
MEASURING LONG CONTEXT DEPENDENCY IN BIRDSONG USING AN ARTIFICIAL NEURAL NETWORK WITH A LONG-LASTING WORKING MEMORY

A PREPRINT

Takashi Morita¹

Hiroki Koda¹

Kazuo Okanoya²

Ryosuke O. Tachibana^{2*}

¹Primate Research Institute, Kyoto University, JAPAN

²Center for Evolutionary Cognitive Sciences, Graduate School of Arts and Sciences, the University of Tokyo, JAPAN

May 9, 2020

ABSTRACT

The production of grammatically and semantically appropriate human language requires reference to non-trivially long history of past utterance, which is referred to as the *context dependency* of human language. Similarly, it is of particular interest to biologists how much effect past behavioral records of individual animals have on their future behavioral decisions. In particular, birdsong serves a representative case to study context dependency in sequential signals produced by animals. Previous studies have suggested that the songs of Bengalese finches (*Lonchura striata* var. *domestica*) exhibited a long dependency on previous outputs, while their estimates were upper-bounded by methodological limitations at that time. This study newly estimated the context dependency in Bengalese finch's song in a more scalable manner using a neural network-based language model, Transformer, whose accessible context length reaches 900 tokens and is thus nearly free from model limitations, unlike the methods adopted in previous studies. A quantitative comparison with a parallel analysis of English sentences revealed that context dependency in Bengalese finch song is much shorter than that in human language but is comparable to human language syntax that excludes semantic factors of dependency. Our findings are in accordance with the previous generalization reported in related studies that birdsong is more homologous to human language syntax than the entire human language, including semantics. Thus, this study supports the hypothesis that human language modules, such as syntax and semantics, evolved from different precursors that are shared with other animals.

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

Significance Statement

- We investigated context dependency in over 10-hour recordings of Bengalese finch songs using a neural network-based language model.
- We proposed an end-to-end unsupervised clustering method of song elements (syllables) into statistically optimal categories.
- Context dependency in the birdsong is shorter than that in English sentences but comparable with the dependency in English syntax excluding semantic factors.

*Corresponding Author: rtachi@gmail.com

1 Introduction

2 Making behavioral decisions based on past information is a crucial to how humans and animals live (Friston, 2003,
3 2010; Friston and Stephan, 2007). This underscores the biological inquiry into how much effect past events have on
4 animal behaviors. These past records are not only limited to observations of the environment and other individuals
5 but also include each individual's behavioral history. A typical example is human language production where the
6 appropriate choice of words to utter depends on previously uttered words/sentences. For example, we can tell whether
7 *was* or *were* is the grammatical option after *The photographs that were taken in the cafe and sent to Mary ___* only
8 if we keep track of the previous words for a sufficient length, at least up to *photographs*, and successfully recognize
9 the two close nouns (*cafe* and *station*) as modifiers rather than the main subject. Similarly, semantically plausible
10 words are selected based on the topic of preceding sentences, as exemplified by the appropriateness of *warship* over
11 *canoe* after “missile” and “navy” are used in the same speech/document. This dependence on the production history is
12 called *context dependency* and is considered a characteristic property of human languages (Harris, 1945; Chomsky,
13 1957; Larson, 2017; Khandelwal et al., 2018; Dai et al., 2019). Birdsongs, in particular, songs by Bengalese finches
14 (*Lonchura striata* var. *domestica*), serve as a representative case study of context dependency in sequential signals
15 produced by non-human animals. Their songs are sound sequences that consist of brief vocal elements called *syllables*
16 (Hosino and Okanoya, 2000; Okanoya, 2004). Previous studies have suggested that birdsongs exhibit non-trivially
17 long dependency on previous outputs (Katahira et al., 2011; Warren et al., 2012; Markowitz et al., 2013). Complex
18 sequential patterns of syllables have been discussed in comparison with human language syntax from the viewpoint of
19 formal linguistics (Okanoya, 2004; Berwick et al., 2011, 2012; Berwick and Chomsky, 2016). Neurological studies
20 have also revealed homological network structures in the vocal production, recognition, and learning of songbirds and
21 humans (Kuypers, 1958; Wild et al., 1997; Prather et al., 2008). Assessing whether birdsongs exhibit long context
22 dependency similar to human language is an important comparative study that several previous studies have addressed
23 using computational methods. Katahira et al. (2011) found that Bengalese finch songs are dependent on more than just
24 one previously uttered syllable. Similarly, Markowitz et al. (2013) report that canary songs exhibit a dependency on
25 the last six chunks of syllables. However, the reported lengths of context dependency were measured using a limited
26 language model—Markov/*n*-gram model—that was only able to access a few recent syllables in the context by design
27 (at most two in Katahira et al. and seven in Markowitz et al.). Thus, it is unclear if those numbers were real dependency
28 lengths in the birdsongs or merely model limitations. Moreover, the use of a limited language model is problematic for
29 comparative studies because human languages are not modeled precisely by a Markov process (Chomsky, 1956; Rabin
30 and Scott, 1959).

31 Powerful language models are now available owing to recent advancements in machine learning, particularly in artificial
32 neural networks (Vaswani et al., 2017; Devlin et al., 2018; Dai et al., 2019). These neural language models are
33 flexible, can be fit to birdsong data, and importantly, can potentially refer to 200–900 syllables from the past when
34 the data include such long dependency (Khandelwal et al., 2018; Dai et al., 2019). In the present study, we assessed
35 a detailed estimate of context dependency length in the songs of Bengalese finches using the Transformer language
36 model, which can exploit the longest context among currently available models and has state-of-the-art architecture in
37 today's natural language processing (Vaswani et al., 2017; Devlin et al., 2018; Dai et al., 2019). We also performed
38 the same dependency analysis on English sentences as a baseline for human language data. On one hand, we found
39 that context dependency in Bengalese finch's song is much shorter than in English sentences. On the other hand, we
40 found that the context dependency in the birdsong is slightly longer than and more comparable to the dependency
41 in English syntax, where each word in the sentence is replaced with its grammatical category—such as NOUN and
42 VERB—and semantic information included in the word is removed. These findings corroborate previous generalizations
43 in comparative studies that birdsongs are more homologous to human language syntax than the entirety of human
44 language including semantics (Berwick et al., 2011; Gibson and Tallerman, 2012; Miyagawa et al., 2013) and provide a
45 new piece of evidence for the hypothesis that human language modules, such as syntax and semantics, evolved from
46 different precursors that are shared with other animals (e.g., birdsongs and alarm calls respectively; Okanoya, 2007;
47 Okanoya and Merker, 2007; Miyagawa et al., 2013, 2014; Nóbrega and Miyagawa, 2015).

48 The adopted analysis of context dependency is straightforward. First, we trained the Transformer language model on
49 sequences of Bengalese finch syllables such that the model predicted upcoming syllables based on previously uttered
50 ones. The trained model can be seen as a simulator of birdsong syntax, whose accessible context can be controlled
51 by researchers, unlike real birds. Second, we measured the difference in predictive performance of the trained model
52 working on full contexts and truncated contexts as depicted in Figure 1b. Intuitively, this difference became smaller
53 as the truncated context got longer and contained more information. We found the maximum length of the truncated
54 contexts where the prediction difference was above a canonical threshold. This length, the effective context length
55 (Khandelwal et al., 2018; Dai et al., 2019), is our estimate of the context dependency in Bengalese finch songs.

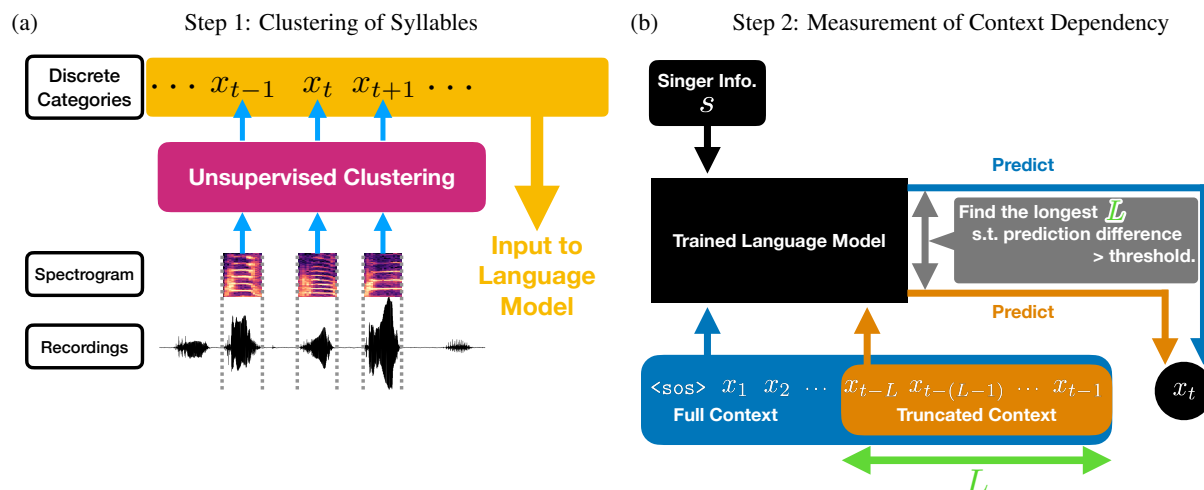


Figure 1: Schematic diagram of the proposed analysis, consisting of two steps. (a) Clustering of Bengalese finch syllables. The recordings were transformed into spectra and categorized by an unsupervised, end-to-end clustering system. The resulting sequences of discrete syllable categories were used as inputs to the language model. (b) Assessment of context dependency in the birdsong. The Transformer language model was fit to the birdsong data and its predictive performance was tested on the full and truncated contexts. The effective context length is the maximum length L of the truncated context such that the prediction difference is greater than an arbitrary threshold (= 1% in perplexity by convention).

56 The syntactic analysis of animal songs is typically performed on sequences of discretely categorized syllables (Payne
 57 and McVay, 1971; Seyfarth et al., 1980; Hosino and Okanoya, 2000; Kojima, 2003; Suzuki et al., 2006; Kakishita et al.,
 58 2007; Markowitz et al., 2013; Kershenbaum et al., 2014, 2016; Sainburg et al., 2019a, but see Katahira et al., 2011;
 59 Morita and Koda, 2019; Sainburg et al., 2019b for categorization-free approaches). The present study also followed
 60 this approach and represented each syllable with a discrete category (Figure 1a). This was useful for comparison with
 61 English text data in the analysis of context dependency. The syllable categories were discovered by a novel end-to-end
 62 clustering method that automatically classified segmented recordings—with variable duration—into an unspecified
 63 number of statistically optimal, individual-invariant categories.

