

# Journal Pre-proof



Title: Bat vocal sequences enhance contextual information independently of syllable order

Y. Amit, Y. Yovel

PII: S2589-0042(23)00543-6

DOI: <https://doi.org/10.1016/j.isci.2023.106466>

Reference: ISCI 106466

To appear in: *ISCIENCE*

Received Date: 1 August 2022

Revised Date: 5 November 2022

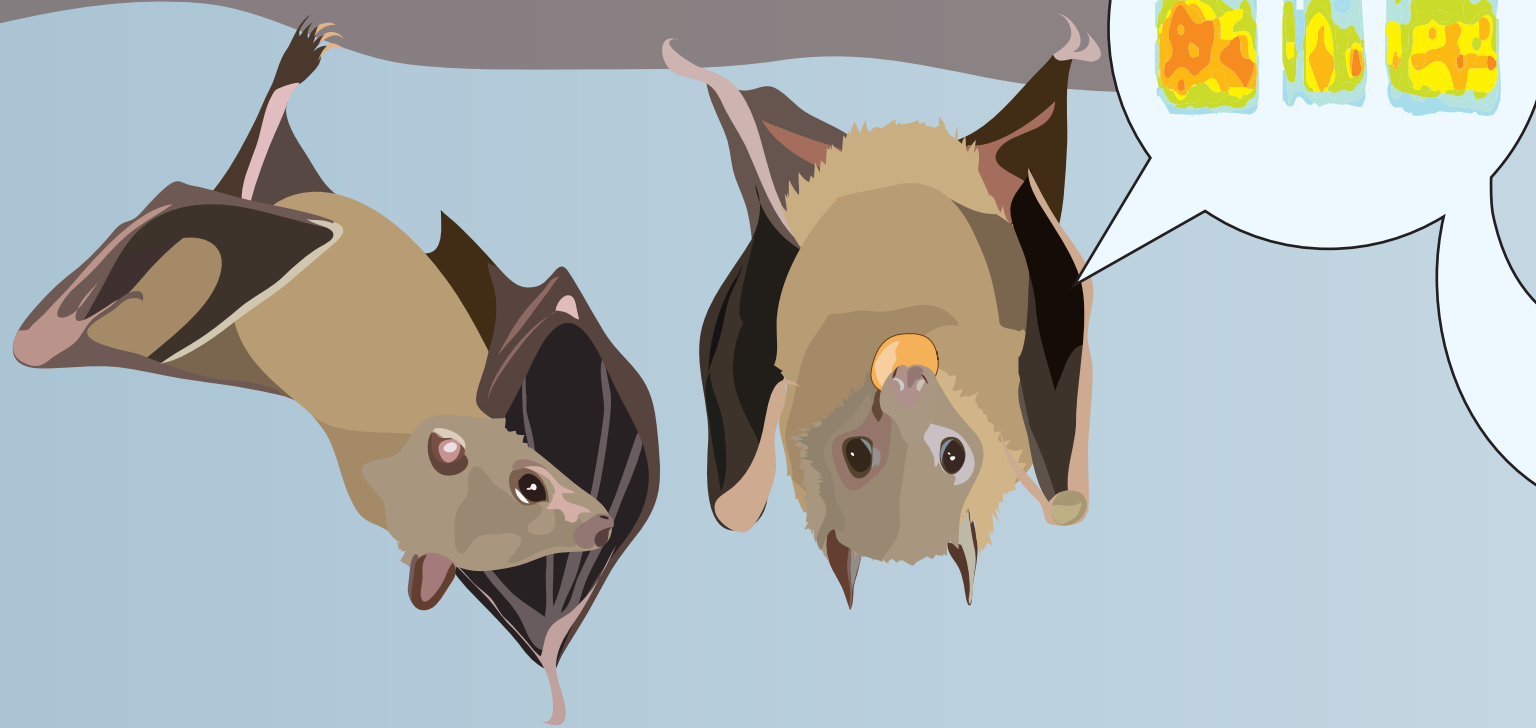
Accepted Date: 17 March 2023

Please cite this article as: Amit, Y., Yovel, Y., Title: Bat vocal sequences enhance contextual information independently of syllable order, *ISCIENCE* (2023), doi: <https://doi.org/10.1016/j.isci.2023.106466>.

