
Contextual and neural representations of sequentially complex animal vocalizations

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Holistically exploring the perceptual and neural representations underlying animal
2 communication has traditionally been very difficult because of the complexity
3 of the underlying signal. We present here a novel set of techniques to project
4 entire communicative repertoires into low dimensional spaces that can be systemat-
5 ically sampled from, exploring the relationship between perceptual representations,
6 neural representations, and the latent representational spaces learned by machine
7 learning algorithms. We showcase this method in one ongoing experiment studying
8 sequential and temporal maintenance of context in songbird neural and perceptual
9 representations of syllables. We further discuss how studying the neural mecha-
10 nisms underlying the maintenance of the long-range information content present in
11 birdsong can inform and be informed by machine sequence modeling.

12 1 Introduction

13 Systems neuroscience has a long history of decomposing the features of complex signals under
14 the assumption that they can be untangled and explored systematically, part-by-part. For example
15 in audition, much of the early work in exploring physiological representations of signals involves
16 playing back sine tones, white noise, or other single-variable-modulated acoustic signals. While
17 these approaches have led to a number of advances in understanding circuits and systems underlying
18 auditory cognition, many neural phenomena cannot be understood without exploring systems in more
19 biologically realistic environments.

20 In many cases, introducing biological realism into controlled experiments requires uncovering the
21 complex feature spaces underlying signals. For example, the neuroscience of human language is based
22 on a rich understanding of the phonological, semantic, and syntactic features of speech and language.
23 In contrast, the communicative spaces of many model organisms in auditory neuroscience are more
24 poorly understood, leading to a very small number of model organisms having the necessary tools for
25 study. In birdsong, for example, biophysical models of song production that have been developed
26 for zebra finches do not capture the dynamics of the dual-syrinx vocal tract of European starlings.
27 More species general approaches to modeling communication would increase the accessibility of
28 more diverse and more systematic explorations of animal communication systems in neuroscience.

29 Here, we propose a method based upon recent advances in generative modeling to explore and
30 exploit the vocal and communicative spaces of animals more generally. We show that unsupervised
31 generative and dimensionality reduction machine learning models can be used to learn a latent
32 representation of various animal communicative systems, requiring few prior assumptions about the
33 animal's communication. We can then sample from these latent spaces to systematically explore
34 neural and perceptual representations of biologically relevant acoustic spaces with complex features.
35 We show that this method is successful in species as diverse as songbirds, primates, insects, cetaceans,
36 and amphibians, and in recording conditions both in the lab and in the field. We demonstrate this

37 technique in one ongoing experimental paradigm exploring the effects of sequential context on neural
38 representations of syllables in songbirds. The method we outline here is currently under development,
39 however, a recent version is available as a resource on GitHub¹.

40 2 Latent representations of animal communication

41 Dimensionality reduction and generative models uncover latent structure and nonlinear features in
42 complex audio and visual data [1]. In recent years, these techniques have played an important role in
43 uncovering motifs in behavioral data [2-4] and circuit computations in neural data [5-7].

44 The methods used to learn latent feature representations are various and have different utility based
45 upon use-case. For example, topologically motivated embedding techniques like UMAP and TSNE
46 are often used for learning data motifs (Fig. 1) but are less useful for generating smooth or novel
47 data manifolds. Generative models like Generative Adversarial Networks (GANs) and Variational
48 Autoencoders (VAEs) are useful for interpolating between and smoothly generating signals, but
49 learn a predefined latent distribution (e.g. uniform or Gaussian) that holds little information about
50 the underlying data. Architectural constraints also result in different representations; for example,
51 convolutional versus recurrent neural networks respect different aspects of data. Because use cases
52 differ when exploring animal communication, we implement and explore a number of these models
53 on audio, such that we produce a series of latent representations that can be used for latent modeling
54 of animal communication in a similar manner as in neural and ethological mapping, as well as stimuli
55 generation for physiological and behavioral probing.