64 2 Results

65 2.1 Clustering of Syllables

66 The context dependency analysis adopted herein was performed on sequences of discrete symbols such as human
 67 language sentences represented by sequences of text words. Thus, we needed a discrete representation of the syllables
 68 of Bengalese finches. The classical approach to this discretization task is based on visual and/or auditory inspection of
 69 human experts (Hosino and Okanoya, 2000; Okanoya, 2004; Kakishita et al., 2007; Katahira et al., 2011; Kershenbaum
 70 et al., 2014; Tachibana et al., 2014). More recently, several researchers explored fully unsupervised classification of
 71 animal vocalization (i.e., getting completely rid of the manual classification) based on acoustic features extracted by
 72 artificial neural networks called *variational autoencoders* (VAEs; Kingma and Welling, 2014; Coffey et al., 2019;
 73 Goffinet et al., 2019; Sainburg et al., 2019b). The present study extended this latter approach and proposed an end-
 74 to-end unsupervised clustering method named the *ABCD-VAE* (whose first four letters stand for the Attention-Based
 75 Categorical sampling with the Dirichlet prior). The proposed method automatically classifies segmented recordings—
 76 with variable duration—into an unspecified number of statistically optimal categories. It also allowed us to exploit
 77 the speaker-normalization technique developed for unsupervised learning of human language from speech recordings
 78 (van den Oord et al., 2017; Chorowski et al., 2019; Dunbar et al., 2019; Tjandra et al., 2019), yielding syllable
 79 classification modulo individual variation.

80 A high-level description of the clustering system is as follows. We trained a RNN to “hear” a syllable in its entirety
 81 (*encoding*) and reproduce it as precisely as possible (*decoding*). To accomplish this task, the RNN must store the
 82 information about the entire syllable in its internal state—represented by a fixed-dimensional vector, analogous to a
 83 fixed number of neurons exhibiting certain activation patterns in the real brain—when it transitions from the encoding
 84 phase to the decoding phase. Thus, this internal state of the RNN can be used as a fixed-dimensional representation of

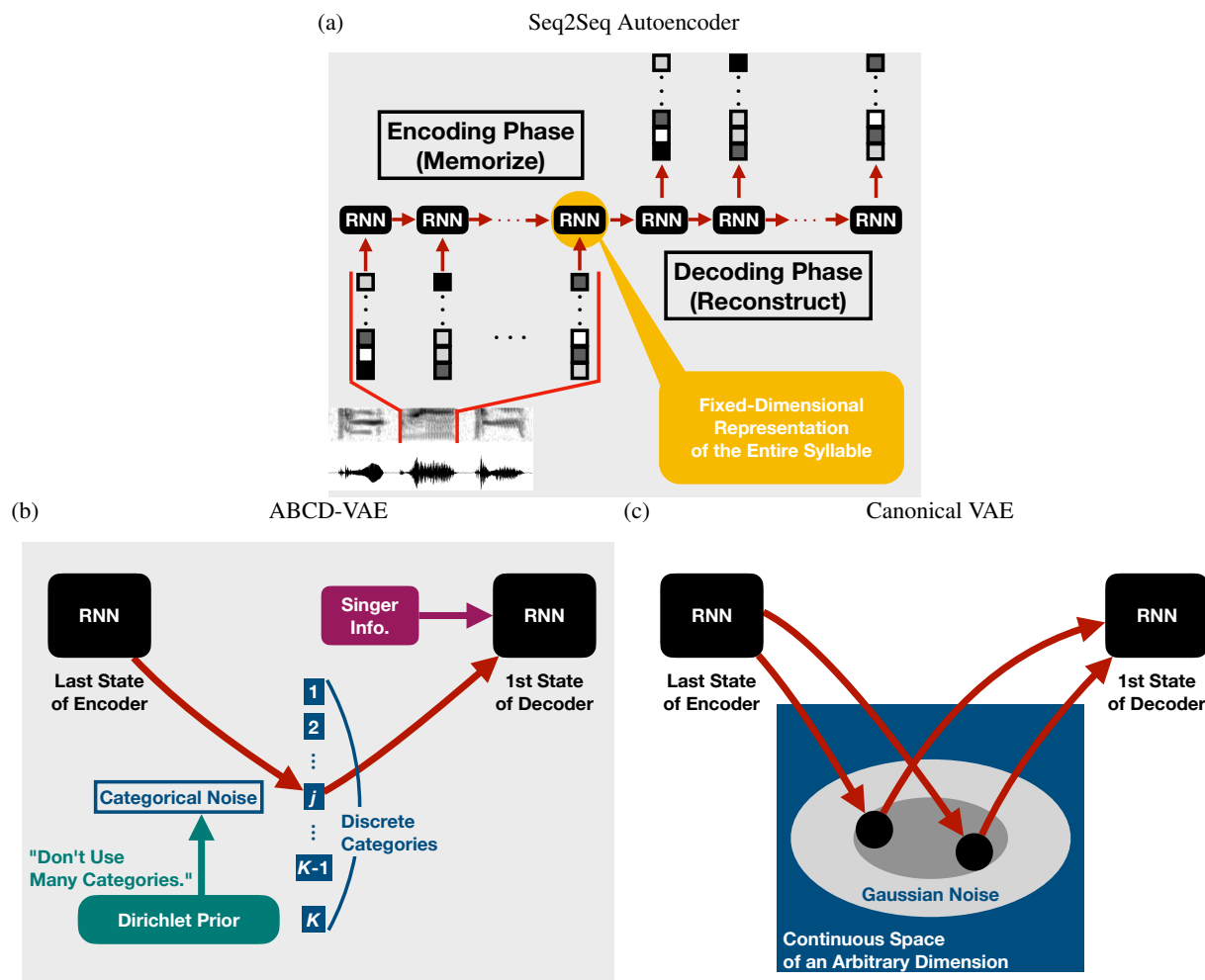


Figure 2: Schematic diagram of the proposed clustering method (combination of a and b, highlighted in light gray). (a) The high-level structure of the sequence-to-sequence (seq2seq) autoencoder. It was trained to first read syllable spectra frame by frame and then reconstruct them as precisely as possible. A fixed-dimensional representation of entire input syllables was obtained in the middle between the encoder and the decoder. (b) The ABCD-VAE that encoded syllables into discrete categories between the encoder and the decoder. A statistically optimal number of categories were detected under an arbitrarily specified upper bound thanks to the Dirichlet prior (Bishop, 2006; O’Donnell, 2015; Little, 2019). The identity of the syllable-uttering individual was informed by the decoder besides the syllable categories. Therefore, individual-specific patterns need not have been encoded in the discrete syllable representation (Jang et al., 2017; van den Oord et al., 2017; Chorowski et al., 2019). (c) The canonical VAE with the Gaussian noise. It was also inserted between the encoder and the decoder and used to increase the interpretability of the continuous-valued features.

85 the syllables. We also installed a discretization module between the encoder and the decoder, and thereby obtained
 86 syllable classifications (Figure 2b). This discretization is the major difference from the previous VAE studies on animal
 87 voices that extracted continuous-valued features from the middle point (using a Gaussian noise, Figure 2c). Moreover,
 88 the proposed method automatically detects the statistically optimal number of syllable categories (under an arbitrary
 89 upper bound, which we set at 128; owing to the Dirichlet prior; Bishop, 2006; O’Donnell, 2015; Little, 2019) and works
 90 in a non-parametric manner compared to the previous end-to-end classification (Jang et al., 2017; van den Oord et al.,
 91 2017; Chorowski et al., 2019; Tjandra et al., 2019).

92 Syllables collected from 18 adult male finches were analyzed (465, 310 syllables in total) and 39 categories—shared
 93 across the individuals—were detected (Figure 3a). The predicted classification captured major spectral patterns
 94 (Figure 3b, e). It ignored individual variations and other minor differences visible in the syllable embeddings obtained
 95 via the canonical continuous-valued VAE (Figure 3d). The unsupervised classification was mostly consistent with
 96 manual annotations assigned by a human expert (Figure 3c; Tachibana et al., 2014).