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2023

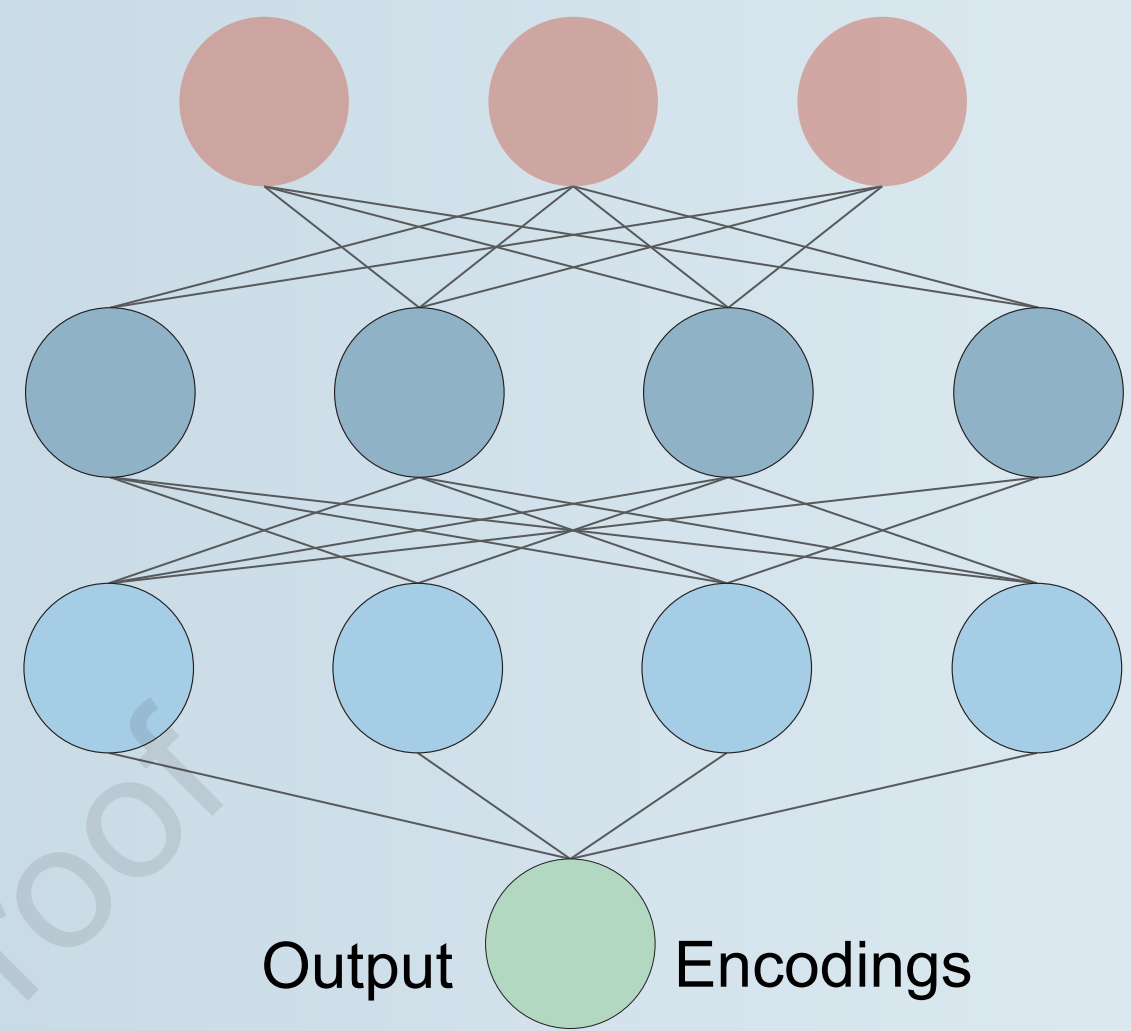
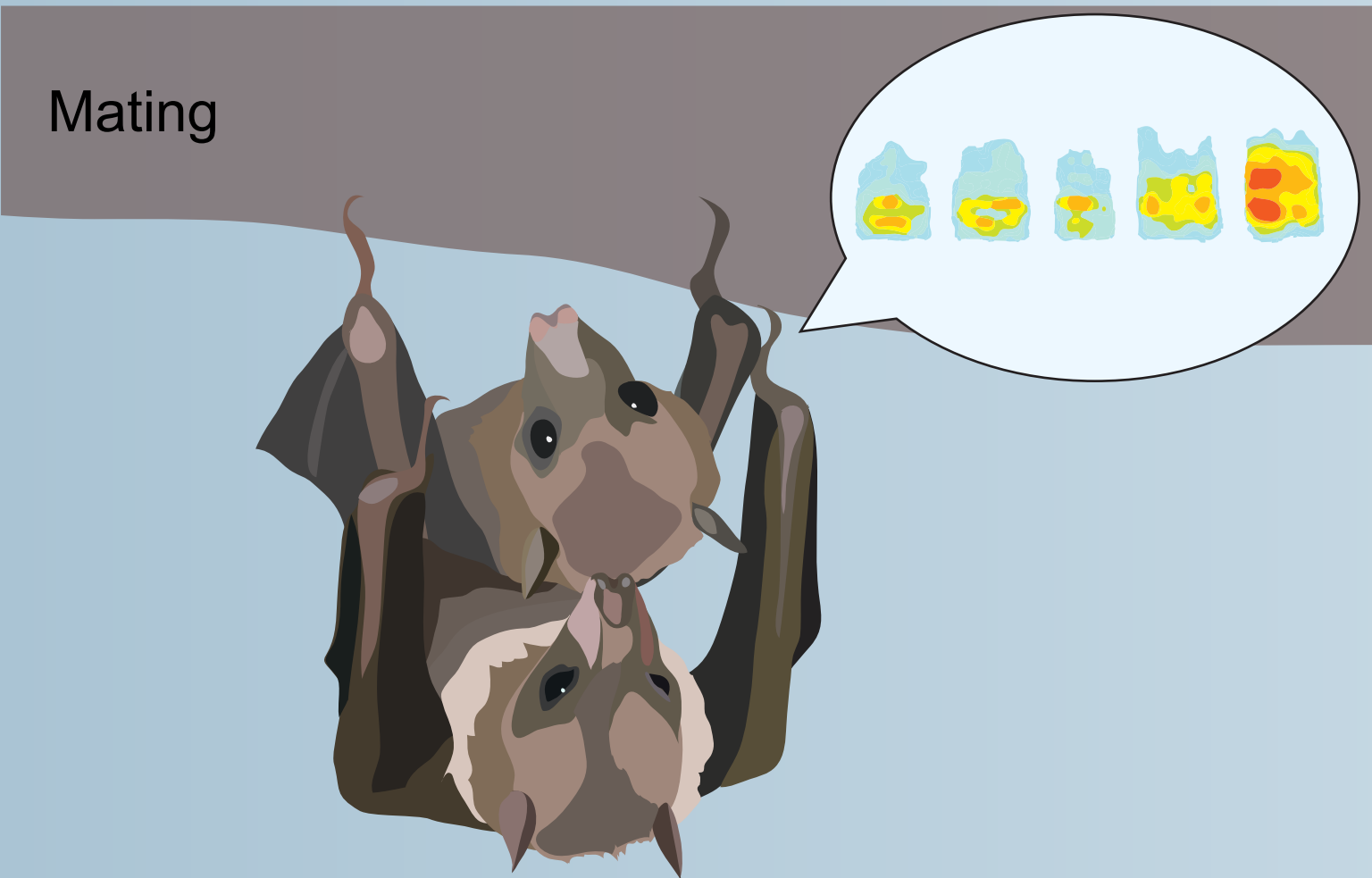
Feeding



