
MEASURING CONTEXT DEPENDENCY IN BIRDSONG USING ARTIFICIAL NEURAL NETWORKS

A PREPRINT

Takashi Morita^{1,2}

Hiroki Koda²

Kazuo Okanoya³

Ryosuke O. Tachibana^{3*}

¹SANKEN, Osaka University, JAPAN

²Primate Research Institute, Kyoto University, JAPAN

³Center for Evolutionary Cognitive Sciences, Graduate School of Arts and Sciences, the University of Tokyo, JAPAN

July 23, 2021

ABSTRACT

Context dependency is a key feature in sequential structures of human language, which requires reference between words far apart in the produced sequence. Assessing how long the past context has an effect on the current status provides crucial information to understand the mechanism for complex sequential behaviors. Birdsongs serve as a representative model for studying the context dependency in sequential signals produced by non-human animals, while previous reports were upper-bounded by methodological limitations. Here, we newly estimated the context dependency in birdsongs in a more scalable way using a modern neural-network-based language model whose accessible context length is sufficiently long. The detected context dependency was beyond the order of traditional Markovian models of birdsong, but was consistent with previous experimental investigations. We also studied the relation between the assumed/auto-detected vocabulary size of birdsong (i.e., fine- vs. coarse-grained syllable classifications) and the context dependency. It turned out that the larger vocabulary (or the more fine-grained classification) is assumed, the shorter context dependency is detected.

Keywords birdsong, context dependency, Bengalese finch, language modeling, discrete variational autoencoder, unsupervised clustering, individual normalization

1 Introduction

2 Making behavioral decisions based on past information is a crucial task in the life of humans and animals [1, 2]. Thus,
3 it is an important inquiry in biology how far past events have an effect on animal behaviors. Such past records are not
4 limited to observations of external environments, but also include behavioral history of oneself. A typical example is
5 human language production; The appropriate choice of words to utter depends on previously uttered words/sentences.
6 For example, we can tell whether ‘was’ or ‘were’ is the grammatical option after a sentence ‘*The photographs that were*
7 *taken in the cafe and sent to Mary ___*’ only if we keep track of the previous words sufficiently long, at least up to
8 ‘*photographs*’, and successfully recognize the two closer nouns (*cafe* and *station*) as modifiers rather than the main
9 subject. Similarly, semantically plausible words are selected based on the topic of preceding sentences, as exemplified
10 by the appropriateness of *olive* over *cotton* after “sugar” and “salt” are used in the same speech/document. Such
11 dependence on the production history is called context dependency and is considered a characteristic property of human
12 languages [3, 4, 5, 6].

*Corresponding Author: rtachi@gmail.com

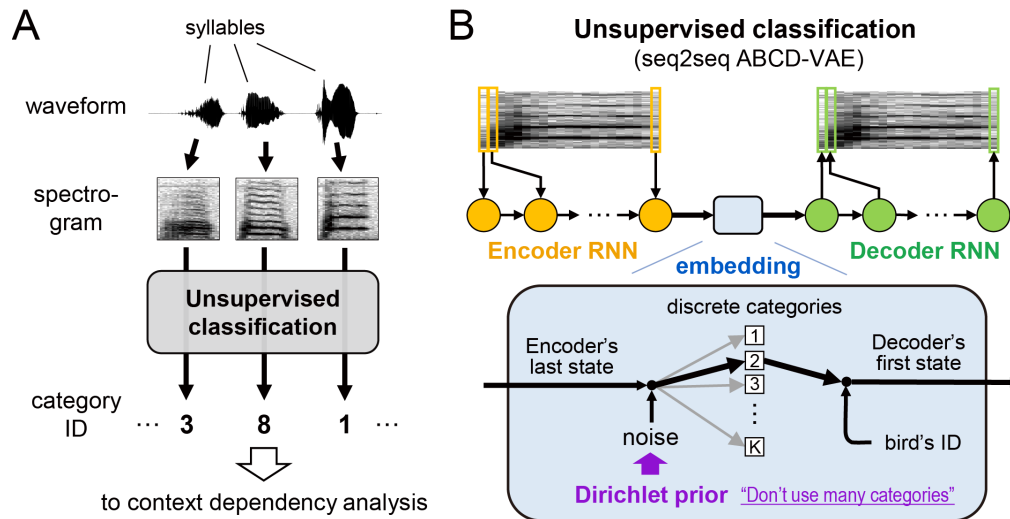


Figure 1. Schematic diagram of newly proposed syllable classification. (A) Each sound waveform segment was converted into the time-frequency representations (spectrograms), and was assigned to one of syllable categories by the unsupervised classification. (B) The unsupervised classification was implemented as a sequence-to-sequence version of the variational autoencoder, consisting of the attention-based categorical sampling with the Dirichlet prior (“seq2seq ABCD-VAE”). The ABCD-VAE encoded syllables into discrete categories between the encoder and the decoder. A statistically optimal number of categories was detected under an arbitrarily specified upper bound thanks to the Dirichlet prior. The identity of the syllable-uttering individual was informed to the decoder besides the syllable categories; Accordingly, individual-specific patterns need not have been encoded in the discrete syllable representation.

13 Birdsongs serve as a representative case study of context dependency in sequential signals produced by non-human
 14 animals. Their songs are sound sequences that consist of brief vocal elements, or *syllables* [7, 8]. Previous studies
 15 have suggested that those birdsongs exhibit non-trivially long dependency on previous outputs [9, 10, 11]. Complex
 16 sequential patterns of syllables have been discussed in comparison with human language syntax from the viewpoint of
 17 formal linguistics [8, 12]. Neurological studies also revealed homological network structures for the vocal production,
 18 recognition, and learning of songbirds and humans [13, 14, 15]. In this line, assessing whether birdsongs exhibit long
 19 context dependency is an important instance in the comparative studies, and several previous studies have addressed
 20 this inquiry using computational methods [16, 9, 11, 17, 18]. However, the reported lengths of context dependency
 21 were often measured using a limited language model (Markov/ n -gram model) that was only able to access a few recent
 22 syllables in the context. Thus, it is unclear if those numbers were real dependency lengths in the birdsongs or merely
 23 model limitations. Moreover, there is accumulating evidence that birdsong sequencing is not precisely modeled by a
 24 Markov process [16, 17].

25 The present study aimed to assess the context dependency in songs of Bengalese finches (*Lonchura striata* var.
 26 *domestica*) using modern techniques for natural language processing. Recent advancements in the machine learning
 27 field, particularly in artificial neural networks, provide powerful language models [19, 6], which can simulate various
 28 time series data without hypothesizing any particular generative process behind them. The neural network-based models
 29 also have a capacity to effectively use information in 200–900 syllables from the past (when the data include such
 30 long dependency) [5, 6], and thus, the proposed analysis no longer suffers from the model limitations in the previous
 31 studies. We performed the context dependency analysis in two steps: unsupervised classification of song syllables
 32 and context-dependent modeling of the classified syllable sequence. The classification enabled flexible modeling
 33 of statistical ambiguity among upcoming syllables, which are not necessarily similar to one another in acoustics.
 34 Moreover, it is preferable to have a common set of syllable categories, which is shared among classifications for all
 35 birds, to represent general patterns in the sequences and also to provide the language model with as big data as possible.
 36 Conventional classification methods depending on manual labeling by human experts could spoil such generality due
 37 to arbitrariness in integrating the category sets across different birds. To satisfy these requirements, we employed a
 38 novel, end-to-end, unsupervised clustering method (“seq2seq ABCD-VAE”, see Fig. 1). Then, we assessed the context
 39 dependency in sequences of the classified syllables by measuring the effective context length [5, 6], which represents
 40 how much portion of the song production history impacts on the prediction performance of a language model. The
 41 language model we used (“Transformer”, see Fig. 4) behaves as a simulator of birdsong production, which exploits the
 42 longest context among currently available models [19, 6].

43 Results

44 *Unsupervised, individual-invariant classification of Bengalese finch syllables*

45 We first converted birdsong syllables into discrete representations, or “labels”. When predicting an upcoming syllable
46 from previous outputs, probable candidates can have non-similar acoustic profiles. For example, “bag” and “beg”
47 in English are similar to each other in terms of phonology but have different syntactic and semantic distributions,
48 belonging to different grammatical categories (noun and verb, respectively). An appropriate language model must
49 assign a more similar probability to syntactically/semantically similar words like “bag” and “wallet” than acoustically
50 similar ones like “bag” and “beg”. Likewise, it is desirable to perform the context dependency analysis of birdsong
51 based on a flexible model of sequence processing so that it can handle ambiguity about possible upcoming syllables
52 that do not necessarily resemble one another from acoustic perspectives. Categorizing continuous-valued signals and
53 predicting the assigned discrete labels based on a categorical distribution is a simple but effective way of achieving such
54 flexible models, especially when paired with deep neural networks [20, 21, 22]. Syllable classification has also been
55 adopted widely in previous studies of birdsong syntax [7, 23, 11, 18].

56 Recent studies have explored fully unsupervised classification of animal vocalization based on acoustic features
57 extracted by an artificial neural network, called variational autoencoder or VAE [24, 25, 26]. We extended this approach
58 and proposed a new end-to-end unsupervised clustering method named ABCD-VAE, which utilizes the attention-based
59 categorical sampling with the Dirichlet prior (see also [27]). This method automatically classifies syllables into an
60 unspecified number of categories in a statistically principled way. It also allowed us to exploit the speaker-normalization
61 technique developed for unsupervised learning of human language from speech recordings [28, 29], yielding syllable
62 classification modulo individual variation. Having common syllable categories across different individuals helps us
63 build a unified model of syllable sequence processing. Individual-invariant classification of syllables is also crucial
64 for deep learning-based analysis that requires a substantial amount of data; i.e., it is hard to collect sufficient data for
65 training separate models on each individual.