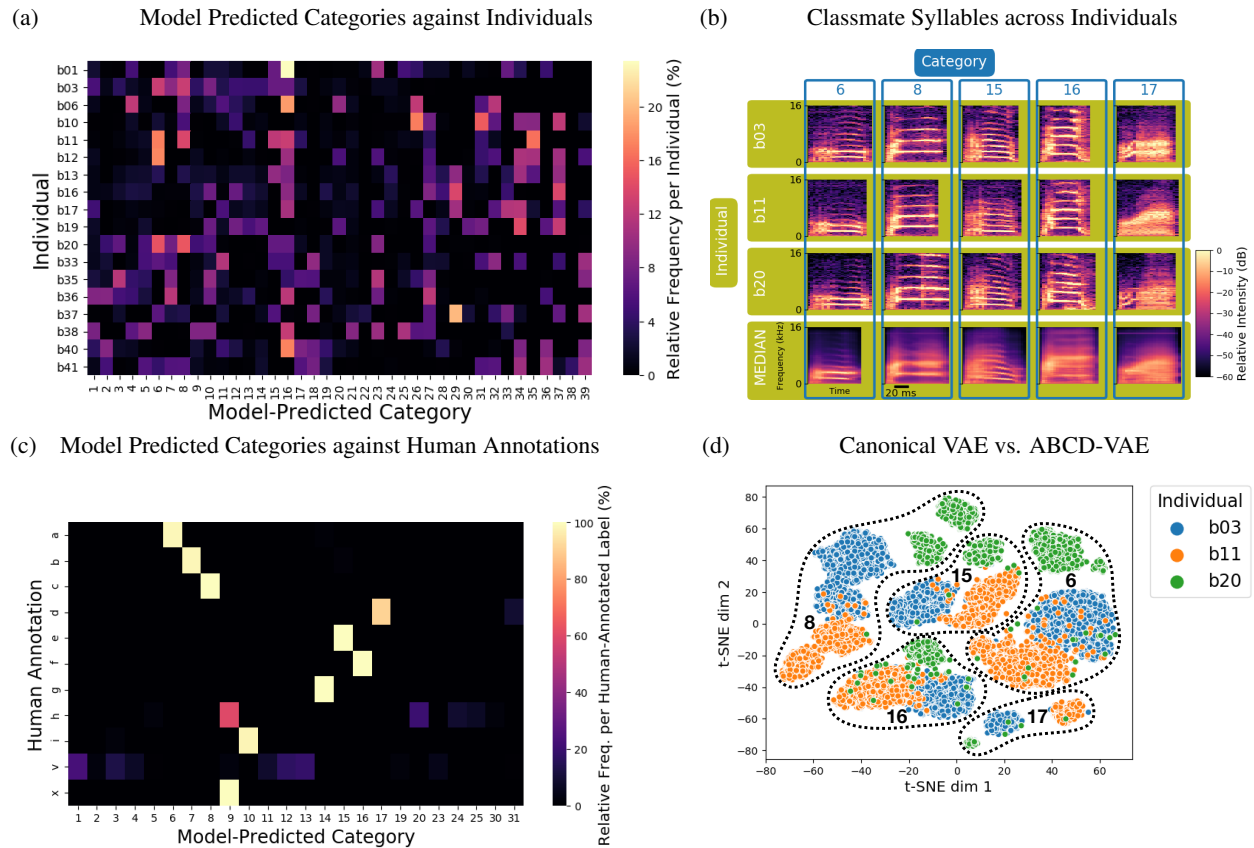
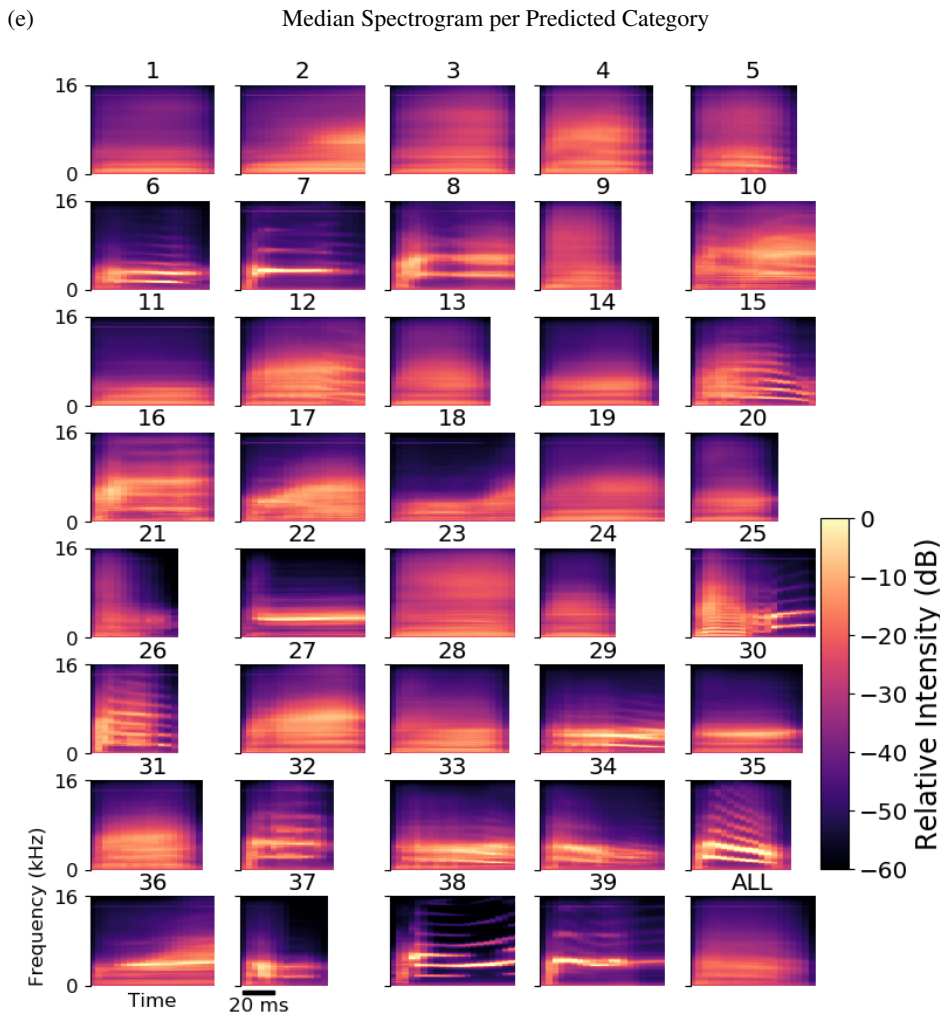


Figure 3: Clustering results of Bengalese finch syllables based on the ABCD-VAE. (a) Relative frequency of syllable categories (columns) per individual (rows). The category indices were renumbered for better visualization because the original values were arbitrarily picked from 128 possible integers and were not contiguous. (b) Syllable spectrograms and their classification across individuals. Syllables in each of the first to third rows (yellow 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 median spectrogram of each category over all the 39 individuals. The examples had the greatest classification probability (> 0.999) among the syllables of the same individual and category. (c) Relative frequency of syllable categories (columns) per label manually annotated by a human expert (Tachibana et al., 2014). Only data from a single individual (b03) were presented because the manual annotations were not shared across individuals. (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 (van der Maaten and Hinton, 2008). The continuous embeddings included notable individual variations represented by colors, whereas the ABCD-VAE classification ignored these individual variations.



97 2.2 Context Dependency

98 The classification results in the previous subsection yielded sequences of discretely represented syllables called *bouts*
99 (9,139 sequences in total). Our objective is to assess context dependency in these bouts. We measured the context
100 dependency via the Transformer language model trained on 9,039 bouts. The remaining 100 bouts were used to
101 score its predictive performance from which the dependency was calculated. The model predictions were provided
102 in the form of the log conditional probability of the test syllables given the preceding ones in the same bout: i.e.,
103 $\log \mathbb{P}(x_t | x_1, \dots)$. We compared the model predictions conditioned on the full context (x_1, \dots, x_{t-1}) and the truncated
104 context $(x_{t-L}, \dots, x_{t-1})$ and found the *statistically* effective context length (SECL) defined by the maximum length of
105 the truncated context wherein the mean prediction difference between the two contexts was significantly greater than
106 the canonical 1% threshold in perplexity (at 0.05 level of significance estimated from 10,000 bootstrapped samples;
107 Khandelwal et al., 2018).

108 In addition to the Bengalese finch bouts, we performed the same analysis on two types of English sentences for
109 comparison (12,327 training sentences and 2,006 test sentences; Silveira et al., 2014). One of them represented the
110 words by their lemma. The other English data reduced the words to the part-of-speech (PoS) categories such as NOUN
111 and VERB (cf. Perfors et al., 2011). This PoS representation removed semantic information from the English sentences;
112 thus, the reported context dependencies were free of semantic factors such as words (at distance) encoding the topic of
113 the sentences (e.g., “missile” and “warship” may co-occur in military documents and be stochastically dependent on
114 each other without a grammatical relation such as subject and object).

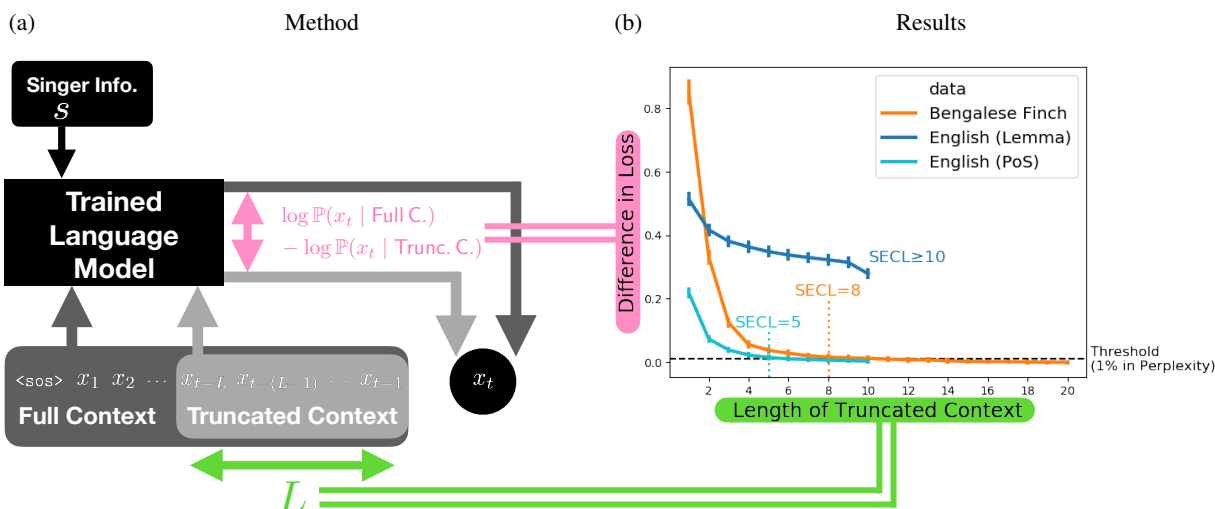


Figure 4: The differences in the mean loss (negative log probability) between the truncated- and full-context predictions. The x-axis corresponds to the length of the truncated context. The error bars show the 90% confidence intervals estimated from 10,000 bootstrapped samples. The loss difference is statistically significant if the lower side of the intervals are above the threshold indicated by the horizontal dashed line.

115 The SECL of the Bengalese finch song was eight (the orange line in Figure 4). In other words, the restriction on the
 116 available context to the one to eight preceding syllables decreased the mean predictive probability from the full-context
 117 baseline by a significantly greater amount than the threshold. However, the difference became marginal when nine
 118 or more syllables were included in the truncated context. On one hand, this is under the SECL of the English lemma
 119 data, which was ten or greater (the dark blue line, achieved the upper bound). On the other hand, the English SECL
 120 decreased to five when we represented the words by the PoS tags and removed the semantic factors (the light blue line).
 121 Hence, the context dependency in Bengalese finch songs is more comparable to that in the English syntax than in the
 122 full English including semantics.

123 3 Discussion

124 This study assessed the context dependency in Bengalese finch’s song to investigate how long individual birds must
 125 remember their previous vocal outputs to generate well-formed song bouts. We addressed this question by fitting a
 126 state-of-the-art language model, Transformer, to the bouts, and evaluating the decline in the model’s performance
 127 upon truncation of the context. We also proposed an end-to-end clustering method of Bengalese finch syllables, the
 128 ABCD-VAE, to obtain discrete inputs for the language model. In the section below, we discuss the results of this
 129 syllable clustering (§3.1) and then move to consider context dependency (§3.2).