Space



Mating

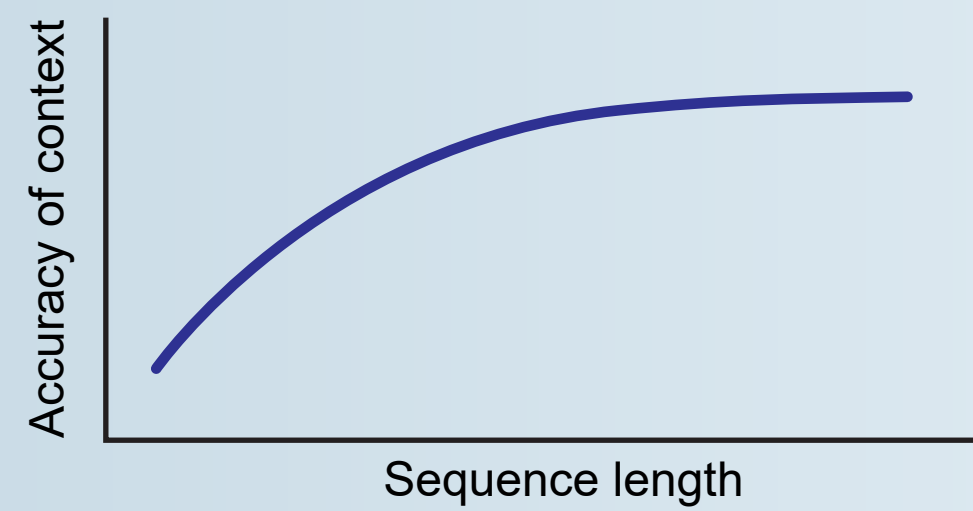


Output Encodings

84663525566852

25786592441711

50113244872541



1 **Title: Bat vocal sequences enhance contextual information independently of syllable order**

2

3 **Authors**

4 Y. Amit<sup>1</sup> & Y. Yovel<sup>1,2,3</sup>

5 **Affiliations**

6 <sup>1</sup>School of Zoology, Faculty of Life sciences, Tel Aviv University, Tel Aviv, Israel.

7 <sup>2</sup>Sagol School of Neuroscience, Tel Aviv University, Tel Aviv, Israel.