66 We used a dataset of Bengalese finches’ songs that was originally recorded for previous studies [30, 31]. Song syllables
67 in the recorded waveform data were detected and segmented by amplitude thresholding. We collected 465,310 syllables
68 in total from 18 adult male birds. Some of these syllables were broken off at the beginning/end of recordings. We
69 filtered out these incomplete syllables, and fed the other 461,994 syllables to the unsupervised classifier (Fig. 1A).
70 The classifier consisted of two concatenated recurrent neural networks (RNNs, see Fig. 1B). We jointly trained the
71 entire network such that the first RNN represented the entirety of each input syllable in its internal state (“encoding”
72 Fig. 1B) and the second RNN restored the original syllable from the internal representation as precisely as possible
73 (“decoding”). The encoded representation of the syllable was mapped to a categorical space (“embedding”) before the
74 decoding process. The number of syllable categories was automatically detected as a statistical optimum owing to the
75 Dirichlet prior [32].

76 As a result, the classifier detected 37 syllable categories in total for all the birds (Fig. 2B). Syllables that exhibited
77 similar acoustic patterns tended to be classified into the same category across different birds (Fig. 2A). All birds
78 produced not all but a part of syllable categories in their songs (Fig. 2C). The syllable repertoire of each bird covered
79 24 to 36 categories (32.39 ± 3.35). The detected syllable vocabulary size was greater than the number of annotation
80 labels used by a human expert (5–14) [30]. Conversely, each category consisted of syllables produced by 7 to 18 birds
81 (15.76 ± 2.91). The detected categories appeared to align with major differences in the spectrotemporal pattern (Fig. 2B).

82 *Quantitative evaluation of syllable classification for Bengalese finch*

83 Speaker-invariant clustering of birdsong syllables should meet at least two desiderata: (i) the resulting classification
84 must keep consistency with the conventional bird-specific classification (i.e., clustered syllables must belong to the
85 same bird-specific class), and (ii) the discovered syllable categories should be anonymized. Regarding (i), we evaluated
86 the alignment of the detected classification with manual annotations by a human expert [30]. We scored the alignment
87 using two metrics. One was Cohen’s Kappa coefficient [33], which has been used to evaluate syllable classifications
88 in previous studies [9, 30]. A problem with this metric is that it requires two classifications to use the same set of
89 categories while our model predictions and human annotations had different numbers of categories and, thus, we needed
90 to force-align each of the model-predicted categories to the most common human-annotated label to use the metric
91 [9]. For example, suppose that the model classified 300 syllables into a category named “A”. If 100 of the syllables
92 in “A” are labeled as “a” by the human annotator and the other 200 are labeled as “b”, then all the syllables in “A”
93 received “b” as their force-aligned label of model predictions. This force-alignment makes the 100 syllables misaligned
94 with their original label “a”. Thus, the force-alignment scores uniformity of syllables within the model-predicted

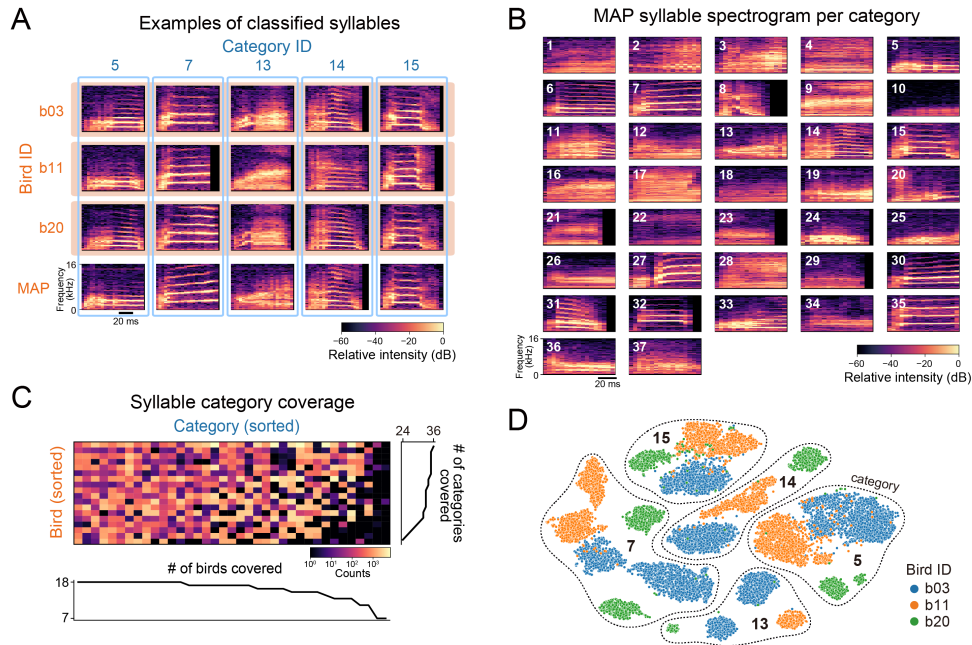


Figure 2. Clustering results of Bengalese finch syllables based on the ABCD-VAE. (A) Syllable spectrograms and their classification across individuals. Syllables in each of the first to third rows (orange box) were sampled from the same individual. Each column (blue frame) corresponds to the syllable categories assigned by the ABCD-VAE. The bottom row provides the spectrogram of each category with the greatest classification probability (MAP: maximum-a-posteriori) over all the individuals. The individual-specific examples also had the greatest classification probability (> 0.999) among the syllables of the same individual and category. (B) Spectrogram of the MAP syllable in each category. (C) Syllable counts per individual bird (rows) and category (columns). The number of non-zero entries is also reported in the line plots. (D) Comparison between syllable embeddings by the canonical continuous-valued VAE with the Gaussian noise (scatter points) and classification by the ABCD-VAE (grouped by the dotted lines). The continuous representation originally had 16 dimensions and was embedded into the 2-dimensional space by t-SNE. The continuous embeddings included notable individual variations represented by colors, whereas the ABCD-VAE classification ignored these individual variations.

95 categories regarding the manual annotations. To get rid of the force-alignment and any other post-processing, we also
 96 evaluated the classification using a more recently developed metric called homogeneity [34]. The homogeneity checks
 97 whether the category-mate syllables according to the ABCD-VAE were annotated with the same manual label (see
 98 the Method for its mathematical definition). Note that the homogeneity does not penalize overclassification (see the
 99 supporting information S1.5 for additional evaluation that takes overclassification into account). For example, suppose
 100 that the ABCD-VAE classified 300 syllables into a category named “A” and another 300 into “B”. The homogeneity is
 101 maximized even if all the 300 syllables in “A” are labeled “a” and all the 300 in “B” are also labeled as “a”. This is
 102 because all the category-mate syllables receive the same label. Instead, the homogeneity penalizes label mismatches
 103 within the model-detected categories, as in the case where 200 of the “A” syllables are labeled “a” and the other 100 are
 104 labeled “b”. Thus, the homogeneity is considered a unified version of Cohen’s kappa plus force-alignment.

105 To assess fulfillment of the second desideratum for ideal clustering (ii), we quantified the speaker-normalization effect
 106 of the ABCD-VAE by measuring the perplexity of speaker identification. We built a simple speaker identification model
 107 based on a syllable category uttered by the target bird, fitting the conditional categorical distribution to 90% of all the
 108 syllables by the maximum likelihood criterion and then evaluating the prediction probabilities on the other 10%. The
 109 prediction probabilities of the test data were averaged in the log scale (= entropy) and then exponentiated to yield the
 110 perplexity. Intuitively, the perplexity tells the expected number of birds among whom we have to guess by chance to
 111 identify the target speaker even after the information about the syllable category uttered by the target bird is provided.
 112 Thus, greater perplexity is an index of successful speaker-normalization.

113 We compared the performance of the ABCD-VAE with baseline scores provided by the combination of the canonical,
 114 continuous-valued VAE (which we call Gauss-VAE) [24, 25, 26] and the Gaussian mixture model (GMM) [35, 32, 36].

Table 1. Quantitative evaluation of the clustering by the ABCD-VAE for Bengalese finch syllables. Cohen’s kappa coefficient and homogeneity evaluated the alignment of the discovered clusters with manual annotations by a human expert. These scores for each individual bird were computed separately and their mean, maximum, and minimum over the individuals were reported since the manual annotation was not shared across individuals (see Method). Additionally, the perplexity of individual identification scored the amount of individuality included in the syllable categories yielded by the ABCD-VAE. The best scores are in boldface (results under the all-birds-together and bird-specific settings were ranked separately).

Method	# of clusters (source)	Cohen’s Kappa mean [min,max]	Homogeneity mean [min,max]	Speaker Perplexity
ABCD-VAE	37	0.8990 [0.7740, 0.9929]	0.9084 [0.7635, 0.9868]	8.0434
Gauss-VAE +	37 (ABCD-VAE)	0.7446 [0.5956, 0.8912]	0.7844 [0.6004, 0.9086]	4.0783
GMM (All-Birds- Together)	14 (manual) ≥128 (auto-detected)	0.6057 [0.4250, 0.8972] 0.8475 [0.5725, 0.9911]	0.6718 [0.5053, 0.8536] 0.8773 [0.6666, 0.9869]	6.7212 1.7112
Gauss-VAE +	37 (ABCD-VAE)	0.9304 [0.6619, 0.9906]	0.9292 [0.6479, 0.9893]	—
GMM (Bird-Specific)	5–14 (manual) 50–109 (auto-detected)	0.7888 [0.5012, 0.9328] 0.9516 [0.7629, 0.9982]	0.8090 [0.4732, 0.9254] 0.9505 [0.7687, 0.9962]	— —

115 This baseline model can be seen as a non-end-to-end version of our clustering method, having distinct optimizations
 116 for feature extraction and clustering. The Gauss-VAE was trained on the same datasets and by the same procedure
 117 as the ABCD-VAE. On the other hand, the GMM was trained in several ways. First, we built both bird-specific and
 118 common models: the former consisted of multiple models, each trained on data collected from a single individual
 119 bird, whereas the latter was a single model trained on the entire data collected from all the birds. The bird-specific
 120 clusterings provide “topline” scores because the gold-standard annotations by the human expert were also defined in a
 121 bird-specific way, and hence, they do not suffer from individual variations included in the Gauss-VAE features. On the
 122 other hand, the all-birds-together classifications tell us how much degree of difficulties exist in the clustering without
 123 end-to-end optimization or speaker normalization and, thus, serve as a baseline. Another kind of variation in the GMMs
 124 we tested was the number of syllable categories. We tested three ways of determining the number: (i) equals to the
 125 results from automatic detection by the ABCD-VAE, (ii) equals to the manual annotations by the human expert, and (iii)
 126 automatically detected from the distribution of syllable features defined by the Gauss-VAE. (i) and (ii) were obtained by
 127 specifying the number of mixture components of the GMM and training the GMM by the maximum likelihood criterion.
 128 On the other hand, (iii) was implemented by Bayesian estimation of active mixture components under the Dirichlet
 129 distribution prior [32].