130 3.1 Clustering of Syllables

131 The clustering of syllables into discrete categories played an essential role in our analysis of context dependency in
 132 Bengalese finch songs, particularly for the comparison to human language in text. Various studies have observed how
 133 fundamental the classification of voice elements is to animal vocalization (Payne and McVay, 1971; Seyfarth et al.,
 134 1980; Hosino and Okanoya, 2000; Kojima, 2003; Suzuki et al., 2006; Kakishita et al., 2007; Markowitz et al., 2013;
 135 Kershenbaum et al., 2016; Sainburg et al., 2019a, but see Katahira et al., 2011; Morita and Koda, 2019; Sainburg et al.,
 136 2019b for categorization-free approaches).

137 Our syllable clustering is based on the AVCD-VAE and features the following advantages over previous approaches. First,
 138 the ABCD-VAE works in a completely unsupervised fashion. The system finds the statistically optimal classification
 139 of syllables instead of generalizing manual labeling of syllables by human annotators (as opposed to Tachibana et al.,
 140 2014). Thus, the obtained results are more objective and reproducible (cf. Janik, 1999). Second, the ABCD-VAE detects
 141 the statistically optimal number of syllable categories rather than pushing syllables into a pre-specified number of
 142 classes (as opposed to Jang et al., 2017; van den Oord et al., 2017; Chorowski et al., 2019). This update is of particular
 143 importance when we know little about the ground truth classification—as in the cases of animal song studies—and need
 144 a more non-parametric analysis. Third, the ABCD-VAE adopted the speaker-normalization technique used for human

145 speech analysis and finds individual-invariant categories of syllables (van den Oord et al., 2017; Chorowski et al., 2019;
146 Tjandra et al., 2019). Finally, the end-to-end clustering by the ABCD-VAE is more optimal than the previous two-step
147 approach—acoustic feature extraction followed by clustering—because the feature extractors are not optimized for
148 clustering and the clustering algorithms are often blind to the optimization objective of the feature extractors (Coffey
149 et al., 2019; Goffinet et al., 2019; Sainburg et al., 2019b). Chorowski et al. (2019) also showed that a similar end-to-end
150 clustering is better at finding speaker-invariant categories in human speech than the two-step approach.

151 It should be noted that the classical manual classification of animal voice was often based on *visual* inspection on
152 the waveforms and/or spectrograms rather than auditory inspection (Payne and McVay, 1971; Katahira et al., 2011;
153 Tachibana et al., 2014). Similarly, previous VAE analyses of animal voice often used a convolutional neural network
154 that processed spectrograms as images of a fixed size (Coffey et al., 2019; Goffinet et al., 2019). By contrast, the present
155 study adopted a RNN (specifically, a version called the long short-term memory, abbreviated as LSTM Hochreiter
156 and Schmidhuber, 1997) to process syllable spectra frame by frame as time series data. Owing to the lack of ground
157 truth as well as empirical limitations on experimental validation, it is difficult to adjudicate on the best neural network
158 architecture for auto-encoding Bengalese finch syllables and other animals' voice. Nevertheless, RNN deserves close
159 attention as a neural/cognitive model of vocal learning. There is a version of RNN called *reservoir computer* that has
160 been developed to model computations in cortical microcircuits (Maass et al., 2002; Natschläger et al., 2003; Jaeger
161 and Haas, 2004). Future studies may replace the LSTM in the ABCD-VAE with a reservoir computer to build a more
162 biologically plausible model of vocal learning (cf. Dehaene et al., 1987). Similarly, we may filter some frequency
163 bands in the input sound spectra to simulate the auditory perception of the target animal (cf. the Mel-frequency cepstral
164 coefficients, MFCCs, are used in human speech analysis; Chung et al., 2016; Chorowski et al., 2019; Tjandra et al.,
165 2019), and/or adopt more anatomically/bio-acoustically realistic articulatory systems for the decoder module (cf. Wang
166 et al., 2020, implemented the source-filter model of vocalization based on an artificial neural network). Such Embodied
167 VAEs would allow constructive investigation of vocal learning beyond mere acoustic analysis.

168 A visual inspection of classification results shows that the ABCD-VAE can discover individual-invariant categories of
169 the Bengalese finch syllables (Figure 3). This speaker-normalization effect is remarkable because the syllables exhibit
170 notable individual variations in the continuous feature space mapped into by the canonical VAE and cross-individual
171 clustering is difficult there (see Figure 3d and the supporting information S1.4; Coffey et al., 2019; Goffinet et al.,
172 2019; Sainburg et al., 2019b). Previous studies on Bengalese finch and other songbirds often assigned distinct sets of
173 categories to syllables of different individuals, presumably because of similar individual variations in the feature space
174 they adopted (Katahira et al., 2011; Markowitz et al., 2013; Tachibana et al., 2014; Kershenbaum et al., 2016; Sainburg
175 et al., 2019b).

176 The end-to-end classification by the ABCD-VAE can be applied to a broad range of studies on animal vocalization,
177 including cases where sequential organization of voice units is not at issue. The limitations of the proposed method
178 are the prerequisite for appropriate voice segmentation as it operates on predefined time series of sound spectra, and a
179 single category is assigned to each time series. Although birdsongs often exhibit clear pauses and researchers use them
180 to define syllable boundaries, appropriate voice segmentation is not necessarily clear for other animals (Kershenbaum
181 et al., 2016; Sainburg et al., 2019b), including human speech (Chiu et al., 2017; Dunbar et al., 2017, 2019; Rao et al.,
182 2017). A possible solution to this problem (in accordance with our end-to-end clustering) is to categorize sounds
183 frame by frame (e.g., by spectrum and MFCC) and merge contiguous classmate frames to define a syllable-like span
184 (Chorowski et al., 2019; Tjandra et al., 2019).

185 **3.2 Context Dependency**

186 According to our analysis of context dependency, Bengalese finches are expected to keep track of up to eight previously
187 uttered syllables—not just one or two—during their singing. This is evidenced by the relatively poor performance of the
188 song simulator conditioned on the truncated context of one to eight syllables compared to the full-context condition. Our
189 findings add a new piece of evidence for long context dependency in Bengalese finch songs found in previous studies.
190 Katahira et al. (2011) showed that there are at least two dependent context lengths. They compared the first order and
191 second order Markov models, which can only access the one and two preceding syllable(s), respectively, and found
192 significant differences between them. A similar analysis was performed on canary songs by Markowitz et al. (2013),
193 with an extended Markovian order (up to seventh). The framework in these studies cannot scale up to assess longer
194 context dependency owing to the empirical difficulty of training higher-order Markov models (Katz, 1987; Kneser and
195 Ney, 1995; Bengio et al., 2001, 2003; Goldwater et al., 2006; Teh, 2006). By contrast, the present study exploited
196 a state-of-the-art neural language model (Transformer) that can effectively combine information from much longer
197 contexts than previous Markovian models and potentially refer up to 900 tokens (Dai et al., 2019). Thus, the dependency
198 length reported in this study is not likely to be upper-bounded by the model limitations and provides a more precise
199 estimation of the real dependency length in a birdsong than previous studies. The long context dependency in Bengalese

200 finch songs is also evidenced by experimental studies. Warren et al. (2012) reported that several pairs of syllable
201 categories had different transitional probability depending on whether or not the same transition pattern occurred in the
202 previous 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
203 categories, the dots represent intervening syllables of an arbitrary length ($\neq A$), and the underline indicates the position
204 of B whose probability is measured. They also found that the probability of such history-dependent transition patterns
205 is harder to modify through reinforcement learning than that of more locally dependent transitions. These results are
206 consistent with our findings. It often takes more than two transitions for syllables to recur (12.17 syllables on average
207 with the SD of 11.30 according to our own bout data, excluding consecutive repetitions); therefore, the dependency on
208 the previous occurrence cannot be captured by memorizing just one or two previously uttered syllable(s).

209 Our study also found that Bengalese finch songs are more comparable to human language syntax than to the entirety
210 of human language including semantics. This was demonstrated by our analysis of English sentences represented
211 by sequences of lemmas and PoS categories. While the lemma-represented English sentences exhibited long context
212 dependency beyond ten words as reported in previous studies (Khandelwal et al., 2018; Dai et al., 2019), the dependency
213 length decreased to five—below the Bengalese finch result—when the PoS representation was used and semantic
214 information was removed from the sentences. The gap between the two versions of English suggests that the major
215 factor of long-distance dependencies in human language is the semantics, not the syntax. This is consistent with previous
216 studies reporting that human language syntax prefers shorter dependency (Gibson, 1998; Futrell et al., 2015). Moreover,
217 comparative studies between birdsong and human language often argue the former’s lack of semantic function (Berwick
218 et al., 2011, 2012; Gibson and Tallerman, 2012; Miyagawa et al., 2013, 2014), without referential variations seen in
219 alarm calls (Seyfarth et al., 1980; Ouattara et al., 2009; Suzuki et al., 2016). This claim led to the hypothesis that human
220 language syntax and semantics evolved from different precursors—sequence-generating system, such as animal song,
221 and information-carrying system such as alarm calls—which were integrated to shape the entirety of human language
222 (Okanoya, 2007; Okanoya and Merker, 2007; Miyagawa et al., 2013, 2014; Nóbrega and Miyagawa, 2015). Our findings
223 are in accordance with this view, providing a novel relative similarity between birdsong and human language syntax
224 compared to the whole linguistic system. Note that this kind of direct comparative study of human language and animal
225 song was not feasible until flexible language models based on neural networks became available.