8 <sup>3</sup>National Research Center for Biodiversity Studies, The Steinhardt Museum of Natural History,  
9 Tel-Aviv University, Tel Aviv, Israel

10

11 **Author List Footnotes**

12 Lead contact: Yossi Yovel, e-mail: [yossiyovel@gmail.com](mailto:yossiyovel@gmail.com)

13

14

15 **Summary:**

16 Many animals, humans included, rely on acoustic vocalizations for communication. The complexity of non-  
17 human vocal communication has been under debate one of the main open questions being: What could  
18 be the function of multi-syllabic vocal sequences? We address this questions by analyzing fruit-bat vocal  
19 communication. We use neural networks to encode the vocalizations, and statistical models to examine  
20 the information conveyed by sequences of vocalizations. We show that fruit bat vocal sequences  
21 potentially convey more contextual information than individual syllables, but that the order of the  
22 syllables within the sequence is unimportant for context. Specifically, sequences are composed of slightly  
23 modified syllables, thus increasing the probability of context-specificity. We note that future behavioral,  
24 e.g., playback experiments are needed in order to validate the biological relevance of our statistical  
25 results. We hypothesize that such sequences might have served as pre-syntax precursors in the evolution  
26 of animal communication.

27

28 **Introduction:**

29 Animals often emit sequences of social vocalizations. The function of such vocal sequences and how they  
30 evolved from single vocalizations is currently unknown. Many previous studies have suggested that vocal  
31 sequences are not random; that is, they are not composed of a random set of syllables from the animal's  
32 repertoire. The regularities defining non-random sequences are often referred to as the 'syntax' of the  
33 animal communication system<sup>1-3</sup>. In its widest definition, as adopted in this paper, animal communication  
34 syntax refers to any system of rules that orders a sequence of signals in a non-random manner<sup>1-4</sup>. More

35 complex communication systems include syntax that affects the meaning of the vocalizations; that is,  
36 communication systems in which syntax and semantics interact<sup>5</sup>. Syntax is thus commonly graded  
37 according to its complexity. At the highest level is compositional syntax, which has only been shown for a  
38 handful of species<sup>5-8</sup>, which combines meaningful units together into sequences that generate novel  
39 meaning.

40 Sequences and their regularities have been studied in birds<sup>3,4,9-11</sup> and in many mammals including  
41 primates<sup>7,8,12,13</sup>, cetaceans<sup>14</sup>, hyraxes<sup>15</sup>, mongoose<sup>16</sup> and bats<sup>17-19</sup>. Many bats rely on vocalizations for  
42 intra-species social communication (e.g.,<sup>20-22</sup>) often emitting sequences of vocalizations. Several previous  
43 studies suggested that bat vocal sequences are not random. One such study showed that Mexican free-  
44 tailed bats emit sequences with different elements when they are directed at a passing bat vs. when they  
45 are uttered spontaneously<sup>23</sup>. Another study focusing on the neural processing of vocal sequences in the  
46 bat auditory cortex, revealed that neurons respond when the animal is exposed to certain sequences of  
47 vocalizations but not to others<sup>17</sup>. A third study examined the ontogeny of the production of bat vocal  
48 sequences, and found a human-like babbling phase in which sequences or vocalizations are uttered by  
49 newborn pups<sup>24</sup>. However, none of these studies examined the potential information that might be  
50 conveyed by sequences of bat vocalizations, which was the goal of the present study.

51 Focusing on the Egyptian fruit bat, we set out to determine the role of the sequence in bat vocal  
52 communication and to obtain new insight into its evolution. Egyptian fruit bats roost in large colonies,  
53 that can be inhabited by thousands of individuals, which frequently emit sequences of vocalizations as  
54 part of their social interactions. Such sequences are composed of series of up to ~20 vocalizations  
55 (henceforth syllables) with (100-200 ms) intervals of silence between them (Figure 1A-B). Sequences are  
56 separated from each other by much longer (at least one second but often many minute) intervals. The  
57 great majority of vocalizations in this species are uttered during agonistic interactions in the colony, where  
58 each sequence accompanies a single agonistic interaction, and yet, manifest different types of information  
59 (Videos S1-S2-S3 which each demonstrate a single interaction in the contexts: feeding, mating and space  
60 respectively). In a previous study carried by our lab, Prat et al. showed that fruit-bat vocalizations contain  
61 information about the identity of the individual emitter, about the context in which they were uttered,  
62 and to some extent also about the outcome of the interaction<sup>25</sup>. Specifically, it was shown that  
63 vocalizations uttered during agonistic interactions over food, space or mating can be distinguished.  
64 However, in that study, the acoustics of the vocalizations were analyzed in short time windows only, and  
65 thus, the importance of the *sequence* for conveying information and their statistical regularities were  
66 never examined. Because in the previous study, we have already demonstrated that short vocal segments  
67 contain considerable information about the identity of the emitter, here, we focus on the contextual  
68 information conveyed by the sequences.

