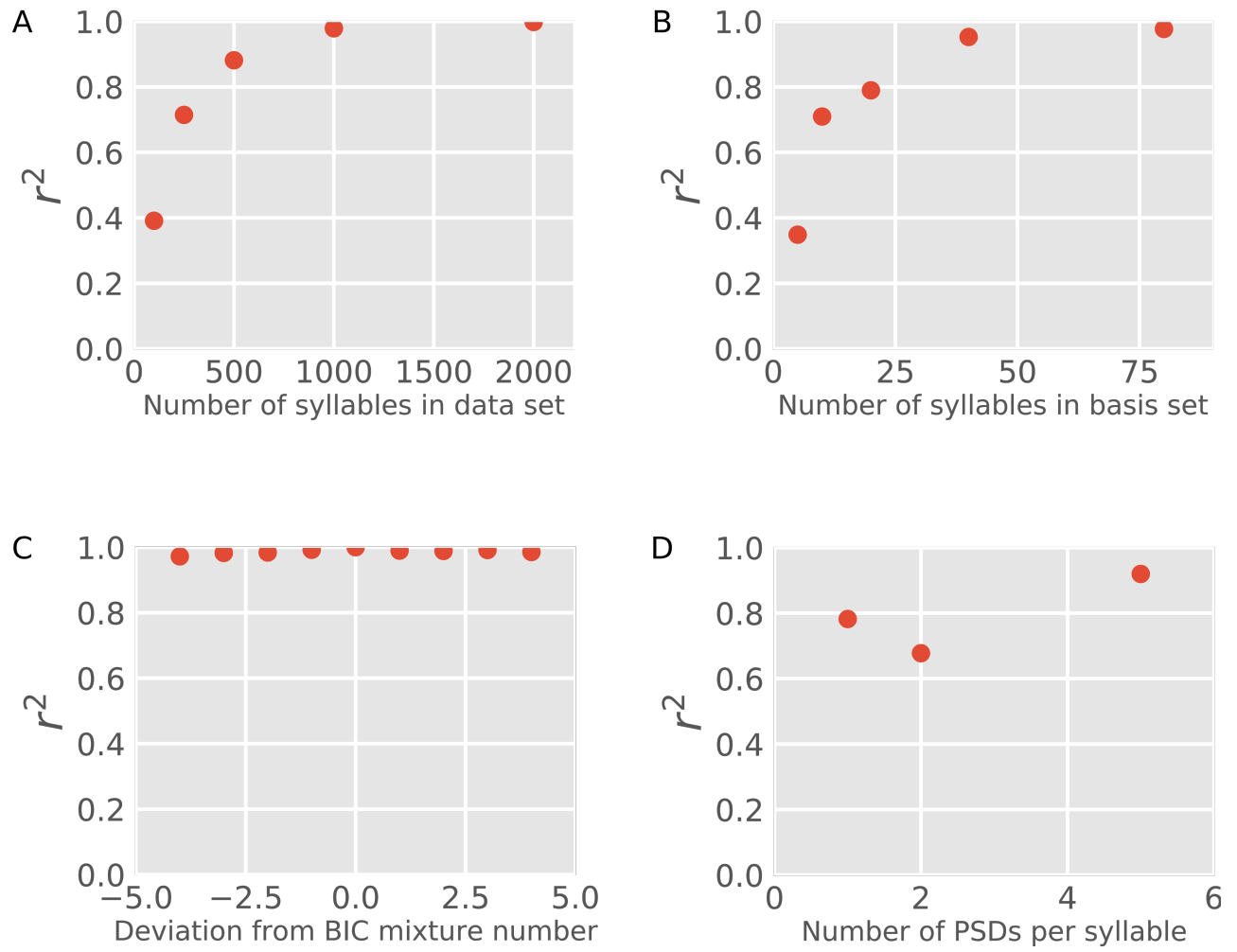


603

604 Fig 5. Quantification of changes to song following deafening. (A) Spectrograms from before (pre), two
605 weeks, and six weeks post deafening for three Zebra finches demonstrate typical disruption to the
606 spectral content of song due to deafening. (B) CE values indicating information missing from post
607 deafening songs relative to baseline reference songs for nine birds at two, four, six, and eight weeks
608 following deafening. Colors indicate bird identity, with green, yellow and blue in panels A and B
609 illustrating data from birds that had small, intermediate and large changes to song spectral structure
610 following deafening.