226 The reported context dependency on eight previous syllables also has an implication for possible models of Bengalese
227 finch syntax. Feasible models should be able to represent the long context efficiently. For example, the simplest and
228 traditional model of the birdsong and voice sequences of other animals—including human language before the deep
229 learning era—is the n -gram model, which exhaustively represents all the possible contexts of length $n - 1$ as distinct
230 conditions (Katz, 1987; Kneser and Ney, 1995; Hosino and Okanoya, 2000; Goldwater et al., 2006; Teh, 2006). This
231 approach, however, requires an exponential number of contexts to be represented in the model. In the worst case, the
232 number of possible contexts is $39^8 = 5,352,009,260,481$ when there are 39 syllable types and the context length is
233 eight as detected in this study. Such an exhaustive representation is not only hard to store and learn—for both real
234 birds and simulators—but also uninterpretable to researchers. Thus, a more efficient representation of the context
235 syllables is required (cf. Morita and Koda, 2020). Katahira et al. (2011) assert that the song syntax of the Bengalese
236 finch can be better described with a lower-order hidden Markov model (Rabiner, 1989; Beal et al., 2002, HMM);
237 than the n -gram model. Moreover, hierarchical language models used in computational linguistics (e.g., probabilistic
238 context-free grammar) are known to allow a more compact description of human language (Perfors et al., 2011) and
239 animal voice sequences (Morita and Koda, 2019) than sequential models like HMM. Another compression possibility is
240 to represent consecutive repetitions of the same syllable categories differently from transitions between heterogeneous
241 syllables (cf. Kershenbaum et al., 2014). This idea is essentially equivalent to the run length encoding of digital signals
242 (e.g., AAABBCDDEEEEEE can be represented as 3A2B1C2D5E where the numbers count the repetitions of the following
243 letter) and is effective for data including many repetitions like Bengalese finch’s song. For the actual implementation in
244 birds’ brains, the long contexts can be represented in a distributed way (Nishikawa et al., 2008): Activation patterns of
245 neuronal ensemble can encode a larger amount of information than the simple sum of information representable by
246 individual neurons, as demonstrated by the achievements of artificial neural networks (Bengio et al., 2001, 2003; Ryeu
247 et al., 2001; Tsuda, 2001; Maass et al., 2002; Jaeger and Haas, 2004; Nishikawa and Okanoya, 2006).

248 While this study discussed context dependency in the context of memory durability required for generating/processing
249 birdsongs (cf. Katahira et al., 2011; Warren et al., 2012; Markowitz et al., 2013), there are different definitions of
250 context dependency designed for different research purposes. Sainburg et al. (2019a) studied the *mutual information*
251 between birdsong syllables—including Bengalese finch ones—appearing at each discrete distance. Following a study
252 on human language by Lin and Tegmark (2017), Sainburg et al. analyzed patterns in the decay of mutual information to
253 diagnose the generative model behind the birdsong data, instead of addressing the question about memory. Importantly,
254 their mutual information analysis cannot replace our model-based analysis to assess the memory-oriented context
255 dependency: Mutual information is a pairwise metric of probabilistic dependence between two tokens (e.g., words in
256 human languages, syllables in birdsongs), and thus, everything in the middle is ignored. To see the problem, suppose

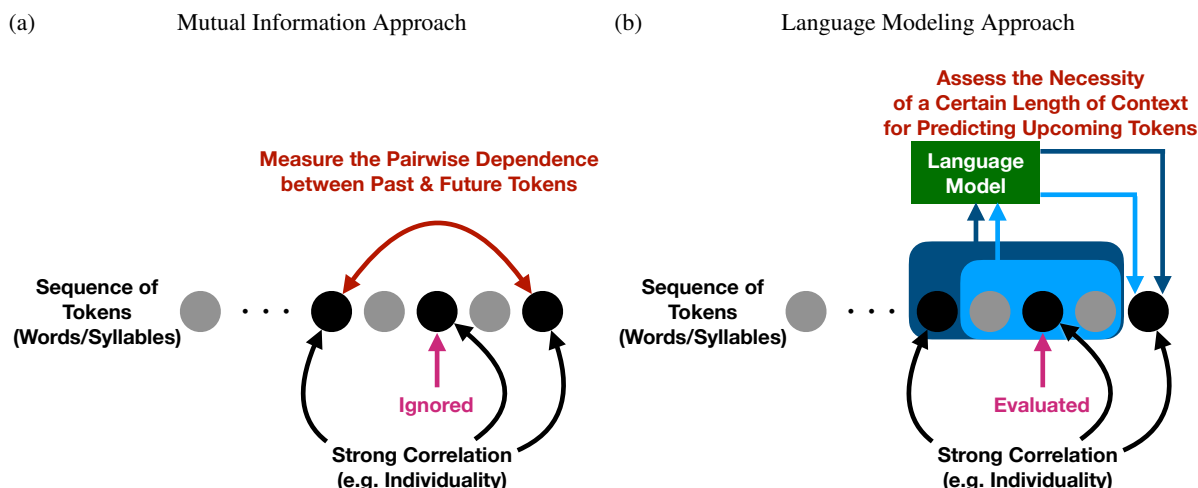


Figure 5: The analysis of context dependency based on the (a) mutual information and (b) language modeling.

257 that some tokens reflect the individuality of the speaker, as depicted in Figure 5a. (See the supporting information
258 S3.1 for a more concrete, mathematical example of this problematic situation. S3.2 introduces other examples that
259 demonstrate difficulties in the mutual information analysis.) Two occurrences of speaker-encoding tokens are dependent
260 on each other regardless of their distance if the other tokens between the two are ignored, and this pairwise dependence
261 is what mutual information accounts for. It should be clear now that such pairwise dependence does not necessarily
262 match the agent-oriented concept of context dependency as the only thing relevant to the song recognition task (or
263 speaker identification in this toy example) is the most recent occurrence of the correlating tokens. By contrast, our
264 language modeling approach captured the agent-oriented concept of context dependency as desired. Dependency on
265 a token in the past is detected if the prediction of upcoming tokens becomes notably more difficult by limiting the
266 available context to the more recent tokens (Figure 5b; Khandelwal et al., 2018; Dai et al., 2019). In other words,
267 reference to a token in the distant past is considered unnecessary if the same information (e.g., speaker identity)
268 is available from more recent tokens. Therefore, the present study complements, rather than repeats/replaces, the mutual
269 information analysis and findings from it.

270 We conclude the present paper by noting that the analysis of context dependency via neural language modeling is
271 not limited to Bengalese finch's song. Since neural networks are universal approximators and potentially fit to any
272 kind of data (Cybenko, 1989; Hornik, 1991; Jin et al., 1995; Maass et al., 2002; Lu et al., 2017), the same analytical
273 method is applicable to other animals' voice sequences (Payne and McVay, 1971; Suzuki et al., 2006; Markowitz et al.,
274 2013; Morita and Koda, 2019). Moreover, the analysis of context dependency can also be performed in principle on
275 other sequential behavioral data besides vocalization, including dance (Frith and Beehler, 1998; Scholes, 2006, 2008)
276 and gestures (van Lawick-Goodall, 1968; de Waal, 1988; Tanner and Byrne, 1996; Liebal et al., 2006). Hence, our
277 method provides a crossmodal research paradigm for inquiry into the effect of past behavioral records on future decision
278 making.

279 4 Materials & Methods

280 4.1 Recording and Segmentation of Bengalese Finch's Song

281 We used the same recordings of Bengalese finch songs that were originally reported in our earlier studies Tachibana
282 et al. (2014, 2015). The data were collected from 18 adult males (>140 days after hatching), each isolated in a birdcage
283 placed inside a soundproof chamber. The microphone (Audio-Technica PRO35) was installed above the birdcages. The
284 output of the microphone was amplified using a mixer (Mackie 402-VLZ3) and digitized through an audio interface
285 (Roland UA-1010/UA-55) at 16-bits with a sampling rate of 44.1 kHz. The recordings were then down-sampled to
286 32 kHz (see Tachibana et al. (2014, 2015) for more information about the recording).

287 Song syllables were segmented from the continuous recordings using the thresholding algorithm proposed in the
288 previous studies (Tachibana et al., 2014, 2015). We defined a sequence of the syllables as a bout if every two adjacent
289 syllables in the sequence were spaced at most 500 msec apart. These segmentation processes yielded 465,310 syllables
290 and 9,139 bouts in total (≈ 10.79 hours).

291 4.2 Clustering of Syllables

292 To perform an analysis parallel to the discrete human language data, we classified the segmented syllables into discrete
293 categories in an unsupervised way. Specifically, we used an end-to-end clustering method, named the sequence-
294 to-sequence ABCD-VAE, that combined (i) neural network-based extraction of syllable features and (ii) Bayesian
295 classification, both of which worked in an unsupervised way (i.e., without top-down selection of acoustic features
296 or manual classification of the syllables). This section provides an overview of our method, with a brief, high-level
297 introduction to the two components. Interested readers are referred to S1 in the supporting information, where we
298 provide more detailed information. One of the challenges to clustering syllables is their variable duration as many of
299 the existing clustering methods require their input to be a fixed-dimensional vector. Thus, it is convenient to represent
300 the syllables in such a format (but see Bellman and Kalaba, 1959; Levenshtein, 1966; Morita and O'Donnell, To appear,
301 for alternative approaches). Previous studies on animal vocalization often used acoustic features like syllable duration,
302 mean pitch, spectral entropy/shape (centroid, skewness, etc.), mean spectrum/cepstrum, and/or Mel-frequency cepstral
303 coefficients at some representative points for the fixed-dimensional representation (Katahira et al., 2011; Tachibana
304 et al., 2014; Mielke and Zuberbühler, 2013; Morita and Koda, 2019). In this study, we took a non-parametric approach
305 based on a sequence-to-sequence (seq2seq) autoencoder (Bowman et al., 2016; Chung et al., 2016; Zhao et al., 2017;
306 Sainburg et al., 2019b). The seq2seq autoencoder is a RNN that first reads the whole spectral sequence of an input
307 syllable frame by frame (*encoding*). The spectral sequence was obtained by the short-term Fourier transform with the
308 8 msec Hanning window and 4 msec stride), and then reconstructs the input spectra (*decoding*; see the schematic
309 diagram of the system provided in Figure 2a). Improving the precision of this reconstruction is the training objective of
310 the seq2seq autoencoder. For successful reconstruction, the RNN must store the information about the entire syllable
311 in its internal state—represented by a fixed-dimensional vector—when it transitions from the encoding phase to the
312 decoding phase. And this internal state of the RNN served as the fixed-dimensional representation of the syllables.
313 We implemented the encoder and decoder RNNs by the LSTM (Hochreiter and Schmidhuber, 1997, the encoder was
314 bidirectional; Schuster and Paliwal, 1997).

315 One problem with the auto-encoded features of the syllables is that the encoder does not guarantee their interpretability.
316 The only thing the encoder is required to do is push the information of the entire syllables into fixed-dimensional
317 vectors, and the RNN decoder is so flexible that it can map two neighboring points in the feature space to completely
318 different sounds. A widely adopted solution to this problem is to introduce Gaussian noise to the features, turning the
319 network into the *variational* autoencoder (VAE, depicted in Figure 2c; Kingma and Welling, 2014; Bowman et al.,
320 2016; Zhao et al., 2017, see also Coffey et al., 2019; Goffinet et al., 2019; Sainburg et al., 2019b for its applications to
321 animal vocalization). Abstracting away from the mathematical details, the Gaussian noise prevents the encoder from
322 representing two dissimilar syllables close to each other. Otherwise, the noisy representation of the two syllables will
323 overlap and the decoder cannot reconstruct appropriate sounds for each.