69 Detecting repeating elements (i.e., categorization of vocalizations) of an animal's communication system  
70 is usually a prerequisite for studying syntax<sup>26-28</sup>. One of the most common methods to achieve this is to  
71 manually scrutinize the recorded vocalizations and to group syllables based on their visual similarities.  
72 This method has been used in numerous studies on song-birds and other species, as well as in most of the  
73 previous bat studies<sup>23,29</sup>. Unlike song-bird vocalizations, fruit-bat (and many other mammalian)  
74 vocalizations are non-tonal and have relatively low fundamental frequencies<sup>25</sup>. They are thus  
75 characterized by numerous noisy harmonics. This makes them especially challenging for categorization,  
76 and thus ill-suited for visual identification of repeatable syllables (see examples in Figure 1A). Here, we  
77 used a combination of deep-learning algorithms and Hidden-Markov-Models (HMMs) in order to embed

78 fruit-bat vocalizations in a lower dimensional feature space, and to examine the order of vocal sequences  
79 and their role in conveying information. We show that while grouping syllables into sequences improves  
80 context classification, the order of the syllables within the sequence, does not affect context classification.  
81 We suggest that such sequences of vocalizations might have appeared early on during the evolution of  
82 animal vocal communication. We note that our analysis is only statistical at this stage, and requires  
83 behavioral experiments for validation.

84

85 **Results:** We adopted a non-supervised deep-learning algorithm to encode the syllables into a lower  
86 dimensional feature space. Specifically, we used a Conditional Variational Autoencoder (CVAE) to encode  
87 the syllables into a 512-dimensions vector. This values of this vector can be thought of as the equivalent  
88 of routinely used acoustic features (e.g., spectral peak). However, when using a neural network (such as  
89 a CVAE), the features usually represent complex spatio-temporal features. Notably, the CVAE was trained  
90 with spectrograms of single syllables while taking the emitter's identity into account (as the condition).  
91 This procedure is common in human speech analysis<sup>30,31</sup> and is crucial for representing inter-individual  
92 variability, which is often the main source of variability in such data-sets. We analyzed recordings of three  
93 female adult fruit-bats recorded continuously for 10 weeks generating a total of 28,847 syllables. This  
94 large data-set allowed us to capture much of the variance in the fruit-bat acoustic system.

95 The feature space produced by the CVAE can be thought of as a multi-dimension description of the  
96 acoustics of the fruit-bat communication system. To scrutinize this feature space, we ran a PCA analysis  
97 on the 512-dimensions and projected the encoded vectors onto the first 40 Principal Components  
98 (accounting for 42% of the variance). We then chose arbitrary vocal-syllables and manipulated them by  
99 moving along each of these 40 PCs in order to examine the effect of each PC-direction on the syllable (in  
100 Figure 1C, we present the effect of the five top PCs to exemplify their action). This analysis revealed that  
101 each PC encompasses multiple spectral and temporal acoustic features and cannot be explained by a  
102 single acoustic parameter. Furthermore, in order to determine acoustic information encoded by our  
103 embedding method, we manipulated random syllables by changing the weight of each PC in steps, and  
104 measured the effect of this manipulation on seven temporal and spectral acoustic features (see Methods).  
105 We found that many of the PCs were correlated with one or more of these seven acoustic features,  
106 demonstrating that the PCs encapsulate acoustic variance (figure 1D-E).

107 The advantage of the CVAE representation in comparison to using specific acoustic features is that it  
108 allows capturing multi-feature acoustic variability. The two most correlated acoustic features were the  
109 temporal roll-off, which is related to the duration of the syllable and the spectral contrast, which is related  
110 to the uniformity of the spectrum (the mean Pearson P-value over all 40 PCs was <0.001 for both of these  
111 acoustic features). Indeed, scrutinizing the effect of the first PC on a randomly-chosen syllable (Figure 1C)  
112 reveals how this PC changes both the duration and the spectral contrast of the syllable (compare the blue  
113 and red lines above and on the side of the spectrograms, representing the duration and spectral  
114 uniformity respectively).

115 In all of the following analyses, we thus used the 40-dimensional vectors (PC-weights) generated by this  
116 method to represent each syllable. Below, we also present all the analyses for a representation of the  
117 vocalizations that is based on a set of specific acoustic features (instead of the CVAE). Next, we sought to  
118 determine whether sequences of vocalizations convey more contextual information than single syllables.  
119 We used annotated sequences of vocalization that were uttered by the bats in one of the three contexts