130 As a result, the ABCD-VAE achieved a greater Kappa coefficient on average than the baseline models without subject-
 131 specific training (Table 1). Moreover, the comparison of the worst-bird scores (“min” in the table) showed that the
 132 ABCD-VAE was more robust than the topline models that were optimized to each bird separately. The ABCD-VAE
 133 achieved “almost perfect agreements” with the human expert ($\kappa > 0.8$) for sixteen of the eighteen birds and “substantial”
 134 agreements ($0.6 < \kappa \leq 0.8$) for the other two [37]. Similarly, the ABCD-VAE outperformed the baseline classifications
 135 in the average and worst-bird homogeneity scores. This result was also competitive with the topline models, especially
 136 regarding the worst-bird score. These results suggest that the syllable categories discovered by the ABCD-VAE
 137 kept consistency with the conventional subject-specific classifications, while the consistency was lost in the other
 138 all-birds-together classifications without speaker-normalization. In the meantime, the ABCD-VAE scored the greatest
 139 individual perplexity, indicating that the discovered syllable categories were more anonymized and individual-invariant
 140 than the baselines (see also Fig. 2D).

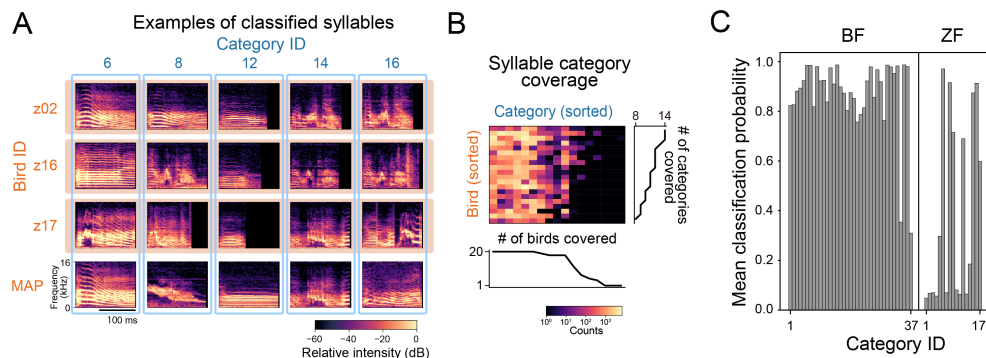


Figure 3. Clustering results of zebra finch syllables based on the ABCD-VAE. (A) Syllable spectrograms and their classification across individuals. Syllables in each of the first to third rows (orange box) were sampled from the same individual. Each column (blue frame) corresponds to the syllable categories assigned by the ABCD-VAE. The bottom row provides the spectrogram of each category with the greatest classification probability (MAP: maximum-a-posteriori) over all the individuals. The individual-specific examples had a top-5 classification probability among the syllables of the same individual and category. (B) Syllable counts per individual bird (rows) and category (columns). The number of non-zero entries is also reported in the line plots. (C) Mean classification probability of Bengalese finch (left) and zebra finch (right) syllables per category.

141 *Unsupervised classification of zebra finch syllables*

142 To further assess the effectiveness/limitations of the ABCD-VAE, the same clustering was performed on zebra finch
 143 syllables (*Taeniopygia guttata*). We collected 237,610 syllables from 20 adult male zebra finches. Again, the data
 144 included incomplete syllables that were broken off at the beginning/end of the syllables, and after filtering out those
 145 incomplete syllables, we fed the remaining 231,792 to the ABCD-VAE.

146 Speaker-normalized classification of zebra finch syllables was not as successful (or interpretable) as that of Bengalese
 147 finch syllables. While the syllables were classified into 17 categories in total (8 to 14 categories covered by a single
 148 bird, mean \pm SD:11.2 \pm 1.77), most of the classifications were not confident; 10 out of the 17 detected categories had a
 149 low mean classification probability under 30% whereas all but two categories of Bengalese finch syllables had a mean
 150 classification probability over 75% (Fig. 3C). Syllables with seemingly major spectral differences were force-aligned
 151 across individuals (Fig. 3A). Specifically, syllables consisting of multiple segments with distinct spectral patterns (or
 152 notes) seem to lack correspondents in different birds' repertoire (e.g., Category 14 and 16).

153 Quantitative evaluation also indicates that the speaker-normalized clustering of zebra finch syllables by the ABCD-VAE
 154 was not as well-aligned with bird-specific human annotations as that of Bengalese finch (Table 2). While the topline
 155 bird-specific models scored about 0.9 of Cohen's kappa coefficient and homogeneity, the scores of the ABCD-VAE
 156 stayed around 0.7. Nevertheless, it is of note that the ABCD-VAE outperformed the baseline all-birds-together models,
 157 except the one that automatically detected the number of categories (and achieved the upper bound at 128). This
 158 auto-detection model achieved high Cohen's kappa and homogeneity by specializing its categories to individual birds
 159 (i.e., by resorting to individual-specific classifications); as a result, the model scored a low individual perplexity,
 160 indicating that each individual was almost completely identifiable from the model-predicted category of a syllable. By
 161 contrast, the ABCD-VAE only used 17 categories and the high individual perplexity indicates that those categories were
 162 anonymized. Looking at each individual bird, the ABCD-VAE yielded "almost perfect agreement" with the manual
 163 annotations ($\kappa > 0.8$) for seven of the twenty birds, "substantial" agreement ($0.6 < \kappa \leq 0.8$) for other seven, and
 164 "moderate agreement" for the remaining two ($0.4 < \kappa \leq 0.6$).

165 *Analysis of context dependency*

166 The classification described above provided us sequences of categorically represented syllables. To assess the context
 167 dependency in the sequence, we then measured differences between syllables predicted from full-length contexts and
 168 truncated contexts. This difference becomes large as the length of the truncated context gets shorter and contains less
 169 information. And, the difference should increase if the original sequence has a longer context dependency (Fig. 4A).
 170 Thus, the context dependency can be quantified as the minimum length of the truncated contexts where the difference
 171 becomes undetectable [5, 6]. For the context-dependent prediction, we employed the Transformer language model

Table 2. Quantitative evaluation of the clustering by the ABCD-VAE for zebra finch syllables. Cohen’s kappa coefficient and homogeneity evaluated the alignment of the discovered clusters with manual annotations by a human expert. These scores for each individual bird were computed separately and their mean, maximum, and minimum over the individuals were reported since the manual annotation was not shared across individuals (see Method). Additionally, the perplexity of individual identification scored the amount of individuality included in the syllable categories yielded by the ABCD-VAE. The best scores are in boldface (results under the all-birds-together and bird-specific settings were ranked separately).

Method	# of clusters (source)	Cohen’s Kappa mean [min,max]	Homogeneity mean [min,max]	Speaker Perplexity
ABCD-VAE	17	0.7097 [0.4413, 0.9288]	0.6793 [0.4972, 0.8718]	12.2834
Gauss-VAE +	17 (ABCD-VAE)	0.6012 [0.2845, 0.9274]	0.6177 [0.3030, 0.8942]	4.3094
GMM (All-Birds- Together)	13 (manual) ≥128 (auto-detected)	0.6102 [0.0401, 0.9741] 0.8938 [0.6843, 0.9915]	0.6315 [0.0433, 0.9609] 0.9016 [0.7643, 0.9894]	5.7021 1.3092
Gauss-VAE +	17 (ABCD-VAE)	0.9579 [0.8847, 0.9938]	0.9545 [0.8828, 0.9905]	—
GMM (Bird-Specific)	4–13 (manual) 18–47 (auto-detected)	0.8762 [0.7915, 0.9744] 0.9812 [0.9360, 1.0000]	0.8623 [0.7056, 0.9607] 0.9782 [0.9274, 1.0000]	— —

172 [19, 6]. Transformer is known to capture long-distance dependency more easily than RNNs since it can directly refer to
 173 any data point in the past at any time while RNNs can only indirectly access past information through their internal
 174 memory [38, 19]. There is also accumulating evidence that Transformer successfully represents latent structures behind
 175 data, such as hierarchies of human language sentences [19, 39, 40].

176 Each sequence included syllables from a single recording. We report the analysis of both Bengalese and zebra finch
 177 songs, even though the classification of zebra finches’ syllables was not as reliable as Bengalese finches’. We obtained
 178 a total of 7,879 sequences of Bengalese finch syllables (each containing 8–338 syllables, 59.06 syllables on average)
 179 and 11,822 sequences of zebra finch syllables (each containing 1–219 syllables, 20.10 syllables on average), and used
 180 7,779 and 11,722 of them respectively to train the Transformer (see Table 3). The remaining 100 sequences were used
 181 to score its predictive performance from which the dependency was calculated. The model predictions were provided of
 182 the log conditional probability of the test syllables (x) given the preceding ones in the same sequence. We compared
 183 the model predictions between the full-context (“Full”, Fig. 4A) and the truncated-context (“Truncated”) conditions.
 184 Then, the context dependency was quantified by a statistical measure of the effective context length [5, 6], which is the
 185 minimum length of the truncated context wherein the mean prediction difference between the two contexts was not
 186 significantly greater than the canonical 1% threshold in perplexity [41].

187 To see the relation between the number of syllable categories and context dependency, we also performed the same
 188 analysis based on more coarse/fine-grained syllable classifications into 10 to 80, 160, and 320 categories. These
 189 classifications were derived from the k-means clustering on the L2-normalized feature vectors of syllables given by the
 190 ABCD-VAE.

191 The statistically effective context length (SECL) of the Bengalese finch song was eight based on the 37 syllable
 192 categories that were automatically detected by the ABCD-VAE (Fig. 4B). In other words, restricting available contexts
 193 to seven or fewer preceding syllables significantly decreased the prediction accuracy compared with the full-context
 194 baseline, while the difference became marginal when eight or more syllables were included in the truncated context.