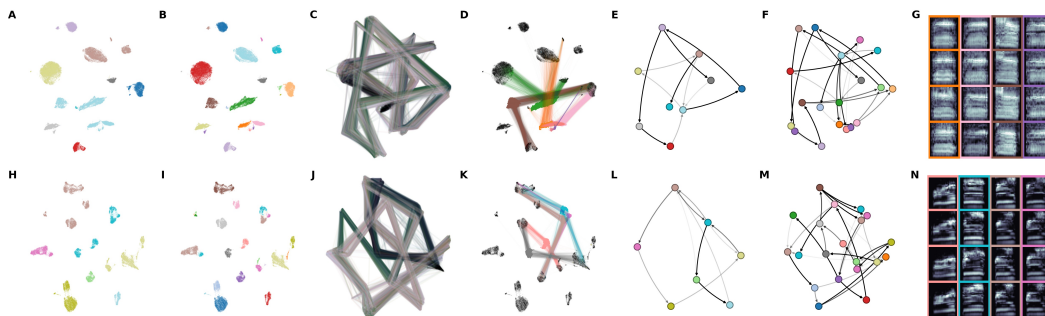


Figure 1: Latent comparisons of hand- and algorithmically-clustered Bengalese finch song. A-H are from a dataset produced by Nicholson et al., [9] and H-N are from a dataset produced by Koumura et al., [10] (A,H) UMAP projections of syllables of Bengalese finch song, colored by a combination of hand labels and supervised labels. (B,I) Algorithmic labels. (C, J) Transitions between syllables, where color represents time. (D,K) Comparing the transitions in one hand-labelled category vs multiple algorithmic labels. (E,L) Markov model from hand labels. (F,M) Markov model from clustered labels. (G,H) Examples of syllables from multiple algorithmic clusters falling under a single hand-labelled cluster.

56 3 Neural and behavioral representations of context

57 A major question in both machine learning and computational cognitive neuroscience is how sequen-
58 tial information can be maintained in order to best predict future events. In machine learning, recent
59 advances such as gated recurrent neural networks, transformer models, and autoregressive models
60 have resulted in progressively improved modelling of sequences, however the flexibility with which
61 sequence models can capture long-range relationships are largely constrained (e.g. [10]) and require
62 vastly different training regimes than the human brain. At the same time, the active maintenance of
63 information in the human brain, through recurrent feedback loops between prefrontal cortex and basal
64 ganglia, show many similarities with sequence models in machine learning [11]. Likewise, songbird
65 basal ganglia and frontal cortex analogous structures actively maintain temporal information, and
66 songbird temporal structure exhibits long-range temporal dependencies that parallel those seen in
67 language [12].

¹Link omitted for anonymity until publication

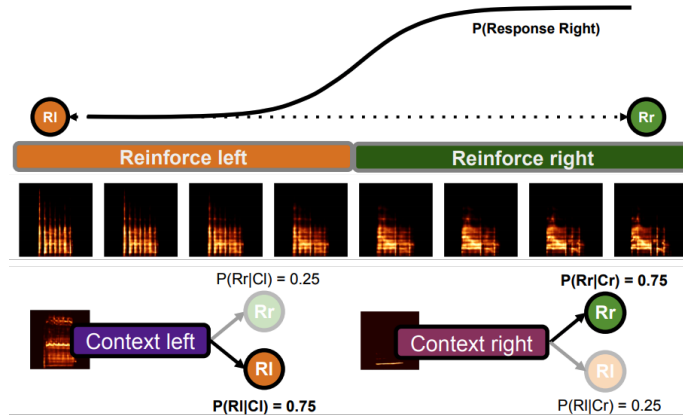


Figure 2: Outline of the context-dependent perception task. Birds are tasked with classifying a smooth morph between syllables generated from a VAE, generating a psychometric function of classification behavior. Sequential-contextual cues that precede the classified syllables are given to bias the psychometric function.

68 In the present experiment we explore how sequential context is maintained in the songbird brain.
 69 To this end, we train a songbird to classify regions of a VAE-generated latent space of song, and
 70 manipulate the perception of those regions of space based upon sequential-contextual information
 71 (Fig 2). Specifically, we interpolate between syllables of European starling song projected into latent
 72 space. We train a starling to classify the left and right halves of the interpolation using an operant-
 73 conditioning apparatus (Fig. 4). We then provide a contextual syllable preceding the classified
 74 syllable that holds predictive information over the classified syllable (Fig 2 bottom). We hypothesize
 75 that the perception of the boundary between the classified stimuli shifts as a function of the context
 76 cue. We model this hypothesis using Bayes rule:

$$\underbrace{P(x_{true} | x_{sensed}, cue)}_{\text{posterior}} \propto \underbrace{P(x_{sensed} | x_{true}, cue)}_{\text{likelihood}} \underbrace{P(x_{true} | cue)}_{\text{prior}}$$

77 When a stimulus varies upon a single dimension x (the interpolation), the perceived value of x is
 78 a function of the true value of x and contextual information (Fig. 3 left). The initial behavioral
 79 results of our experiment confirmed our hypotheses (Fig. 3). We additionally performed acute
 80 extracellular neural recordings from two auditory regions in songbird brain (NCM, CMM) during
 81 stimulus playback under anesthesia. We found single neuron responses to stimuli vary continuously
 82 as a function of interpolation point.