120 (most commonly observed in our colony): fighting over food – when an individual attempts to scrounge  
121 from another individual; over space – when a bat enters the individual space of another bat; or before  
122 mating, when a female responds aggressively to mating attempts. We will refer to these three contexts  
123 as feeding, space and mating respectively. We trained a Multi-Variate-Gaussian-HMM model with three  
124 hidden states representing the three contexts noted above (note that this HMM was trained using a  
125 supervised approach, see Methods). We trained the HMM model with 326 sequences comprising a total  
126 of 2953 syllables. We divided each sequence into all possible n-grams (yielding a total of 12,900 n-grams).  
127 We then tested the HMM's context classification on sequences with increasing length (between 1-7  
128 syllable n-grams). The HMM model was able to identify the context in which the vocalizations were  
129 uttered far above chance level (Figure 2A, the Balanced Accuracy - BA - for sequences of seven syllables  
130 was  $66 \pm 9\%$  vs. 33% by chance, specifically  $63 \pm 17$ ,  $68 \pm 16$   $69 \pm 19\%$  for the feeding, space and mating  
131 contexts). These results show mean  $\pm$  SD for an 8-fold cross validation procedure in which 87.5% of  
132 sequences are used for training and the rest for testing each time. Notably, context classification improves  
133 when the sequences contain more syllables (overall and at least in two contexts - feeding and space). That  
134 is, the longer the sequence, the more information it conveys about the context ( $P = 1.2 \times 10^{-10}$ , GLM with  
135 the accuracy set as the explained variable, the number of syllables and the context set as fixed factors,  
136 and the cross validation iteration as a random effect, see Supplementary Results 1). The differences  
137 between contexts were also significant, with feeding interactions recognized significantly less than the  
138 other two. We controlled for the effect of dividing the sequences into n-grams by training an HMM  
139 without this division (i.e., on the original sequences only). When doing so using an 8-fold cross validation  
140 we obtained a similar performance,  $61 \pm 10$ ,  $63 \pm 19$   $83 \pm 14\%$  for the feeding, space and mating contexts and  
141 an overall BA of  $66 \pm 10\%$ . We also tested the overall performance for each individual separately (after  
142 training the HMM model on all data together), which revealed a similar average performance for the three  
143 individuals – 55, 70 and 71% (in comparison to a chance level of 33%).

144 We then performed another control, in which we switched syllables between all sequences (across  
145 contexts) keeping their position in the sequence (e.g., we permuted all of the position 2 syllables between  
146 the sequences but always kept them in position 2, without changing any other parts of the training-testing  
147 procedure). In this case, longer sequences did not provide more contextual information validating the  
148 hypothesis that a random assembly of syllables would not convey contextual information (Figure 2B,  
149 average accuracy was at chance level,  $P = 0.63$ , GLM with the same variables as above).

150 We next examined whether the order of the syllables within a sequence contributes to context  
151 classification. To this end, we permuted the internal order of syllables within sequences and we then  
152 trained the same supervised context HMM classifier (as above) with an 8-fold cross-validation. This  
153 internal permutation did not affect the context classification performance of the HMM, suggesting that  
154 syllable-order does not contribute to conveying contextual information. Context classification results in  
155 this case were identical to those of the original data with an accuracy of  $63 \pm 17$ ,  $68 \pm 16$  and  $69 \pm 19$  for the  
156 feeding, space and mating contexts and an overall BA of  $66 \pm 10$  (Figure 2C).

157 To determine whether the model we trained can represent a form of compositional syntax, in which  
158 syllables with certain meanings (i.e., context) are combined into sequences to generate new meanings,  
159 we tested the (above-noted) HMM model on each of the syllables within the sequences separately (i.e.,  
160 on 1-grams) and compared their classified context to the context of the entire sequence. We found that  
161 the classified syllable context was the same as the context of the entire sequence negating compositional  
162 syntax. Specifically, more than 80% of the individual syllables were classified as belonging to the same



163 context as the entire sequence. Thus, we conclude that, from a statistical point of view, individual syllables  
164 convey the same contextual information as the sequence, but because they are not identical acoustically,  
165 the sequence conveys more contextual information than a single syllable alone (see additional discussion  
166 below).

167 To determine whether the ‘simple’ acoustic features that we extracted can also provide contextual  
168 information, we ran the same context-HMM model on these features (instead of the VAE embedding),  
169 either using each feature separately or using all seven features together. This analysis revealed that even  
170 a low dimensional acoustic representation of the syllables already provides contextual information, and  
171 that using all seven features together provides similar contextual information to that when using the VAE  
172 embedding (the overall balanced accuracy was  $64\pm 10\%$  vs.  $66\pm 6\%$  for the seven acoustic vs. the CVAE  
173 features, Figure 2D). Note that space vocalizations did not classify well when using acoustic features  
174 ( $<50\%$ ) suggesting that the CVAE represents the different contexts better on average. Note also, that  
175 sequences conveyed more contextual information than individual syllables also when using an acoustic  
176 feature-based representation ( $P < 6 \cdot 10^{-6}$ , GLM as above, see Supplementary Results 2).