195 When syllables were classified into more fine-grained categories, the difference between the model predictions based
 196 on the truncated and full contexts became smaller (Fig. 4B; $p < 0.001$ according to the linear regression of the loss
 197 difference on the number of syllable categories and the length of truncated contexts, both in the log scale). That is, the

Table 3. The size of the training and test data used in the neural language modeling of Bengalese and zebra finch songs. The “SECL” portion of the test syllables was used to estimate the SECL. The numbers of syllables in parentheses report the incomplete syllables that were broken off at the start/end of recordings, which were labeled with a distinct symbol.

Species	Usage	# of sequences	# of syllables	
			Total	SECL
Bengalese Finch	Training (incomplete)	7,779	458,992 (3,275)	—
	Test (incomplete)	100	6,557 (41)	4,657 (36)
Zebra Finch	Training (incomplete)	11,722	234,674 (5,763)	—
	Test (incomplete)	100	2,936 (55)	1,536 (49)

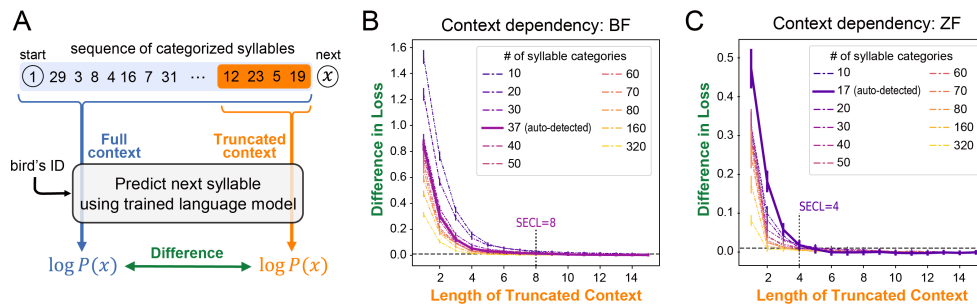


Figure 4. (A) Schematic diagram of the evaluation metric. Predictive probability of each categorized syllable (denoted by x) was computed using the trained language model, conditioned on the full and truncated contexts consisting of preceding syllables (highlighted in blue and orange, respectively). The logarithmic difference of the two predictive probabilities was evaluated, and SECL was defined by the minimum length of the truncated context wherein the prediction difference is not statistically significantly greater than a canonical threshold. (B) The differences in the mean loss (negative log probability) between the truncated- and full-context predictions of Bengalese finch songs and (C) zebra finch songs. The x-axis corresponds to the length of the truncated context. The error bars show the 90% confidence intervals estimated from 10,000 bootstrapped samples. The loss difference is statistically significant if the lower side of the intervals are above the threshold indicated by the horizontal dashed line.

198 context dependency traded off with the number of syllable categories. When 160 or 320 categories were assumed, the
 199 SECL of the Bengalese finch songs decreased to 5.

200 Zebra finch songs showed the same trade-off between the number of syllable categories and context dependency.
 201 Although the SECL of zebra finches based on the syllable classification via ABCD-VAE was four and shorter than that
 202 of Bengalese finches, the difference between the model predictions based on the truncated and full contexts became
 203 smaller as the number of syllable categories increased (Fig. 4C; $p < 0.001$ according to the linear regression of the loss
 204 difference on the number of syllable categories and the length of truncated contexts, both in the log scale).

205 Discussion

206 This study assessed the context dependency in Bengalese finch’s song to investigate how long individual birds must
 207 remember their previous vocal outputs to generate well-formed songs. We addressed this question by fitting a state-of-
 208 the-art language model, Transformer, to the syllable sequences, and evaluating the decline in the model’s performance
 209 upon truncation of the context. We also proposed an end-to-end clustering method of Bengalese finch syllables, the
 210 ABCD-VAE, to obtain discrete inputs for the language model. In the section below, we discuss the results of this
 211 syllable clustering and then move to consider context dependency.

212 *Clustering of syllables*

213 The clustering of syllables into discrete categories played an essential role in our analysis of context dependency in
214 Bengalese finch songs, particularly for the comparison to human language in text. Various studies have observed how
215 fundamental the classification of voice elements is to animal vocalization [42, 43, 7, 11, 44, 18].

216 Our syllable clustering is based on the ABCD-VAE [27] and features the following advantages over previous approaches.
217 First, the ABCD-VAE works in a completely unsupervised fashion. The system finds a classification of syllables from
218 scratch instead of generalizing manual labeling of syllables by human annotators [30]. Thus, the obtained results are
219 more objective and reproducible [45]. Second, the ABCD-VAE automatically detects the number of syllable categories
220 in a statistically grounded way (following the Bayesian optimality under the Dirichlet prior) rather than pushing syllables
221 into a pre-specified number of classes [46, 28, 29]. This update is of particular importance when we know little about the
222 ground truth classification—as in the cases of animal song studies—and need a more non-parametric analysis. Third, the
223 ABCD-VAE adopted the speaker-normalization technique used for human speech analysis and finds individual-invariant
224 categories of syllables [28, 29]. Finally, the end-to-end clustering by the ABCD-VAE is more statistically principled
225 than the previous two-step approach—acoustic feature extraction followed by clustering—because the distinct feature
226 extractors are not optimized for clustering and the clustering algorithms are often blind to the optimization objective of
227 the feature extractors [25, 26]. We consider that such a mismatch led the combination of Gauss-VAE and GMM to
228 detect greater numbers of syllable categories than the ABCD-VAE and manual annotations, even when the clustering
229 was specialized for each individual bird and not disturbed by individual variations (see Table 1). Chorowski et al. [29]
230 also showed that a similar end-to-end clustering is better at finding speaker-invariant categories in human speech than
231 the two-step approach.

232 We acknowledge that discrete representation of data is not the only way of removing individual variations; previous
233 studies have also explored individual normalization on continuous-valued features using deep neural networks. Varia-
234 tional fair autoencoders (VFAE), for example, use speaker embeddings as background information of VAE (in both the
235 encoder and decoder while the ABCD-VAE only fed the speaker information to the decoder) [47]. As the authors note,
236 however, the use of background information does not completely remove individual variations in the extracted features
237 because continuous-valued features can distinguish infinitely many patterns (in principle) and do not have a strong
238 bottleneck effect like discrete categories, making V(F)AE lose motivation to remove individual variations from the
239 features (see also our supporting information S1.4). Accordingly, VFAE has another learning objective that minimizes
240 distances between feature vectors averaged within each speaker. More recently, researchers started to use adversarial
241 training to remove individual and other undesirable variations [48]. In adversarial training, an additional classifier
242 module is installed in the model, and that classifier attempts to *identify* the individual from the corresponding feature
243 representation. The rest of the model is trained to *deceive* the individual classifier into misclassification by anonymizing
244 the encoded features. Both VFAE and adversarial training are compatible with the ABCD-VAE and future studies may
245 combine these methods to achieve stronger speaker-normalization effects. Note, however, that those normalization
246 techniques would not yield speaker-invariant categories if there are no such categories; different individuals may
247 exhibit completely different syllable repertoires and force alignment across individuals can be inappropriate in such
248 cases. Specifically, we suspect that simply adopting other normalization methods would not lead to a more reliable
249 classification of zebra finch syllables modulo speaker variations, unless we find more appropriate segmentation.

250 It should be noted that the classical manual classification of animal voice was often based on *visual* inspection on the
251 waveforms and/or spectrograms rather than auditory inspection [42, 9, 30]. Similarly, previous VAE analyses of animal
252 voice often used a convolutional neural network that processed spectrograms as images of a fixed size [25, 26]. By
253 contrast, the present study adopted a RNN [49] to process syllable spectra frame by frame as time series data. Owing to
254 the lack of ground truth as well as empirical limitations on experimental validation, it is difficult to adjudicate on the
255 best neural network architecture for auto-encoding Bengalese finch syllables and other animals' voice. Nevertheless,
256 RNN deserves close attention as a neural/cognitive model of vocal learning. There is a version of RNN called *reservoir*
257 *computer* that has been developed to model computations in cortical microcircuits [50, 51]. Future studies may replace
258 the LSTM in the ABCD-VAE with a reservoir computer to build a more biologically plausible model of vocal learning
259 [52]. Similarly, we may filter some frequency bands in the input sound spectra to simulate the auditory perception of
260 the target animal [29], and/or adopt more anatomically/bio-acoustically realistic articulatory systems for the decoder
261 module [53]. Such Embodied VAEs would allow constructive investigation of vocal learning beyond mere acoustic
262 analysis.

263 A visual inspection of classification results shows that the ABCD-VAE can discover individual-invariant categories
264 of the Bengalese finch syllables (Fig. 2), which was also supported by their alignment with human annotations and
265 low individuality in the classified syllables (Table 1). This speaker-normalization effect is remarkable because the
266 syllables exhibit notable individual variations in the continuous feature space mapped into by the canonical VAE and
267 cross-individual clustering is difficult there [25, 26, 54]. Previous studies on Bengalese finch and other songbirds

268 often assigned distinct sets of categories to syllables of different individuals, presumably because of similar individual
269 variations in the feature space they adopted [9, 11, 30, 44].

270 By contrast, speaker-normalized clustering of zebra finch syllables was less successful, as evidenced by the lower
271 classification probability (Fig. 3B) and consistency with speaker-specific manual annotations (Table 1) than that of
272 Bengalese finch syllables. A visual inspection of category-mate syllables across individuals suggests that one major
273 challenge for finding individual-invariant categories is the complex syllables that exhibit multiple elements, or ‘notes’,
274 without clear silent intervals (gaps; Fig. 3A). Such complex syllables may be better analyzed by segmenting them
275 into smaller vocal units [55, 12, 56], and the prerequisite for appropriate voice segmentation is a major limitation of
276 the proposed method because the unclarity of segment boundaries in low-level acoustic spaces is a common problem
277 in analyses of vocalization, especially of mammals’ vocalization [44], including human speech [57, 58]. A possible
278 solution to this problem (in accordance with our end-to-end clustering) is to categorize sounds frame by frame (e.g., by
279 spectrum and MFCCs) and merge contiguous classmate frames to define a syllable-like span [29, 27, 59, 60].