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

- 327 i The VAE is not trained for the sake of clustering, and the entire distribution of the encoded features may not
328 be friendly to existing clustering methods.
- 329 ii The encoded features often include individual differences and do not exhibit inter-individually clusterable
330 distribution (see Figure 3d and the supporting information S1.4).

331 To solve these problems, this study adopted the ABCD-VAE, which encoded data into discrete categories with a
332 categorical noise under the Dirichlet prior, and performed end-to-end clustering of syllables within the VAE (Figure 2b).
333 The ABCD-VAE married discrete autoencoding techniques (Jang et al., 2017; van den Oord et al., 2017; Chorowski
334 et al., 2019) and the Bayesian clustering popular in computational linguistics and cognitive science (e.g., Anderson,
335 1990; Kurihara and Sato, 2004, 2006; Teh et al., 2006; Kemp et al., 2007; Goldwater et al., 2009; Feldman et al., 2013;
336 Kamper et al., 2017; Morita and O'Donnell, To appear). It has the following advantages over the Gaussian VAE +
337 independent clustering (whose indices, except iii, correspond to the problems with the Gaussian VAE listed above):

- 338 i Unlike the Gaussian VAE, the ABCD-VAE is optimized for clustering, aiming at optimal discrete encoding of
339 the syllables.
- 340 ii The ABCD-VAE can exploit a speaker-normalization technique that has proven effective for discrete VAEs:
341 The “Speaker Info.” is fed directly to the decoder (Figure 2b), and thus individual-specific patterns need not be
342 encoded in the discrete features (van den Oord et al., 2017; Chorowski et al., 2019; Tjandra et al., 2019, this is
343 also the framework adopted in the ZeroSpeech 2019, a competition on unsupervised learning of spoken human
344 languages; Dunbar et al., 2019).

Table 1: The size of the training and test data used in the neural language modeling of Bengalese finch songs and the English language. The “SECL” portion of the test syllables was used to estimate the SECL (see §4.4).

Data type	Usage	# of bouts/sentences	# of syllables/words	
			Total	SECL
Bengalese finch	Training	9,039	458,753	—
	Test	100	6,557	4,657
English	Training	12,327	179,456	—
	Test	2,006	21,759	8,833

345 iii Thanks to the Dirichlet prior, the ABCD-VAE can detect the optimal number of categories on its own (under
346 an arbitrarily specified upper bound; Bishop, 2006; O’Donnell, 2015; Little, 2019). This is the major update
347 from the previous discrete VAEs that eat up all the categories available (Jang et al., 2017; van den Oord et al.,
348 2017; Chorowski et al., 2019).

349 Note that the ABCD-VAE can still measure the similarity/distance between two syllables by the cosine similarity of
350 their latent representation immediately before the computation of the classification probability (i.e., logits; cf. Mikolov
351 et al., 2013; Deng et al., 2018).

352 4.3 Language Modeling

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

355 The analysis of context dependency made use of a neural language model based on the current state-of-the-art
356 architecture, Transformer (Vaswani et al., 2017; Al-Rfou et al., 2018; Dai et al., 2019). We trained the language model
357 on 9,039 bouts, containing 458,753 syllables (Table 1). These training data were defined by the complement of the 100
358 test bouts that were selected in the following way so that they were long enough (i) and at least one bout per individual
359 singer was included (ii):

- 360 i The bouts containing 20 or more syllables were selected as the candidates.
- 361 ii For each of the 18 finches, one bout was uniformly randomly sampled among those uttered by that finch.
- 362 iii The other 82 bouts were uniformly randomly sampled from the remaining candidates.

363 The training objective was to estimate the probability of the whole bouts \mathbf{x} conditioned on the information about the
364 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 the
365 singer on its own. Hence, the estimated context dependency did not comprise the correlation among syllables with
366 individuality, which would not count as a major factor especially from a generative point of view.

367 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
368 of the left-to-right processor, predicting each syllable x_t conditioned on the preceding context $\langle \text{sos} \rangle, x_1, \dots, x_{t-1}$,
369 where $\langle \text{sos} \rangle$ stands for the special category marking the start of the bout. See the supporting information S2 for details
370 on the model parameters and training procedure.

371 4.4 Measuring Context Dependencies

372 After training the language model, we estimated how much of the context x_1, \dots, x_{t-1} was used effectively for the
373 model to predict the upcoming syllable x_t in the test data. Specifically, we wanted to know the longest length L of the
374 truncated context x_{t-L}, \dots, x_{t-1} such that the prediction of x_t conditioned on the truncated context was worse (with
375 at least 1% greater perplexity) than the prediction based on the full context (Figure 1b). This context length L is called
376 the *effective context length* (ECL) of the trained language model (Khandelwal et al., 2018).

377 One potential problem with the ECL estimation using the Bengalese finch data was that the test data was much smaller
378 in size than the human language corpora used in the previous study. In other words, the perplexity, from which the ECL
379 was estimated, was more likely to be affected by sampling error. To obtain a more reliable result, we bootstrapped the
380 test data (10,000 samples) and used the five percentile of the bootstrapped differences between the truncated and full
381 context predictions. We call this bootstrapped version of ECL the *statistically effective context length* (SECL).

382 It is more appropriate to estimate the SECL by evaluating the same set of syllables across different lengths of the
383 truncated contexts. Accordingly, only those that were preceded by 20 or more syllables (including <sos>) in the test
384 bouts were used for the analysis (4.657 syllables in total, Table 1).

385 4.5 English Data

386 For comparison, we also estimated the SECL of the language model trained on English data. The data were constructed
387 from the Universal Dependencies English Web Treebank (the training and test portions; Silveira et al., 2014). The
388 database consists of textual English sentences and each word is annotated with the lemma and PoS category. We
389 constructed two versions of training and test data using these lemma and PoS representations of the words: Words
390 may exhibit correlation with one another due to their semantics (e.g., same topic) when they are coded as the lemma.
391 By contrast, the PoS representation of words removes such semantic information, and allowed us to assess the purely
392 syntactic dependencies among the words (cf. Perfors et al., 2011). Note that this semantics-free data may serve as a
393 more appropriate baseline for the study of birdsongs, whose variation is considered not to encode different meanings
394 (Okanoya, 2007; Okanoya and Merker, 2007; Berwick et al., 2011, 2012; Gibson and Tallerman, 2012; Miyagawa et al.,
395 2013, 2014) unlike alarm calls (Seyfarth et al., 1980; Ouattara et al., 2009; Suzuki et al., 2016).

396 The words that were preceded by ten or more tokens (including <sos>) in the test data sentences were used to estimate
397 the SECL. Accordingly, the upper bound on the SECL (=10) was lower than in the analysis of the Bengalese finch data
398 (=20). The reason for the different settings is that the English sentences were shorter than the Bengalese finch bouts:
399 The quartiles of the bout lengths were 22, 44, and 68, while those of the sentence lengths were 7, 14, and 22 (where
400 both the training and test data were included).

401 Acknowledgments

402 This study was funded by MEXT Grant-in-aid for Scientific Research on Innovative Areas #4903 (Evolinguistic;
403 JP17H06380) and JSPS Grant-in-Aid for Scientific Research C (JP19KT0023).

404 References

- 405 Al-Rfou, R., Choe, D., Constant, N., Guo, M., and Jones, L. (2018). Character-level language modeling with deeper
406 self-attention.
- 407 Anderson, J. R. (1990). *The adaptive character of thought*. Studies in cognition. L. Erlbaum Associates, Hillsdale, NJ.
- 408 Beal, M. J., Ghahramani, Z., and Rasmussen, C. E. (2002). The infinite hidden Markov model. In Dietterich, T. G.,
409 Becker, S., and Ghahramani, Z., editors, *Advances in Neural Information Processing Systems 14*, pages 577–584.
410 MIT Press.
- 411 Bellman, R. and Kalaba, R. (1959). On adaptive control processes. *IRE Transactions on Automatic Control*, 4(2):1–9.
- 412 Bengio, Y., Ducharme, R., and Vincent, P. (2001). A neural probabilistic language model. In Leen, T. K., Dietterich,
413 T. G., and Tresp, V., editors, *Advances in Neural Information Processing Systems 13*, pages 932–938. MIT Press.
- 414 Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. (2003). A neural probabilistic language model. *Journal of*
415 *Machine Learning Research*, 3:1137–1155.
- 416 Berwick, R., Beckers, G., Okanoya, K., and Bolhuis, J. (2012). A bird’s eye view of human language evolution.
417 *Frontiers in Evolutionary Neuroscience*, 4:5.
- 418 Berwick, R. C. and Chomsky, N. (2016). *Why Only Us: Language and Evolution*. MIT Press.
- 419 Berwick, R. C., Okanoya, K., Beckers, G. J., and Bolhuis, J. J. (2011). Songs to syntax: the linguistics of birdsong.
420 *Trends in Cognitive Science*, 15(3):113–121.
- 421 Bishop, C. M. (2006). *Pattern recognition and machine learning*. Information science and statistics. Springer, New
422 York.
- 423 Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A., Jozefowicz, R., and Bengio, S. (2016). Generating sentences from a
424 continuous space. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*.
- 425 Chiu, C.-C., Sainath, T. N., Wu, Y., Prabhavalkar, R., Nguyen, P., Chen, Z., Kannan, A., Weiss, R. J., Rao, K.,
426 Gonina, E., Jaitly, N., Li, B., Chorowski, J., and Bacchiani, M. (2017). State-of-the-art speech recognition with
427 sequence-to-sequence models.
- 428 Chomsky, N. (1956). Three models for the description of language. *IRE Transactions on Information Theory*, 2:113 –
429 124.