177

178

#### 179 **Discussion:**

180 We found that vocal sequences uttered by fruit-bats convey more contextual information than single  
181 vocalizations. This suggests that the syllables used in each context arise from a different (multi-modal)  
182 acoustic distribution. Notably, there is much overlap between the distributions of the features of syllables  
183 of different contexts (whether we used the CVAE or the simple acoustic features). Indeed, when plotting  
184 any of the features that we tested, they were always part of a continuous distribution rather than  
185 distributed in clusters. Fruit-bat vocalizations thus do not seem to form separate ‘words’ (although it is  
186 also possible that we are not describing them in the relevant feature space of the bat). We thus suggest  
187 that longer sequences convey more contextual information because uttering more vocalizations increases  
188 the chances of producing a distinct context-specific syllable (i.e., from the non-overlapping margins of the  
189 distribution of the two contexts, see schematic in Figure 2E). Note that, when using an HMM-like model  
190 to classify context, concatenating multiple *identical* syllables would not convey more information about  
191 context. Because we found that the order of syllables within a sequence can be randomized without  
192 affecting context classification, we do not refer to fruit-bat sequences as characterized by syntax. While  
193 our results also refute the hypothesis that fruit-bat sequences could be considered a form of  
194 compositional syntax, we do not suggest that bats or even fruit-bats cannot use compositional syntax, as  
195 might be revealed by future studies applying different feature space or different statistics<sup>7</sup>. We thus  
196 describe a system in which animals combine elements (i.e., syllables) that are already informative on their  
197 own to form sequences that convey the same context as the individual syllables, but that combining them  
198 improves the transmission of information (more than repeating them). We note that it is likely that  
199 sequences also provide other information, which we did not test here, such as, regarding the arousal level  
200 or motivation of the emitting animal.

201 In the next paragraph we offer a speculative hypothesis regarding for the evolution of such sequences.  
202 We hypothesize that this form of vocal sequences might be common in animals and might be a precursor  
203 in the evolution of syntax in animal communication (Figure 3). Let us imagine the ancestral fruit-bat colony

204 in which the most common social interaction includes fighting over position in the cluster, and the vocal  
205 repertoire comprises of only a single syllable, which we will refer to as 'Move'. One could imagine that at  
206 higher arousal levels, an excited bat would repeat this syllable several times, uttering a sequence such as:  
207 Move-Move-Move. Such repeated signaling due to urgency is familiar to any pet holder and has also been  
208 documented in non-vocal communication, for instance, in orangutans<sup>32</sup>. In the next phase, the n-  
209 repetition of the syllable might slightly change depending on the context of the interaction. For instance,  
210 when fighting over food the sequence might become Move-Mov-Mov and later perhaps Meve-Mov-Mev.  
211 This could be a result of the arousal level in this specific context (e.g., fighting while mating is more  
212 vigorous than fighting over place) or it could be a result of a physiological constraint, e.g., holding fruit in  
213 the mouth or calling while flying necessitates shortening the syllables. Over time, a sequence structure  
214 similar to the one we describe above might evolve in which a single syllable conveys contextual  
215 information, while a sequence of syllables conveys more information about the same context, because of  
216 the higher chance that one such syllable will be context-distinct. Eventually a communication system will  
217 evolve in which the syllables in the sequence slightly differ from one another and the syllables in  
218 sequences of different contexts derive from different but overlapping distributions. This is somewhat  
219 reminiscent of a process termed 'affixation' shown in primates, in which alarm syllables are modified (e.g.,  
220 elongated) based on motivation and context, leading to a change in their meaning<sup>13</sup>. Notably, several  
221 species of bats including Egyptian fruit bats have been shown to be vocal learners, i.e., they can modify  
222 their vocalizations based on exposure to sounds produced by others. Although vocal learning has mostly  
223 been studied in the context of individual syllables, it could also assist the establishment of certain  
224 sequences as well as the introduction of new variability into sequences.

225 Note that our case differs from what is sometimes referred to as 'graded syntax' where the combination  
226 of syllables signals the degree of agitation in a specific context<sup>6</sup>, because in our case, sequences convey  
227 different contexts (and not a single one). A system such as we describe here might be a precursor for  
228 evolving ordered sequences - or syntax – in which syllables within a sequence are not ordered randomly,  
229 as seems to be the case in fruit-bats. However, much more comparative research is needed in order to  
230 support these ideas.

231 An alternative hypothesis regarding the evolution of sequences with syntax is the lexical constraint  
232 hypothesis<sup>8,33</sup>, suggesting that when a species continuously increases the number of different syllables it  
233 utters, it will reach a point where further additions become uneconomical compared to combining already  
234 existing syllables, either due to production limits or memory limits. We find this hypothesis appealing from  
235 a theoretical point of view, but also suggest that it ignores the fact that animal communication systems  
236 probably evolve from a single or a few syllables<sup>34</sup>, which are thus likely to become first concatenated into  
237 sequences (of identical syllables), and only later modified to convey information. Many simple extant  
238 animal communication systems, such as dog barking, are mostly based on a single syllable that is modified  
239 occasionally based on arousal and other conditions. It is course also possible that different species have  
240 taken different evolutionary routes.

241 Encoding the acoustic properties of fruit-bat vocalizations using a neural-network auto-encoder to  
242 represent the syllables has revealed new insight into the complexity of fruit-bat communication.  
243 Acoustically, we show that both formant-like features and phoneme-like features exist in fruit-bat  
244 vocalizations. This is revealed for instance in PC 3, which seems to both add and remove low frequency  
245 formant-like structures (see red ellipses in Figure 1C) and also to add and remove temporal phoneme-like  
246 features (see orange ellipses in Figure 1C).