280 *Context dependency*

281 According to our analysis of context dependency, Bengalese finches are expected to keep track of up to eight previously
282 uttered syllables—not just one or two—during their singing. This is evidenced by the relatively poor performance of
283 the song simulator conditioned on the truncated context of one to seven syllables compared to the full-context condition.
284 Similarly, we estimated that the production of zebra finch’s songs is dependent on four previously uttered syllables.
285 Our findings add a new piece of evidence for long context dependency in Bengalese finch songs found in previous
286 studies. Katahira et al. [9] showed that the dependent context length was at least two. They compared the first order and
287 second order Markov models, which can only access the one and two preceding syllable(s), respectively, and found
288 significant differences between them. A similar analysis was performed on canary songs by Markowitz et al. [11], with
289 an extended Markovian order (up to seventh). The framework in these studies cannot scale up to assess longer context
290 dependency owing to the empirical difficulty of training higher-order Markov models [61, 62]. By contrast, the present
291 study exploited a state-of-the-art neural language model (Transformer) that can effectively combine information from
292 much longer contexts than previous Markovian models and potentially refer up to 900 tokens [6]. Thus, the dependency
293 length reported in this study is less likely to be upper-bounded by the model limitations and provides a more precise
294 estimation (or at least a tighter lower-bound) of the real dependency length in a birdsong than previous studies.

295 The long context dependency on eight previous syllables in Bengalese finch songs is also evidenced by experimental
296 studies. Bouchard and Brainard [63] found that activities of Bengalese finches’ HVC neurons in response to listening
297 to a syllable x_t encoded the probability of the preceding syllable sequence x_{t-L}, \dots, x_{t-1} (i.e., context) given x_t ,
298 or $p(x_{t-L}, \dots, x_{t-1} | x_t)$. They reported that the length L of the context encoded by HVC neurons (that exhibited
299 strong activities to the bird’s own song) reached 7–10 syllables, which is consistent with the dependency length of
300 eight syllables estimated in the present study. Warren et al. [10] also provided evidence for long context dependency
301 from a behavioral experiment. They reported that several pairs of syllable categories of Bengalese finch songs had
302 different transitional probabilities depending on whether or not the same transition pattern occurred in the previous
303 opportunity. In other words, $\mathbb{P}(B | AB \dots A\underline{}) \neq \mathbb{P}(B | AC \dots A\underline{})$ where A, B, C are distinct syllable categories,
304 the dots represent intervening syllables of an arbitrary length ($\neq A$), and the underline indicates the position of B
305 whose probability is measured. Moreover, they found that the probability of such history-dependent transition patterns
306 is harder to modify through reinforcement learning than that of more locally dependent transitions. These results are
307 consistent with our findings. It often takes more than two transitions for syllables to recur (12.24 syllables on average
308 with the SD of 11.02 according to our own Bengalese finch data, excluding consecutive repetitions); therefore, the
309 dependency on the previous occurrence cannot be captured by memorizing just one or two previously uttered syllable(s).

310 There is also a previous study that suggests a longer context dependency in Bengalese finch songs than estimated in this
311 study (i.e., $\gg 8$). Sainburg et al. [18] studied the mutual information between birdsong syllables—including Bengalese
312 finch ones—appearing at each discrete distance. They analyzed patterns in the decay of mutual information to diagnose
313 the generative model behind the birdsong data, and reported that birdsongs were best modeled by a combination of
314 a hierarchical model that is often adopted for human language sentences and a Markov process: subsequences of
315 the songs were generated from a Markov process and those subsequences were structured into a hierarchy. Mutual
316 information decayed exponentially in the local Markov domain, but the decay slowed down and followed the power-law
317 as the inter-syllable distance became large. Sainburg et al. estimated that this switch in the decay pattern occurred when
318 the inter-syllable distance was around 24 syllables. This estimated length was substantially longer than our estimated
319 context dependency on eight syllables. The difference between the two results might be attributed to several factors.
320 First, the long-distance mutual information may not be useful for the specific task of predicting upcoming syllables that
321 defined the context dependency here and in the previous studies based on language modeling. It is possible that all the
322 information necessary for the task is available locally while the mutual information does not asymptote in the local

323 domain (see S3 for concrete examples). Another possible factor responsible for the longer context dependency detected
324 by Sainburg et al. is that their primary analysis was based on long-sequence data concatenating syllables recorded in a
325 single day (amounting to 2,693–34,588 syllables, 11,985.56 on average, manually annotated with 16–26 labels per
326 individual). Importantly, they also showed that the bimodality of mutual information decay in the Bengalese finch song
327 became less clear when the analysis was performed on bouts (consisting of 8–398 syllables, 80.98 on average). Since
328 our data was more akin to the latter, potential long dependency in the hierarchical domain might be too weak to be
329 detected in the language modeling-based analysis.

330 We also found that the greater number of syllable categories is assumed, the shorter context length becomes sufficient to
331 predict upcoming syllables. We attribute this result to the minor acoustic variations among syllables that are ignored as a
332 noise in the standard clustering or manual classification but encoded in the fine-grained classifications. When predicting
333 upcoming syllables based on the fine-grained categories, the model has to identify the minor acoustic variations encoded
334 by the categories. And the identification of such minor variations improved by referring to the local context, rather than
335 syllables far apart from the prediction target. This increases the importance of the local context compared to predictions
336 of more coarse-grained categories.

337 The reported context dependency on previous syllables also has an implication for possible models of birdsong syntax.
338 Feasible models should be able to represent the long context efficiently. For example, the simplest and traditional
339 model of the birdsong and voice sequences of other animals—including human language before the deep learning
340 era—is the n -gram model, which exhaustively represents all the possible contexts of length $n - 1$ as distinct conditions
341 [61, 62, 7]. This approach, however, requires an exponential number of contexts to be represented in the model. In
342 the worst case, the number of possible contexts in Bengalese finch songs is $37^8 = 3,512,479,453,921$ when there
343 are 37 syllable types and the context length is eight as detected in this study. While the effective context length
344 can be shortened if birds had a larger vocabulary size, the number of logically possible contexts remains huge (e.g.,
345 $160^5 = 104,857,600,000$). Such an exhaustive representation is not only hard to store and learn—for both real birds
346 and simulators—but also uninterpretable to researchers. Thus, a more efficient representation of the context syllables
347 is required [64]. Katahira et al. [9] assert that the song syntax of the Bengalese finch can be better described with
348 a lower-order hidden Markov model [65] than the n -gram model. Moreover, hierarchical language models used in
349 computational linguistics (e.g., probabilistic context-free grammar) are known to allow a more compact description
350 of human language [66] and animal voice sequences [67] than sequential models like HMM. Another compression
351 possibility is to represent consecutive repetitions of the same syllable categories differently from transitions between
352 heterogeneous syllables [16, 17] (see also [68] for neurological evidence for different treatments of heterosyllabic
353 transitions and homosyllabic repetitions). This idea is essentially equivalent to the run length encoding of digital signals
354 (e.g., AAABBCDDEEEEEE can be represented as 3A2B1C2D5E where the numbers count the repetitions of the following
355 letter) and is effective for data including many repetitions like Bengalese finch’s song. For the actual implementation in
356 birds’ brains, the long contexts can be represented in a distributed way [69]: Activation patterns of neuronal ensemble
357 can encode a larger amount of information than the simple sum of information representable by individual neurons, as
358 demonstrated by the achievements of artificial neural networks [50, 51, 70].

359 We conclude the present paper by noting that the analysis of context dependency via neural language modeling is not
360 limited to Bengalese/zebra finch’s song. Since neural networks are universal approximators and potentially fit to any
361 kind of data [71, 72], the same analytical method is applicable to other animals’ voice sequences [42, 11, 67], given
362 reasonable segmentation and classification of sequence components like syllables. Moreover, the analysis of context
363 dependency can also be performed in principle on other sequential behavioral data besides vocalization, including dance
364 [73, 74] and gestures [75, 76]. Hence, our method provides a crossmodal research paradigm for inquiry into the effect
365 of past behavioral records on future decision making.

366 **Materials and methods**

367 *Recording and preprocessing*

368 We used the same recordings of Bengalese finch songs that were originally reported in our earlier studies [30, 31]. The
369 data were collected from 18 Bengalese finches, each isolated in a birdcage placed inside a soundproof chamber. All the
370 birds were adult males (>140 days after hatching). All but two birds were obtained from commercial breeders, and the
371 other two birds (bird ID: b10 and b20) were raised in laboratory cages. Note that one bird (b20) was a son of another
372 (b03), and learned its song from the father bird. No other birds had any explicit family relationship. The microphone
373 (Audio-Technica PRO35) was installed above the birdcages. The output of the microphone was amplified using a mixer
374 (Mackie 402-VLZ3) and digitized through an audio interface (Roland UA-1010/UA-55) at 16-bits with a sampling
375 rate of 44.1 kHz. The recordings were then down-sampled to 32 kHz [30, 31]. Recording process was automatically

376 started upon detection of vocalization and terminated when no voice was detected for 500–1000 msec (the threshold
377 was adjusted for individual birds). Thus, the resulting recordings roughly corresponded to bout-level sequences, and we
378 used them as the sequence unit for the analysis of context dependency.

379 An additional dataset for song recordings of 20 zebra finches were kindly provided by Prof. Kazuhiro Wada (Hokkaido
380 University). The recording was performed in the same procedure as previously reported [77, 78].

381 Song syllables were segmented from the continuous recordings using the thresholding algorithm proposed in the
382 previous studies [30, 31]. The original waveforms were first bandpass-filtered at 1–8 kHz. Then, we obtained their
383 amplitude envelope via full-wave rectification and lowpass-filtered it at 200 Hz. Syllable onsets and offsets were
384 detected by thresholding this amplitude envelope at a predefined level, which was set at 6–10 SD above the mean of
385 the background noise level (the exact coefficient of the SD was adjusted for individual birds). The mean and SD of
386 background noise were estimated from the sound level histogram. Sound segments detected from this thresholding
387 algorithm were sometimes too close to their neighbors (typically separated by a <5 msec interval), and such coalescent
388 segments were reidentified as a single syllable, by lower-bounding possible inter-syllable gaps at 3–13 msec for
389 Bengalese finches and 3–10 msec for zebra finches (both adjusted for individual birds). Finally, extremely short sound
390 segments were discarded as noise, by setting a lower bound on possible syllable durations at 10–30 ms for Bengalese
391 finches and 5–30 msec for zebra finches (adjusted for individual birds). These segmentation processes yielded 465,310
392 Bengalese finch syllables (≈ 10.79 hours) and 237,610 zebra finch syllables (≈ 7.72 hours) in total.