- 430 Chomsky, N. (1957). *Syntactic Structures*. Mouton and Co., The Hague.
- 431 Chorowski, J., Weiss, R. J., Bengio, S., and van den Oord, A. (2019). Unsupervised speech representation learning using
432 wavernet autoencoders. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(12):2041–2053.
- 433 Chung, Y.-A., Wu, C.-C., Shen, C.-H., yi Lee, H., and Lee, L.-S. (2016). Audio word2vec: Unsupervised learning of
434 audio segment representations using sequence-to-sequence autoencoder. In *INTERSPEECH*, pages 765–769.
- 435 Coffey, K. R., Marx, R. G., and Neumaier, J. F. (2019). DeepSqueak: a deep learning-based system for detection and
436 analysis of ultrasonic vocalizations. *Neuropsychopharmacology*, 44(5):859–868.
- 437 Cybenko, G. (1989). Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals and
438 Systems*, 2(4):303–314.
- 439 Dai, Z., Yang, Z., Yang, Y., Carbonell, J., Le, Q. V., and Salakhutdinov, R. (2019). Transformer-XL: Attentive language
440 models beyond a fixed-length context.
- 441 de Waal, F. B. (1988). The communicative repertoire of captive bonobos (*Pan Paniscus*), compared to that of
442 chimpanzees. *Behaviour*, 106(3-4):183–251.
- 443 Dehaene, S., Changeux, J. P., and Nadal, J. P. (1987). Neural networks that learn temporal sequences by selection.
444 *Proceedings of the National Academy of Sciences*, 84(9):2727–2731.
- 445 Deng, J., Guo, J., Xue, N., and Zafeiriou, S. (2018). Arcface: Additive angular margin loss for deep face recognition.
- 446 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). BERT: Pre-training of deep bidirectional transformers for
447 language understanding. arXiv:1810.04805.
- 448 Dunbar, E., Algayres, R., Karadayi, J., Bernard, M., Benjumea, J., Cao, X.-N., Miskic, L., Dugrain, C., Ondel, L.,
449 Black, A. W., Besacier, L., Sakti, S., and Dupoux, E. (2019). The Zero Resource Speech Challenge 2019: TTS
450 without T. In *Proceedings of Interspeech 2019*, pages 1088–1092.
- 451 Dunbar, E., Cao, X. N., Benjumea, J., Karadayi, J., Bernard, M., Besacier, L., Anguera, X., and Dupoux, E. (2017). The
452 zero resource speech challenge 2017. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop
453 (ASRU)*, pages 323–330.
- 454 Feldman, N. H., Goldwater, S., Griffiths, T. L., and Morgan, J. L. (2013). A role for the developing lexicon in phonetic
455 category acquisition. *Psychological Review*, 120(4):751–778.
- 456 Friston, K. (2003). Learning and inference in the brain. *Neural Networks*, 16(9):1325–1352. Neuroinformatics.
- 457 Friston, K. (2010). The free-energy principle: a unified brain theory? *Nature Reviews Neuroscience*, 11:127–138.
- 458 Friston, K. J. and Stephan, K. E. (2007). Free-energy and the brain. *Synthese*, 159(3):417–458.
- 459 Frith, C. B. and Beehler, B. M. (1998). *The Birds of Paradise: Paradisaeidae*. Bird Families of the World. Oxford
460 University Press, Oxford.
- 461 Futrell, R., Mahowald, K., and Gibson, E. (2015). Large-scale evidence of dependency length minimization in 37
462 languages. *Proceedings of the National Academy of Sciences*, 112(33):10336–10341.
- 463 Gibson, E. (1998). Linguistic complexity: locality of syntactic dependencies. *Cognition*, 68(1):1–76.
- 464 Gibson, K. R. and Tallerman, M. (2012). *The Oxford Handbook of Language Evolution*. Oxford University Press.
- 465 Goffinet, J., Mooney, R., and Pearson, J. (2019). Inferring low-dimensional latent descriptions of animal vocalizations.
466 *bioRxiv*.
- 467 Goldwater, S., Griffiths, T. L., and Johnson, M. (2006). Interpolating between types and tokens by estimating power-law
468 generators. In Weiss, Y., Schölkopf, B., and Platt, J. C., editors, *Advances in Neural Information Processing Systems
469 18*, pages 459–466, Cambridge, MA. MIT Press.
- 470 Goldwater, S., L Griffiths, T., and Johnson, M. (2009). A Bayesian framework for word segmentation: Exploring the
471 effects of context. *Cognition*, 112:21–54.
- 472 Harris, Z. S. (1945). Discontinuous morphemes. *Language*, 21(3):121–127.
- 473 Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- 474 Hornik, K. (1991). Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257.
- 475 Hosino, T. and Okanoya, K. (2000). Lesion of a higher-order song nucleus disrupts phrase level complexity in bengalese
476 finches. *Neuroreport*, 11(10):2091–2095.
- 477 Jaeger, H. and Haas, H. (2004). Harnessing nonlinearity: Predicting chaotic systems and saving energy in wireless
478 communication. *Science*, 304(5667):78–80.

- 479 Jang, E., Gu, S., and Poole, B. (2017). Categorical reparameterization with Gumbel-softmax. In *5th International Con-*
480 *ference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings.*
- 481 Janik, V. M. (1999). Pitfalls in the categorization of behaviour: a comparison of dolphin whistle classification methods.
482 *Animal Behaviour*, 57(1):133–143.
- 483 Jin, L., Gupta, M. M., and Nikiforuk, P. N. (1995). Universal approximation using dynamic recurrent neural networks:
484 discrete-time version. In *Proceedings of ICNN'95 - International Conference on Neural Networks*, volume 1, pages
485 403–408.
- 486 Kakishita, Y., Sasahara, K., Nishino, T., Takahasi, M., and Okanoya, K. (2007). Pattern extraction improves automata-
487 based syntax analysis in songbirds. *Lecture Notes in Artificial Intelligence*, 4828:320–332.
- 488 Kamper, H., Jansen, A., and Goldwater, S. (2017). A segmental framework for fully-unsupervised large-vocabulary
489 speech recognition. *Computer Speech & Language*, 46:154–174.
- 490 Katahira, K., Suzuki, K., Okanoya, K., and Okada, M. (2011). Complex sequencing rules of birdsong can be explained
491 by simple hidden Markov processes. *PLoS ONE*, 6(9):1–9.
- 492 Katz, S. M. (1987). Estimation of probabilities from sparse data for the language model component of a speech
493 recognizer. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 35(3):400–401.
- 494 Kemp, C., Perfors, A., and Tenenbaum, J. (2007). Learning overhypotheses with hierarchical Bayesian models.
495 *Developmental Science*, 10(3):307–321.
- 496 Kershenbaum, A., Blumstein, D. T., Roch, M. A., Akçay, Ç., Backus, G., Bee, M. A., Bohn, K., Cao, Y., Carter, G.,
497 Căsar, C., Coen, M., DeRuiter, S. L., Doyle, L., Edelman, S., Ferrer-i Cancho, R., Freeberg, T. M., Garland, E. C.,
498 Gustison, M., Harley, H. E., Huetz, C., Hughes, M., Hyland Bruno, J., Ilany, A., Jin, D. Z., Johnson, M., Ju, C.,
499 Karnowski, J., Lohr, B., Manser, M. B., McCowan, B., Mercado, E., Narins, P. M., Piel, A., Rice, M., Salmi, R.,
500 Sasahara, K., Sayigh, L., Shiu, Y., Taylor, C., Vallejo, E. E., Waller, S., and Zamora-Gutierrez, V. (2016). Acoustic
501 sequences in non-human animals: a tutorial review and prospectus. *Biological Reviews*, 91(1):13–52.
- 502 Kershenbaum, A., Bowles, A. E., Freeberg, T. M., Jin, D. Z., Lameira, A. R., and Bohn, K. (2014). Animal vocal
503 sequences: not the Markov chains we thought they were. *Proceedings of the Royal Society of London B: Biological*
504 *Sciences*, 281(1792).
- 505 Khandelwal, U., He, H., Qi, P., and Jurafsky, D. (2018). Sharp nearby, fuzzy far away: How neural language models
506 use context. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1:*
507 *Long Papers)*, pages 284–294, Melbourne, Australia. Association for Computational Linguistics.
- 508 Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. The International Conference on Learning
509 Representations (ICLR) 2014.
- 510 Kneser, R. and Ney, H. (1995). Improved backing-off for n-gram language modeling. In *Proceedings of the IEEE*
511 *International Conference on Acoustics, Speech and Signal*, volume 1, pages 181–184.
- 512 Kojima, S. (2003). *A Search for the Origin of Human Speech: Auditory and Vocal Functions of Chimpanzee*. Trans
513 Pacific Press and Kyoto University Press, Rosanna, Melbourne; Kyoto.
- 514 Kurihara, K. and Sato, T. (2004). An application of the variational Bayesian approach to probabilistic context-free
515 grammars. In *International Joint Conference on Natural Language Processing Workshop Beyond Shallow Analyses*.
- 516 Kurihara, K. and Sato, T. (2006). Variational Bayesian grammar induction for natural language. In Sakakibara, Y.,
517 Kobayashi, S., Sato, K., Nishino, T., and Tomita, E., editors, *Grammatical Inference: Algorithms and Applications:*
518 *8th International Colloquium, ICGI 2006, Tokyo, Japan, September 20-22, 2006. Proceedings*, pages 84–96. Springer
519 Berlin Heidelberg, Berlin, Heidelberg.
- 520 Kuypers, H. G. J. M. (1958). Corticobulbar connexions to the pons and lower brain-stem in man: an anatomical study.
521 *Brain*, 81(3):364–388.
- 522 Larson, B. (2017). Long distance dependencies. Oxford Bibliographies.
- 523 Levenshtein, V. I. (1966). Binary codes capable of correcting deletions, insertions and reversals. *Soviet Physics Doklady*,
524 10(8):707–710.
- 525 Liebal, K., Pika, S., and Tomasello, M. (2006). Gestural communication of orangutans (*Pongo pygmaeus*). *Gesture*,
526 6(1):1–38.
- 527 Lin, H. W. and Tegmark, M. (2017). Critical behavior in physics and probabilistic formal languages. *Entropy*, 19(7):299.
- 528 Little, M. A. (2019). *Machine Learning for Signal Processing: Data Science, Algorithms, and Computational Statistics*.
529 Oxford University Press.