611

Fig 6



613

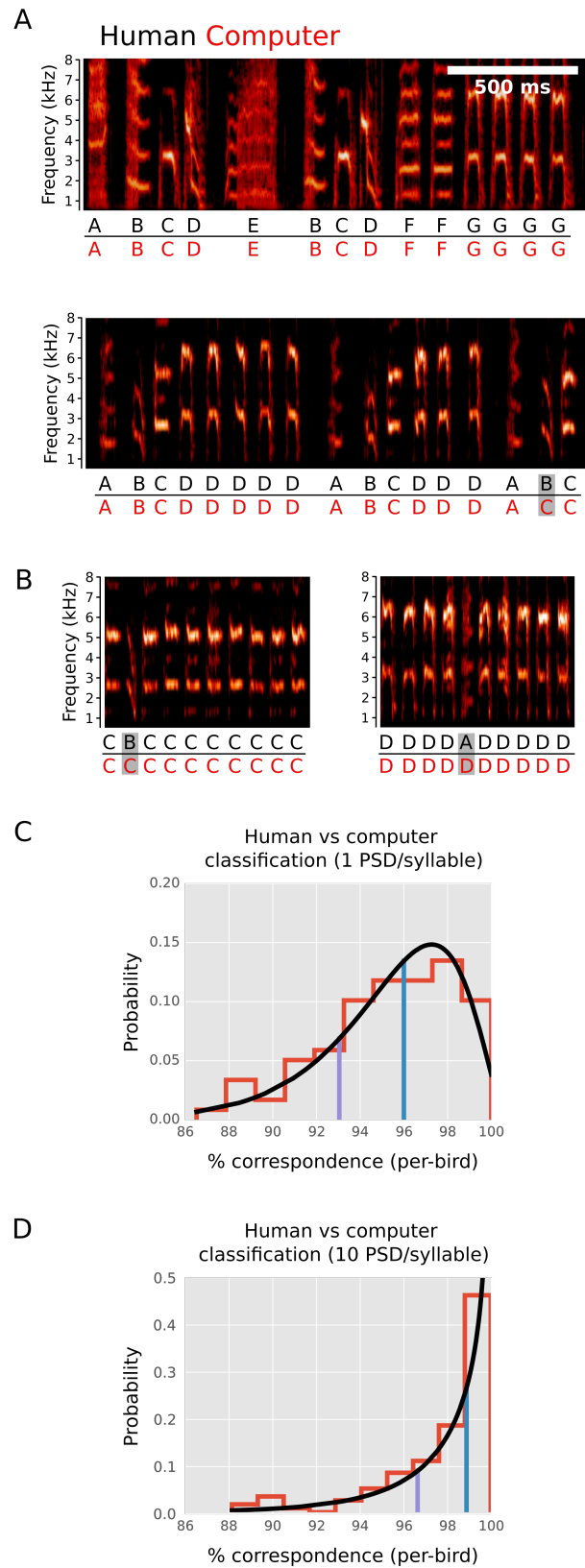
614 Fig 6. Contrast entropy song similarity measures are robust to variation in defined parameters.

615 (A) Plot of r^2 values for correlations between CE calculated using a range of input data sizes and CE
616 calculated using 3000 syllables of input data. (B) Plot of r^2 values for correlations between CE
617 calculated using a range of basis set sizes and CE calculated using a 160 syllable basis set. (C) Plot of r^2
618 values for correlations between CE calculated using the number of mixture components (k) determined
619 by BIC ($nBIC$) and CE calculated using a number of mixture components ranging from $nBIC-4$ to
620 $nBIC+4$. (D) Plot of r^2 values for correlations between CE calculated using a single PSD representation
621 of a syllable and CE calculated using a 10 PSD representation.

622

623

Fig 7



624
625 Fig 7. GMM derived syllable classifications are correlated with human syllable classifications. (A)
626 Examples of labels assigned to two songs by human inspection (black) and GMM (red). For many
627 birds, there were no differences between human assigned and GMM assigned labels (e.g. upper panel).
628 However, for some birds, there were discrepancies (e.g. gray box, lower panel). (B) Erroneous GMM
629 classifications can be identified by inspection of spectrograms for groups of syllables assigned to a
630 given Gaussian mixture. Illustrated here are two examples of groups of syllables assigned to individual
631 Gaussian mixtures, where it is apparent in each case that a single syllable (gray boxes) is miss-
632 classified relative to human assignment. For 90 animals, the number of miss-classified syllables was
633 determined by such human inspection of groups of syllables that were assigned to each Gaussian
634 mixture. (C) Distribution of the percent of correctly classified syllables (per-bird) is shown in red with
635 a gamma distribution fit to these data shown in black. 50% of animals had greater than 96% correctly
636 classified syllables (blue line) while 80% had more than 93% correctly classified syllables (purple
637 line). (D) Distribution of the percent of correctly classified syllables per bird is shown as in C, but here
638 with categorization carried out in which the input representation of each syllable to the GMM includes
639 10 PSDs evenly spaced over the duration of the syllable, rather than a single PSD for the entire
640 syllable. Using this richer representation of syllable, 50% of animals had more than 99% correctly
641 classified syllables (blue line) while 80% had more than 96% correctly classified syllables (purple
642 line).
643
644