393 *Clustering of syllables*

394 To perform an analysis parallel to the discrete human language data, we classified the segmented syllables into discrete
395 categories in an unsupervised way. Specifically, we used an end-to-end clustering method, named the seq2seq ABCD-
396 VAE, that combined (i) neural network-based extraction of syllable features and (ii) Bayesian classification, both of
397 which worked in an unsupervised way (i.e., without top-down selection of acoustic features or manual classification
398 of the syllables). This section provides an overview of our method, with a brief, high-level introduction to the two
399 components. Interested readers are referred to S1 in the supporting information, where we provide more detailed
400 information. One of the challenges to clustering syllables is their variable duration as many of the existing clustering
401 methods require their input to be a fixed-dimensional vector. Thus, it is convenient to represent the syllables in such a
402 format [79, 80]. Previous studies on animal vocalization often used acoustic features like syllable duration, mean pitch,
403 spectral entropy/shape (centroid, skewness, etc.), mean spectrum/cepstrum, and/or Mel-frequency cepstral coefficients at
404 some representative points for the fixed-dimensional representation [9, 30, 67]. In this study, we took a non-parametric
405 approach based on a sequence-to-sequence (seq2seq) autoencoder [81]. The seq2seq autoencoder is a RNN that first
406 reads the whole spectral sequence of an input syllable frame by frame (*encoding*; the spectral sequence was obtained
407 by the short-term Fourier transform with the 8 msec Hanning window and 4 msec stride), and then reconstructs the
408 input spectra (*decoding*; see the schematic diagram of the system provided in the upper half of Fig. 1B). Improving
409 the precision of this reconstruction is the training objective of the seq2seq autoencoder. For successful reconstruction,
410 the RNN must store the information about the entire syllable in its internal state—represented by a fixed-dimensional
411 vector—when it transitions from the encoding phase to the decoding phase. And this internal state of the RNN served
412 as the fixed-dimensional representation of the syllables. We implemented the encoder and decoder RNNs by the LSTM
413 [49].

414 One problem with the auto-encoded features of the syllables is that the encoder does not guarantee their interpretability.
415 The only thing the encoder is required to do is push the information of the entire syllables into fixed-dimensional vectors,
416 and the RNN decoder is so flexible that it can map two neighboring points in the feature space to completely different
417 sounds. A widely adopted solution to this problem is to introduce Gaussian noise to the features, turning the network
418 into the *variational* autoencoder [24, 81, 82]. Abstracting away from the mathematical details, the Gaussian noise
419 prevents the encoder from representing two dissimilar syllables close to each other. Otherwise, the noisy representation
420 of the two syllables will overlap and the decoder cannot reconstruct appropriate sounds for each.

421 The Gaussian VAE represents the syllables as real-valued vectors of an arbitrary dimension, and researchers need to
422 apply a clustering method to these vectors in order to obtain discrete categories. This two-step analysis has several
423 problems:

- 424 i The VAE is not trained for the sake of clustering, and the entire distribution of the encoded features may not
425 be friendly to existing clustering methods.
- 426 ii The encoded features often include individual differences and do not exhibit inter-individually clusterable
427 distribution (see Figure 2D and the supporting information S1.4).

428 To solve these problems, this study adopted the ABCD-VAE, which encoded data into discrete categories with a
429 categorical noise under the Dirichlet prior, and performed end-to-end clustering of syllables within the VAE (Fig. 1B).
430 The ABCD-VAE married discrete autoencoding techniques [46, 28, 29] and the Bayesian clustering popular in
431 computational linguistics and cognitive science [35, 36]. It has the following advantages over the Gaussian VAE +
432 independent clustering (whose indices, except iii, correspond to the problems with the Gaussian VAE listed above):

433 i Unlike the Gaussian VAE, the ABCD-VAE includes clustering in its learning objective, aiming at statistically
434 grounded discrete encoding of the syllables.

435 ii The ABCD-VAE can exploit a speaker-normalization technique that has proven effective for discrete VAEs:
436 The “Speaker Info.” is fed directly to the decoder (Fig. 1B), and thus individual-specific patterns need not be
437 encoded in the discrete features [28, 29].

438 iii Thanks to the Dirichlet prior, the ABCD-VAE can detect the statistically grounded number of categories on its
439 own [32]. This is the major update from the previous discrete VAEs that eat up all the categories available
440 [46, 28, 29].

441 Note that the ABCD-VAE can still measure the similarity/distance between two syllables by the cosine similarity of
442 their latent representation immediately before the computation of the classification probability (i.e., logits).

443 The original category indices assigned by the ABCD-VAE were arbitrarily picked up from 128 possible integers and
444 not contiguous. Accordingly, the category indices reported in this paper were renumbered for better visualization.

445 *Other clustering methods*

446 Clustering results of the ABCD-VAE were evaluated in comparison with baselines and topline provided by the
447 combination of feature extraction by the Gaussian VAE [24, 25, 26] and clustering on the VAE features by GMM
448 [35, 32, 36]. The number K of GMM clusters was either predetermined or auto-detected. The former fit K multivariate
449 Gaussian distributions by the expectation maximization algorithm while the latter was implemented by Bayesian
450 inference with the Dirichlet distribution prior, approximated by mean-field variational inference. Since a single run
451 of the expectation maximization and variational inference only achieved a local optimum, the best among 100 runs
452 with random initialization was adopted as the clustering results. We used the scikit-learn implementation of GMMs
453 (GaussianMixture and BayesianGaussianMixture) [83]. The default parameter values were used unless otherwise
454 specified above.

455 In the analysis of context dependency, we obtained fine-/coarse-grained classifications of syllables based on the features
456 extracted immediately before the computation of classification logits by the ABCD-VAE. The ABCD-VAE computes
457 the classification probability based on the inner-product of those features and the reference vector of each category.
458 Thus, we can compute the similarity among syllables by their cosine in the feature space, and accordingly, we applied
459 k-means clustering on the L2-normalized features. We again adopted the scikit-learn implementation of k-means
460 clustering [83].

461 *Evaluation metrics of syllable clustering*

462 The syllable classification yielded by the ABCD-VAE was evaluated by its alignment with manual annotation by a
463 human expert. We used two metrics to score the alignment: Cohen’s Kappa coefficient [33] and homogeneity [34].
464 Cohen’s Kappa coefficient is a normalized index for the agreement rate between two classifications, and has been used
465 to evaluate syllable classifications in previous studies [9, 30]. One drawback of using this metric is that it only works
466 when the two classifications use the same set of categories. This requirement was not met in our case, as the model
467 predicted classification and human annotation had different numbers of categories, and we needed to force-align each
468 of the model-predicted categories to the most common human-annotated label to compute Cohen’s Kappa [9]. On the
469 other hand, the second metric, homogeneity, can score alignment between any pair of classifications, even with different
470 numbers of categories. Homogeneity is defined based on the desideratum that each of the predicted clusters should
471 only contain members of a single ground truth class. Mathematically, violation of this desideratum is quantified by the
472 conditional entropy of the distribution of ground truth classes \mathcal{C} given the predicted clusters \mathcal{K} :

$$\text{homogeneity}(\mathcal{C}, \mathcal{K}) := \begin{cases} 1 & H(\mathcal{C}) = 1 \\ 1 - \frac{H(\mathcal{C}|\mathcal{K})}{H(\mathcal{C})} & \text{Otherwise} \end{cases} \quad (1)$$

$$H(\mathcal{C} | \mathcal{K}) := - \sum_{k \in \mathcal{K}} \sum_{c \in \mathcal{C}} \frac{|c \cap k|}{N} \log \frac{|c \cap k|}{|k|} \quad (2)$$

$$H(\mathcal{C}) := - \sum_{c \in \mathcal{C}} \frac{|c|}{N} \log \frac{|c|}{N} \quad (3)$$

473 where N denotes the total number of data points, and $|c \cap k|$ is the number of data that belong to the ground truth
474 class c and the model-predicted category k . The non-conditional entropy $H(\mathcal{C})$ normalizes the homogeneity so that it
475 ranges between 0 and 1. As we noted in the Result section, homogeneity does not penalize overclassification, so it is
476 often combined with another evaluation metric for scoring overclassification, called completeness, and constitutes a
477 more comprehensive metric named V-measure [34]. We report the completeness and V-measure scores of the syllable
478 clustering results in the supporting information S1.5.

479 *Language modeling*

480 After the clustering of the syllables, each sequence, $\mathbf{x} := (x_1, \dots, x_T)$, was represented as a sequence of discrete
481 symbols, x_t . We performed the analysis of context dependency on these discrete data.

482 The analysis of context dependency made use of a neural language model based on the current state-of-the-art
483 architecture, Transformer [19, 6]. We trained the language model on 7,779 sequences of Bengalese finch syllables
484 (amounting to 458,753 syllables in total; see Table 3) and 11,722 sequences of zebra finch syllables (234,674 syllables
485 in total). These training data were defined by the complement of the 100 test sequences that were selected in the
486 following way so that they were long enough (i) and at least one sequence per individual singer was included (ii):

- 487 i The sequences containing 15 or more syllables were selected as the candidates.
- 488 ii For each of the 18 Bengalese finches and 20 zebra finches, one sequence was uniformly randomly sampled
489 among the candidates uttered by that finch.
- 490 iii The other 82/80 sequences were uniformly randomly sampled from the remaining candidates.

491 The training objective was to estimate the probability of the whole sequences \mathbf{x} conditioned on the information about
492 the individual s uttering \mathbf{x} : That is, $\mathbb{P}(\mathbf{x} | s)$. Thanks to the background information s , the model did not need to infer
493 the singer on its own. Hence, the estimated context dependency did not comprise the correlation among syllables with
494 individuality, which would not count as a major factor especially from a generative point of view.