- 530 Lu, Z., Pu, H., Wang, F., Hu, Z., and Wang, L. (2017). The expressive power of neural networks: A view from the
531 width. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors,
532 *Advances in Neural Information Processing Systems 30*, pages 6231–6239. Curran Associates, Inc.
- 533 Maass, W., Natschläger, T., and Markram, H. (2002). Real-time computing without stable states: A new framework for
534 neural computation based on perturbations. *Neural Computation*, 14(11):2531–2560.
- 535 Markowitz, J. E., Ivie, E., Kligler, L., and Gardner, T. J. (2013). Long-range order in canary song. *PLOS Computational*
536 *Biology*, 9(5):1–12.
- 537 Mielke, A. and Zuberbühler, K. (2013). A method for automated individual, species and call type recognition in
538 free-ranging animals. *Animal Behaviour*, 86(2):475–482.
- 539 Mikolov, T., Chen, K., Corrado, G. S., and Dean, J. (2013). Efficient estimation of word representations in vector space.
540 arXiv:1301.3781.
- 541 Miyagawa, S., Berwick, R., and Okanoya, K. (2013). The emergence of hierarchical structure in human language.
542 *Frontiers in Psychology*, 4:71.
- 543 Miyagawa, S., Ojima, S., Berwick, R. C., and Okanoya, K. (2014). The integration hypothesis of human language
544 evolution and the nature of contemporary languages. *Frontiers in Psychology*, 5:564.
- 545 Morita, T. and Koda, H. (2019). Superregular grammars do not provide additional explanatory power but allow for a
546 compact analysis of animal song. *Royal Society Open Science*, 6(7):190139. Preprinted in arXiv:1811.02507.
- 547 Morita, T. and Koda, H. (2020). Difficulties in analysing animal song under formal language theory framework:
548 comparison with metric-based model evaluation. *Royal Society Open Science*, 7(2):192069.
- 549 Morita, T. and O’Donnell, T. J. (To appear). Statistical evidence for learnable lexical subclasses in Japanese. *Linguistic*
550 *Inquiry*. Accepted with major revisions.
- 551 Natschläger, T., Markram, H., and Maass, W. (2003). Computer models and analysis tools for neural microcircuits. In
552 Kötter, R., editor, *Neuroscience Databases: A Practical Guide*, pages 123–138. Springer US, Boston, MA.
- 553 Nishikawa, J., Okada, M., and Okanoya, K. (2008). Population coding of song element sequence in the Bengalese finch
554 hvc. *European Journal of Neuroscience*, 27(12):3273–3283.
- 555 Nishikawa, J. and Okanoya, K. (2006). Dynamical neural representation of song syntax in bengalese finch: a model
556 study. *Ornithological Science*, 5(1):95–103.
- 557 Nóbrega, V. A. and Miyagawa, S. (2015). The precedence of syntax in the rapid emergence of human language in
558 evolution as defined by the integration hypothesis. *Frontiers in Psychology*, 6:271.
- 559 O’Donnell, T. J. (2015). *Productivity and reuse in language : a theory of linguistic computation and storage*. MIT
560 Press, Cambridge, MA; London, England.
- 561 Okanoya, K. (2004). Song syntax in Bengalese finches: proximate and ultimate analyses. *Advances in the Study of*
562 *Behavior*, 34:297–345.
- 563 Okanoya, K. (2007). Language evolution and an emergent property. *Current Opinion in Neurobiology*, 17(2):271–276.
564 Cognitive neuroscience.
- 565 Okanoya, K. and Merker, B. (2007). Neural substrates for string-context mutual segmentation: A path to human
566 language. In Lyon, C., Nehaniv, C. L., and Cangelosi, A., editors, *Emergence of Communication and Language*,
567 pages 421–434. Springer London, London.
- 568 Ouattara, K., Lemasson, A., and Zuberbühler, K. (2009). Campbell’s monkeys use affixation to alter call meaning.
569 *PLOS ONE*, 4(11):1–7.
- 570 Payne, R. S. and McVay, S. (1971). Songs of humpback whales. *Science*, 173(3997):585–597.
- 571 Perfors, A., Tenenbaum, J. B., and Regier, T. (2011). The learnability of abstract syntactic principles. *Cognition*,
572 118(3):306–338.
- 573 Prather, J. F., Peters, S., Nowicki, S., and Mooney, R. (2008). Precise auditory–vocal mirroring in neurons for learned
574 vocal communication. *Nature*, 451(7176):305–310.
- 575 Rabin, M. O. and Scott, D. (1959). Finite automata and their decision problems. *IBM Journal of Research and*
576 *Development*, 3(2):114–125.
- 577 Rabiner, L. R. (1989). A tutorial on hidden Markov models and selected applications in speech recognition. *Proceedings*
578 *of the IEEE*, 77:257–286.

- 579 Rao, K., Sak, H., and Prabhavalkar, R. (2017). Exploring architectures, data and units for streaming end-to-end speech
580 recognition with rnn-transducer. In *2017 IEEE Automatic Speech Recognition and Understanding Workshop, ASRU*
581 *2017, Okinawa, Japan, December 16-20, 2017*, pages 193–199.
- 582 Ryeu, J. K., Aihara, K., and Tsuda, I. (2001). Fractal encoding in a chaotic neural network. *Phys. Rev. E*, 64:046202.
- 583 Sainburg, T., Theilman, B., Thielk, M., and Gentner, T. Q. (2019a). Parallels in the sequential organization of birdsong
584 and human speech. *Nature Communications*, 10(3636).
- 585 Sainburg, T., Thielk, M., and Gentner, T. Q. (2019b). Latent space visualization, characterization, and generation of
586 diverse vocal communication signals. *bioRxiv*.
- 587 Scholes, E. I. (2006). Courtship Ethology of Carola’s Parotia (Parotia Carolae). *The Auk*, 123(4):967–990.
- 588 Scholes, E. I. (2008). Evolution of the courtship phenotype in the bird of paradise genus Parotia (Aves: Paradisaeidae):
589 homology, phylogeny, and modularity. *Biological Journal of the Linnean Society*, 94(3):491–504.
- 590 Schuster, M. and Paliwal, K. K. (1997). Bidirectional recurrent neural networks. *IEEE Transactions on Signal*
591 *Processing*, 45(11):2673–2681.
- 592 Seyfarth, R., Cheney, D., and Marler, P. (1980). Monkey responses to three different alarm calls: evidence of predator
593 classification and semantic communication. *Science*, 210(4471):801–803.
- 594 Silveira, N., Dozat, T., de Marneffe, M.-C., Bowman, S., Connor, M., Bauer, J., and Manning, C. (2014). A gold standard
595 dependency corpus for English. In *Proceedings of the Ninth International Conference on Language Resources and*
596 *Evaluation (LREC’14)*, pages 2897–2904, Reykjavik, Iceland. European Language Resources Association (ELRA).
- 597 Suzuki, R., Buck, J. R., and Tyack, P. L. (2006). Information entropy of humpback whale songs. *The Journal of the*
598 *Acoustical Society of America*, 119(3):1849–1866.
- 599 Suzuki, T. N., Wheatcroft, D., and Griesser, M. (2016). Experimental evidence for compositional syntax in bird calls.
600 *Nature Communications*, 7(1):10986.
- 601 Tachibana, R. O., Koumura, T., and Okanoya, K. (2015). Variability in the temporal parameters in the song of the
602 bengalese finch (*Lonchura striata* var. *domestica*). *Journal of Comparative Physiology A*, 201(12):1157–1168.
- 603 Tachibana, R. O., Oosugi, N., and Okanoya, K. (2014). Semi-automatic classification of birdsong elements using a
604 linear support vector machine. *PLOS ONE*, 9(3):1–8.
- 605 Tanner, J. E. and Byrne, R. W. (1996). Representation of action through iconic gesture in a captive lowland gorilla.
606 *Current Anthropology*, 37(1):162–173.
- 607 Teh, Y. W. (2006). A hierarchical Bayesian language model based on Pitman-Yor processes. In *Proceedings of*
608 *the 21st International Conference on Computational Linguistics and the 44th Annual Meeting of the Association*
609 *for Computational Linguistics*, ACL-44, pages 985–992, Stroudsburg, PA, USA. Association for Computational
610 Linguistics.
- 611 Teh, Y. W., Jordan, M. I., Beal, M. J., and Blei, D. M. (2006). Hierarchical Dirichlet processes. *Journal of the American*
612 *Statistical Association*, 101(476):1566–1581.
- 613 Tjandra, A., Sisman, B., Zhang, M., Sakti, S., Li, H., and Nakamura, S. (2019). VQVAE unsupervised unit discovery
614 and multi-scale Code2Spec inverter for Zerospeech Challenge 2019.
- 615 Tsuda, I. (2001). Toward an interpretation of dynamic neural activity in terms of chaotic dynamical systems. *Behavioral*
616 *and Brain Sciences*, 24(5):793–810.
- 617 van den Oord, A., Vinyals, O., and Kavukcuoglu, K. (2017). Neural discrete representation learning. In Guyon, I.,
618 Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *Advances in Neural*
619 *Information Processing Systems 30*, pages 6306–6315. Curran Associates, Inc.
- 620 van der Maaten, L. and Hinton, G. (2008). Visualizing high-dimensional data using t-sne. *Journal of Machine Learning*
621 *Research*, 9:2579–2605.
- 622 van Lawick-Goodall, J. (1968). The behaviour of free-living chimpanzees in the Gombe stream reserve. *Animal*
623 *Behaviour Monographs*, 1:161–311.
- 624 Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017).
625 Attention is all you need. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and
626 Garnett, R., editors, *Advances in Neural Information Processing Systems 30*, pages 5998–6008. Curran Associates,
627 Inc.
- 628 Wang, X., Takaki, S., and Yamagishi, J. (2020). Neural source-filter waveform models for statistical parametric speech
629 synthesis. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 28:402–415.

- 630 Warren, T. L., Charlesworth, J. D., Tumer, E. C., and Brainard, M. S. (2012). Variable sequencing is actively maintained
631 in a well learned motor skill. *Journal of neuroscience*, 32(44):15414–15425.
- 632 Wild, J., Li, D., and Eagleton, C. (1997). Projections of the dorsomedial nucleus of the intercollicular complex (dm) in
633 relation to respiratory-vocal nuclei in the brainstem of pigeon (*columba livia*) and zebra finch (*taeniopygia guttata*).
634 *Journal of Comparative Neurology*, 377(3):392–413.
- 635 Zhao, T., Zhao, R., and Eskenazi, M. (2017). Learning discourse-level diversity for neural dialog models using
636 conditional variational autoencoders. In *Proceedings of the 55th Annual Meeting of the Association for Computational*
637 *Linguistics (Volume 1: Long Papers)*, pages 654–664. Association for Computational Linguistics.