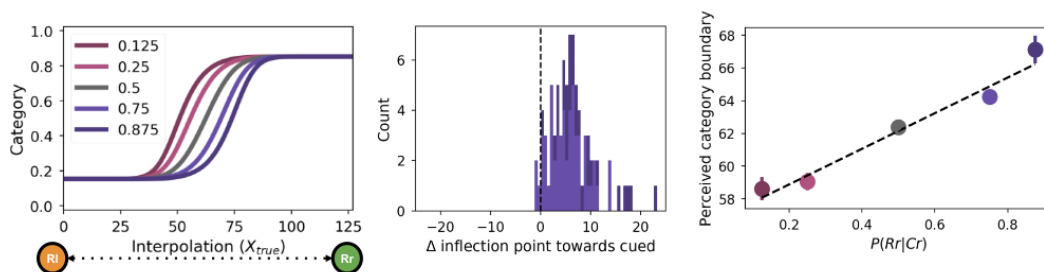


Figure 3: Results for behavior. (left) The Bayesian model predicts that the inflection point of the psychometric classification function will be shifted as a function of the context cue (colored lines). (center) We found that in nearly all behavioral conditions, the inflection point was shifted in the direction of the cue. (right) We found a significant correlation between the cue probability, and the fit categorical boundary using a 4-parameter logistic function ($r(223) = 0.60, p < 0.001$).

83 Because this experiment requires active behavior and maintenance of information, we are currently
 84 implementing a chronic version of this experiment, where neural populations are recorded during

85 active behavior. For this chronic version of the experiment, we designed a custom Raspberry-Pi-based
 86 behavioral apparatus that interfaces with Open Ephys for extracellular acute physiological recordings
 87 (Fig. 4).

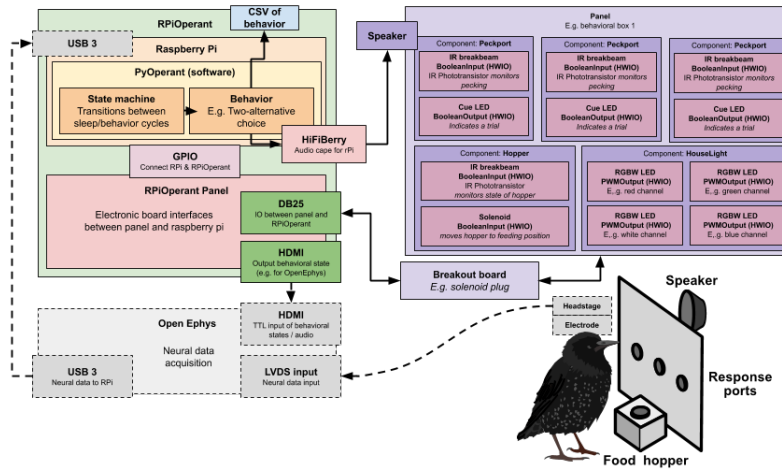


Figure 4: Rig designed for the current study. Behaviors are controlled through a custom interface based upon the Raspberry Pi. Neural data is recorded through an Open Ephys board.

88 Future work will scale the current behavior up to the more complex sequences produced by birds in
 89 a naturalistic setting. Using this paradigm, we hope to directly explore the computational circuits
 90 involved in the complex long-range patterns that emerge in songbird communication. We expect that
 91 understanding how the long-range hierarchical organization of songbird communication is maintained
 92 and represented will be informative in generating new methodologies for maintaining long-range
 93 information content in machine learning and AI. Likewise, we expect that improvements in sequence
 94 modeling in machine learning will provide testable hypotheses for mechanisms underlying neural
 95 circuits of context representation.

96 References

- 97 [1] Sainburg et al., (2018) *Generative adversarial interpolative autoencoding: adversarial training on latent*
 98 *space interpolations encourage convex latent distributions* arXiv:1807.06650
- 99 [2] Brown et al., (2018) *Ethology as a physical science* 10.1038/s41567-018-0093-0
- 100 [3] Berman et al., (2016) *Predictability and hierarchy in Drosophila behavior* 10.1073/pnas.1607601113
- 101 [4] Wiltschko et al., (2015) *Mapping Sub-Second Structure in Mouse Behavior* 10.1016/j.neuron.2015.11.031
- 102 [5] Briggman et al., (2005) *Optical Imaging of Neuronal Populations During Decision-Making* 10.1126/sci-
 103 *ence.1103736*
- 104 [6] Churchland et al., (2012) *Neural population dynamics during reaching* 10.1038/nature11129
- 105 [7] Gao et al., (2017) *A theory of multineuronal dimensionality, dynamics and measurement* 10.1101/214262
- 106 [8] Niven and Kao (2019) *Probing Neural Network Comprehension of Natural Language Arguments*
 107 arXiv:1907.07355
- 108 [9] Nicholson et al., (2017) *Bengalese Finch song repository* 10.6084/m9.figshare.4805749.v5
- 109 [10] Koumura (2016) *BirdsongRecognition* 10.6084/m9.figshare.3470165.v1
- 110 [11] O'Reilly and Frank (2006) *Making working memory work: a computational model of learning in the*
 111 *prefrontal cortex and basal ganglia* 10.1162/089976606775093909
- 112 [12] Sainburg et al., (2019) *Parallels in the sequential organization of birdsong and human speech*
 113 10.1038/s41467-019-11605-y