495 The joint probability, $\mathbb{P}(\mathbf{x} | s)$, was factorized as $\mathbb{P}(\mathbf{x} | s) = \prod_{t=1}^T \mathbb{P}(x_t | x_1, \dots, x_{t-1}, s)$, and, the model took a form
496 of the left-to-right processor, predicting each syllable x_t conditioned on the preceding context $\langle \text{sos} \rangle, x_1, \dots, x_{t-1}$,
497 where $\langle \text{sos} \rangle$ stands for the special category marking the start of the sequence. See the supporting information S2 for
498 details on the model parameters and training procedure.

499 While the VAE training excluded incompletely recorded syllables positioned at the beginning/end of recordings, we
500 included them in the language modeling by assigning them with a distinct category. This corresponds to the replacement
501 of non-frequent words with the “unk(nown)” label in natural language processing.

502 *Measuring context dependencies*

503 After training the language model, we estimated how much of the context x_1, \dots, x_{t-1} was used effectively for the
504 model to predict the upcoming syllable x_t in the test data. Specifically, we wanted to know the longest length L of the
505 truncated context x_{t-L}, \dots, x_{t-1} such that the prediction of x_t conditioned on the truncated context was worse (with
506 at least 1% greater perplexity) than the prediction based on the full context (Fig. 4A). This context length L is called the
507 *effective context length* (ECL) of the trained language model [5].

508 One potential problem with the ECL estimation using the birdsong data was that the test data was much smaller in
509 size than the human language corpora used in the previous study. In other words, the perplexity, from which the ECL
510 was estimated, was more likely to be affected by sampling error. To obtain a more reliable result, we bootstrapped the
511 test data (10,000 samples) and used the five percentile of the bootstrapped differences between the truncated and full

512 context predictions. Note that the bootstrapping was performed *after* the predictive probability of the test syllables was
513 computed, so there was no perturbation in the available contexts or any other factors affecting the language model. We
514 call this bootstrapped version of ECL the *statistically effective context length* (SECL). It is more appropriate to estimate
515 the SECL by evaluating the same set of syllables across different lengths of the truncated contexts. Accordingly, only
516 those that were preceded by 15 or more syllables (including <so>) in the test sequences were used for the analysis
517 (4,918 syllables of Bengalese finches and 1,536 syllables of zebra finches; see Table 3).

518 Supporting information

519 **S1–3 Supplementary Methods & Discussion** Detailed information of the proposed methods and comparison of the
520 language-modeling and information-theoretic approaches to context dependency.

521 Acknowledgments

522 This study was supported by MEXT Grant-in-aid for Scientific Research on Innovative Areas #4903 (Evolinguistic;
523 JP17H06380) to HK and KO, JSPS Grant-in-Aid for Scientific Research C (JP19KT0023) to ROT, and for Early-
524 Career Scientists (21K17805) to TM, and the JST Core Research for Evolutional Science and Technology 17941861
525 (JPMJCR17A4) to HK. We deeply thank Prof. Kazuhiro Wada in Hokkaido University for providing the zebra finch
526 recording dataset. We also gratefully acknowledge the support of the Academic Center for Computing and Media
527 Studies, Kyoto University, regarding the use of their supercomputer system.

References

- [1] Friston K. Learning and inference in the brain. *Neural Networks*. 2003;16(9):1325–1352. doi:<https://doi.org/10.1016/j.neunet.2003.06.005>.
- [2] Friston K. The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*. 2010;11:127–138.
- [3] Chomsky N. *Syntactic Structures*. The Hague: Mouton and Co.; 1957.
- [4] Larson B. *Long Distance Dependencies*; 2017. Oxford Bibliographies.
- [5] Khandelwal U, He H, Qi P, Jurafsky D. Sharp Nearby, Fuzzy Far Away: How Neural Language Models Use Context. In: *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Melbourne, Australia: Association for Computational Linguistics; 2018. p. 284–294. Available from: <https://www.aclweb.org/anthology/P18-1027>.
- [6] Dai Z, Yang Z, Yang Y, Carbonell J, Le Q, Salakhutdinov R. Transformer-XL: Attentive Language Models beyond a Fixed-Length Context. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics; 2019. p. 2978–2988.
- [7] Hosino T, Okanoya K. Lesion of a higher-order song nucleus disrupts phrase level complexity in Bengalese finches. *Neuroreport*. 2000;11(10):2091–2095.
- [8] Okanoya K. Song syntax in Bengalese finches: proximate and ultimate analyses. *Advances in the Study of Behavior*. 2004;34:297–345.
- [9] Katahira K, Suzuki K, Okanoya K, Okada M. Complex Sequencing Rules of Birdsong Can be Explained by Simple Hidden Markov Processes. *PLoS ONE*. 2011;6(9):1–9. doi:10.1371/journal.pone.0024516.
- [10] Warren TL, Charlesworth JD, Tumer EC, Brainard MS. Variable sequencing is actively maintained in a well learned motor skill. *Journal of neuroscience*. 2012;32(44):15414–15425. doi:10.1523/JNEUROSCI.1254-12.2012.
- [11] Markowitz JE, Ivie E, Kligler L, Gardner TJ. Long-range Order in Canary Song. *PLOS Computational Biology*. 2013;9(5):1–12. doi:10.1371/journal.pcbi.1003052.
- [12] Berwick RC, Okanoya K, Beckers GJL, Bolhuis JJ. Songs to syntax: the linguistics of birdsong. *Trends in Cognitive Science*. 2011;15(3):113–121.
- [13] Kuypers HGJM. Corticobulbar connexions to the pons and lower brain-stem in man: an anatomical study. *Brain*. 1958;81(3):364–388. doi:10.1093/brain/81.3.364.

- [14] Wild JM, Li D, Eagleton C. Projections of the dorsomedial nucleus of the intercollicular complex (DM) in relation to respiratory-vocal nuclei in the brainstem of pigeon (*Columba livia*) and zebra finch (*Taeniopygia guttata*). *Journal of Comparative Neurology*. 1997;377(3):392–413. doi:10.1002/(SICI)1096-9861(19970120)377:3<392::AID-CNE7>3.0.CO;2-Y.
- [15] Prather JF, Peters S, Nowicki S, Mooney R. Precise auditory–vocal mirroring in neurons for learned vocal communication. *Nature*. 2008;451(7176):305–310. doi:10.1038/nature06492.
- [16] Jin DZ, Kozhevnikov AA. A Compact Statistical Model of the Song Syntax in Bengalese Finch. *PLOS Computational Biology*. 2011;7(3):1–19. doi:10.1371/journal.pcbi.1001108.
- [17] Kershenbaum A, Bowles AE, Freeberg TM, Jin DZ, Lameira AR, Bohn K. Animal vocal sequences: not the Markov chains we thought they were. *Proceedings of the Royal Society of London B: Biological Sciences*. 2014;281(1792). doi:10.1098/rspb.2014.1370.
- [18] Sainburg T, Theilman B, Thielk M, Gentner TQ. Parallels in the sequential organization of birdsong and human speech. *Nature Communications*. 2019;10(3636). doi:10.1038/s41467-019-11605-y.
- [19] Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, et al. Attention is All you Need. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc.; 2017. p. 5998–6008.
- [20] van den Oord A, Kalchbrenner N, Kavukcuoglu K. *Pixel Recurrent Neural Networks*; 2016.
- [21] van den Oord A, Dieleman S, Zen H, Simonyan K, Vinyals O, Graves A, et al.. *WaveNet: A Generative Model for Raw Audio*; 2016.
- [22] Dhariwal P, Jun H, Payne C, Kim JW, Radford A, Sutskever I. *Jukebox: A Generative Model for Music*; 2020.
- [23] Okanoya K. Language evolution and an emergent property. *Current Opinion in Neurobiology*. 2007;17(2):271–276. doi:https://doi.org/10.1016/j.conb.2007.03.011.
- [24] Kingma DP, Welling M. *Auto-Encoding Variational Bayes*; 2014. *The International Conference on Learning Representations (ICLR) 2014*.
- [25] Coffey KR, Marx RG, Neumaier JF. DeepSqueak: a deep learning-based system for detection and analysis of ultrasonic vocalizations. *Neuropsychopharmacology*. 2019;44(5):859–868. doi:10.1038/s41386-018-0303-6.
- [26] Goffinet J, Mooney R, Pearson J. Inferring low-dimensional latent descriptions of animal vocalizations. *bioRxiv*. 2019;doi:10.1101/811661.
- [27] Morita T, Koda H. Exploring TTS without T Using Biologically/Psychologically Motivated Neural Network Modules (*ZeroSpeech 2020*). In: *Proceedings of Interspeech 2020*; 2020. p. 4856–4860.
- [28] van den Oord A, Vinyals O, Kavukcuoglu K. *Neural Discrete Representation Learning*. In: Guyon I, Luxburg UV, Bengio S, Wallach H, Fergus R, Vishwanathan S, et al., editors. *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc.; 2017. p. 6306–6315.
- [29] Chorowski J, Weiss RJ, Bengio S, van den Oord A. *Unsupervised Speech Representation Learning Using WaveNet Autoencoders*. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*. 2019;27(12):2041–2053. doi:10.1109/TASLP.2019.2938863.
- [30] Tachibana RO, Oosugi N, Okanoya K. *Semi-Automatic Classification of Birdsong Elements Using a Linear Support Vector Machine*. *PLOS ONE*. 2014;9(3):1–8. doi:10.1371/journal.pone.0092584.
- [31] Tachibana RO, Koumura T, Okanoya K. Variability in the temporal parameters in the song of the Bengalese finch (*Lonchura striata* var. *domestica*). *Journal of Comparative Physiology A*. 2015;201(12):1157–1168. doi:10.1007/s00359-015-1046-z.
- [32] Bishop CM. *Pattern recognition and machine learning*. Information science and statistics. New York: Springer; 2006.
- [33] Cohen J. A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*. 1960;20(1):37–46. doi:10.1177/001316446002000104.
- [34] Roseberg A, Hirschberg J. V-Measure: A conditional entropy-based external cluster evaluation measure. In: *Proceedings of the 2007 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning (EMNLP-CoNLL)*. Association for Computational Linguistics; 2007. p. 410–420.
- [35] Anderson JR. *The adaptive character of thought*. Studies in cognition. Hillsdale, NJ: L. Erlbaum Associates; 1990.
- [36] Feldman NH, Goldwater S, Griffiths TL, Morgan JL. A Role for the Developing Lexicon in Phonetic Category Acquisition. *Psychological Review*. 2013;120(4):751–778.

- [37] Landis JR, Koch GG. The Measurement of Observer Agreement for Categorical Data. *Biometrics*. 1977;33(1):159–174.
- [38] Hochreiter S, Bengio Y, Frasconi P, Schmidhuber J. Gradient Flow in Recurrent Nets: The Difficulty of Learning Long-Term Dependencies. In: Kolen JF, Kremer SC, editors. *A Field Guide to Dynamical Recurrent Networks*. Wiley-IEEE Press; 2001. p. 237–243.
- [39] Abnar S, Zuidema W. Quantifying Attention Flow in Transformers. In: *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Online: Association for Computational Linguistics; 2020. p. 4190–4197.
- [40] Manning CD, Clark K, Hewitt J, Khandelwal U, Levy O. Emergent linguistic structure in artificial neural networks trained by self-supervision. *Proceedings of the National Academy of Sciences*. 2020;117(48):30046–30054. doi:10.1073/pnas.1907367117.
- [41] Khalighinejad B, Cruzatto da Silva G, Mesgarani N. Dynamic Encoding of Acoustic Features in Neural Responses to Continuous Speech. *Journal of Neuroscience*. 2017;37(8):2176–2185. doi:10.1523/JNEUROSCI.2383-16.2017.
- [42] Payne RS, McVay S. Songs of Humpback Whales. *Science*. 1971;173(3997):585–597. doi:10.1126/science.173.3997.585.
- [43] Seyfarth R, Cheney D, Marler P. Monkey responses to three different alarm calls: evidence of predator classification and semantic communication. *Science*. 1980;210(4471):801–803. doi:10.1126/science.7433999.
- [44] Kershenbaum A, Blumstein DT, Roch MA, Akçay Ç, Backus G, Bee MA, et al. Acoustic sequences in non-human animals: a tutorial review and prospectus. *Biological Reviews*. 2016;91(1):13–52. doi:10.1111/brv.12160.
- [45] Janik VM. Pitfalls in the categorization of behaviour: a comparison of dolphin whistle classification methods. *Animal Behaviour*. 1999;57(1):133–143. doi:https://doi.org/10.1006/anbe.1998.0923.
- [46] Jang E, Gu S, Poole B. Categorical Reparameterization with Gumbel-Softmax. In: *Proceedings of the 5th International Conference on Learning Representations (ICLR)*; 2017.
- [47] Louizos C, Swersky K, Li Y, Welling M, Zemel RS. The Variational Fair Autoencoder. In: Bengio Y, LeCun Y, editors. *Proceedings of the 4th International Conference on Learning Representations (ICLR)*; 2016.
- [48] Ganin Y, Ustinova E, Ajakan H, Germain P, Larochelle H, Laviolette F, et al. Domain-Adversarial Training of Neural Networks. *Journal of Machine Learning Research*. 2016;17(59):1–35.
- [49] Hochreiter S, Schmidhuber J. Long short-term memory. *Neural Computation*. 1997;9(8):1735–1780.
- [50] Maass W, Natschläger T, Markram H. Real-Time Computing Without Stable States: A New Framework for Neural Computation Based on Perturbations. *Neural Computation*. 2002;14(11):2531–2560. doi:10.1162/089976602760407955.
- [51] Jaeger H, Haas H. Harnessing Nonlinearity: Predicting Chaotic Systems and Saving Energy in Wireless Communication. *Science*. 2004;304(5667):78–80. doi:10.1126/science.1091277.
- [52] Dehaene S, Changeux JP, Nadal JP. Neural networks that learn temporal sequences by selection. *Proceedings of the National Academy of Sciences*. 1987;84(9):2727–2731. doi:10.1073/pnas.84.9.2727.
- [53] Wang X, Takaki S, Yamagishi J. Neural Source-Filter Waveform Models for Statistical Parametric Speech Synthesis. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*. 2020;28:402–415. doi:10.1109/TASLP.2019.2956145.
- [54] Sainburg T, Thielk M, Gentner TQ. Latent space visualization, characterization, and generation of diverse vocal communication signals. *bioRxiv*. 2019;doi:10.1101/870311.
- [55] Williams H, Staples K. Syllable chunking in zebra finch (*Taeniopygia guttata*) song. *Journal of Comparative Psychology*. 1992;106(3):278–286. doi:10.1037/0735-7036.106.3.278.
- [56] Lachlan RF, van Heijningen CAA, ter Haar SM, ten Cate C. Zebra Finch Song Phonology and Syntactical Structure across Populations and Continents—A Computational Comparison. *Frontiers in Psychology*. 2016;7:980. doi:10.3389/fpsyg.2016.00980.
- [57] Chiu CC, Sainath TN, Wu Y, Prabhavalkar R, Nguyen P, Chen Z, et al.. *State-of-the-art Speech Recognition With Sequence-to-Sequence Models*; 2017.
- [58] Dunbar E, Cao XN, Benjumea J, Karadayi J, Bernard M, Besacier L, et al. The Zero Resource Speech Challenge 2017. In: *2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*; 2017. p. 323–330.
- [59] van Niekerk B, Nortje L, Kamper H. Vector-quantized neural networks for acoustic unit discovery in the ZeroSpeech 2020 challenge. In: *Proceedings of Interspeech 2020*; 2020. p. 4836–4840.

- [60] Baevski A, Hsu WN, Conneau A, Auli M. Unsupervised Speech Recognition; 2021.
- [61] Katz SM. Estimation of Probabilities from Sparse Data for the Language Model Component of a Speech Recognizer. *IEEE Transactions on Acoustics, Speech, and Signal Processing*. 1987;35(3):400–401.
- [62] Kneser R, Ney H. Improved Backing-off for N-gram Language Modeling. In: *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal*. vol. 1; 1995. p. 181–184.
- [63] Bouchard KE, Brainard MS. Neural Encoding and Integration of Learned Probabilistic Sequences in Avian Sensory-Motor Circuitry. *Journal of Neuroscience*. 2013;33(45):17710–17723. doi:10.1523/JNEUROSCI.2181-13.2013.
- [64] Morita T, Koda H. Difficulties in analysing animal song under formal language theory framework: comparison with metric-based model evaluation. *Royal Society Open Science*. 2020;7(2):192069. doi:10.1098/rsos.192069.
- [65] Rabiner LR. A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings of the IEEE*. 1989;77:257–286.
- [66] Perfors A, Tenenbaum JB, Regier T. The learnability of abstract syntactic principles. *Cognition*. 2011;118(3):306–338.
- [67] Morita T, Koda H. Superregular grammars do not provide additional explanatory power but allow for a compact analysis of animal song. *Royal Society Open Science*. 2019;6(7):190139. doi:10.1098/rsos.190139.
- [68] Fujimoto H, Hasegawa T, Watanabe D. Neural Coding of Syntactic Structure in Learned Vocalizations in the Songbird. *Journal of Neuroscience*. 2011;31(27):10023–10033. doi:10.1523/JNEUROSCI.1606-11.2011.
- [69] Nishikawa J, Okada M, Okanoya K. Population coding of song element sequence in the Bengalese finch HVC. *European Journal of Neuroscience*. 2008;27(12):3273–3283. doi:10.1111/j.1460-9568.2008.06291.x.
- [70] Nishikawa J, Okanoya K. Dynamical neural representation of song syntax in Bengalese Finch: a model study. *Ornithological Science*. 2006;5(1):95–103. doi:10.2326/osj.5.95.
- [71] Cybenko G. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and Systems*. 1989;2(4):303–314. doi:10.1007/BF02551274.
- [72] Jin L, Gupta MM, Nikiforuk PN. Universal approximation using dynamic recurrent neural networks: discrete-time version. In: *Proceedings of ICNN'95 - International Conference on Neural Networks*. vol. 1; 1995. p. 403–408.
- [73] Frith CB, Beehler BM. *The Birds of Paradise: Paradisaeidae*. *Bird Families of the World*. Oxford: Oxford University Press; 1998.
- [74] Scholes EI. Courtship Ethology of Carola's Parotia (*Parotia Carolae*). *The Auk*. 2006;123(4):967–990. doi:10.1093/auk/123.4.967.
- [75] van Lawick-Goodall J. The Behaviour of Free-living Chimpanzees in the Gombe Stream Reserve. *Animal Behaviour Monographs*. 1968;1:161–311. doi:https://doi.org/10.1016/S0066-1856(68)80003-2.
- [76] Tanner JE, Byrne RW. Representation of Action Through Iconic Gesture in a Captive Lowland Gorilla. *Current Anthropology*. 1996;37(1):162–173. doi:10.1086/204484.
- [77] Mori C, Wada K. Audition-Independent Vocal Crystallization Associated with Intrinsic Developmental Gene Expression Dynamics. *Journal of Neuroscience*. 2015;35(3):878–889. doi:10.1523/JNEUROSCI.1804-14.2015.
- [78] Hayase S, Wada K. Singing activity-driven Arc expression associated with vocal acoustic plasticity in juvenile songbird. *European Journal of Neuroscience*. 2018;48(2):1728–1742. doi:10.1111/ejn.14057.
- [79] Bellman R, Kalaba R. On adaptive control processes. *IRE Transactions on Automatic Control*. 1959;4(2):1–9. doi:10.1109/TAC.1959.1104847.
- [80] Levenshtein VI. Binary Codes Capable of Correcting Deletions, Insertions and Reversals. *Soviet Physics Doklady*. 1966;10(8):707–710.
- [81] Bowman SR, Vilnis L, Vinyals O, Dai A, Jozefowicz R, Bengio S. Generating Sentences from a Continuous Space. In: *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*; 2016.
- [82] Zhao T, Zhao R, Eskenazi M. Learning Discourse-level Diversity for Neural Dialog Models using Conditional Variational Autoencoders. In: *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Association for Computational Linguistics; 2017. p. 654–664. Available from: <http://aclweb.org/anthology/P17-1061>.
- [83] Pedregosa F, Varoquaux G, Gramfort A, Michel V, Thirion B, Grisel O, et al. Scikit-learn: Machine Learning in Python. *Journal of Machine Learning Research*. 2011;12:2825–2830.