247 Both syntax and semantics were traditionally thought to be unique to human language, but have since  
248 been shown to exist to some degree in other animal species<sup>5</sup>. It has been suggested that compositional  
249 syntax evolved when callers and receivers share an interest in exchanging information<sup>6</sup>. We accept this  
250 hypothesis, and suggest how the use of sequences could have evolved even in a social structure in which  
251 individuals typically do not operate as a group<sup>35,36</sup>, but only roost together in aggregations. We have  
252 uncovered a simple form of sequences that conveys contextual information in fruit-bats, despite the lack  
253 of clearly distinguishable syllables and of order within the sequence. Our statistical analysis should be  
254 followed by behavioral experiments in order to validate our findings. This study, however, has touched  
255 upon one of the fundamental questions in animal communication, namely, what is the basic unit of  
256 information while demonstrating a system in which a sequence of multiple units exemplifies the  
257 information that is already conveyed by a single syllable. Such sequences might have served as precursors  
258 for sequences with more developed regularities.

259

### 260 **Limitations of the study:**

261 One major limitation of this study is that the features extracted by the VAE neural network that we used  
262 to encode bat vocalizations might not be the optimal ones. The bat's brain has probably evolved over a  
263 long time period to extract information from social vocalizations. Similar to our VAE, the brain is a non-  
264 linear machine, but the encoding that it uses might be completely different from ours and probably  
265 extracts much more information. A second, and related limitation of this study is the lack of behavioral  
266 evidence to support our statistical findings. Behavioral validation is essential in order to prove that our  
267 findings are relevant for the animals.

268

269

270

271

272 **Figure 1. Acoustic representation of bat vocalizations using neural networks.** (A) Four representative  
273 sequences of fruit bat vocalizations uttered in two contexts. See typical interactions in Videos S1-S2-S3.  
274 (B) The distribution of the number of syllables in fruit-bat vocal sequences. (C) The effect of the first top  
275 five PCs on a random syllable is presented (PC weight increases from left to right). The blue and red lines  
276 above the first row of spectrograms depict the temporal and spectral envelopes (computed by projecting  
277 the spectrogram on the X or Y axes, respectively). These two envelopes are proxies of the temporal roll-  
278 off and the spectral contrast respectively, and it can be seen how moving along PC1 (from left to right)  
279 elongates the syllable and flattens the spectrum, thus reducing spectral contrast. The orange and red  
280 ellipses in the fourth row demonstrate the addition / removal of a temporal phoneme-like feature and a  
281 low frequency formant-like spectral feature, respectively. (D) The correlation of the first five PCs with  
282 seven acoustic features (X-axis, see Methods) revealed that the temporal roll-off and the spectral contrast  
283 were most correlated – see examples in panel (E), where we varied the PC weight and examined the effect  
284 on these two acoustic features.

285

286 **Figure 2. Sequences of information. (A-D)** HMM classification (on the test set only) as a function of the  
 287 number of syllables (X axis) for three contexts (color-coded -see legend). Black line shows the balanced  
 288 accuracy for all three. (A) Original Data. (C) Permuted sequences where syllables are randomly moved  
 289 between sequences but their position within the sequence remains the same. Note that the 1-grams were  
 290 not permuted and thus provide the same information as in 'A'. (C) Permuted sequences where the order  
 291 of the syllables within the sequences was randomly shuffled. Results are identical to in 'A'. (D) Sequences  
 292 represented by seven acoustic features (instead of VAE's). (E) A schematic suggesting why sequences  
 293 contribute to context conveyance. The red and blue shaded areas represent hypothetical distributions of  
 294 several (hypothetical) features for two different behavioral contexts. The numbers represent the order of  
 295 syllables taken from the two sequences shown above the distributions. Despite much overlap between  
 296 the distributions, some syllables within the sequence (e.g., 3 blue and 5 red) will fall near the margins of  
 297 the distribution making classification easier. The schematic depicts one feature, but the feature space is  
 298 actually multi-dimensional.

299

300 **Figure 3. A conceptual framework for the evolution of animal vocal sequences.** We hypothesize that  
 301 single vocalizations ('Move') first evolved into sequences of identical vocalizations, and then modified  
 302 into sequences of slightly different context-specific syllables.

303

#### 304 **Supplementary video titles:**

305 Supplementary videos 1: an example of a feeding interaction including the accompanying vocalizations,  
 306 related to Figure 1A.

307 Supplementary videos 2: an example of a mating interaction including the accompanying vocalizations,  
 308 related to Figure 1A.

309 Supplementary videos 3: an example of an interaction regarding space including the accompanying  
 310 vocalizations, related to Figure 1A.

311

#### 312 **Methods**

313 **Data:** The data include recordings of 3,601 communication sequences (accounting for a total of 28,847  
 314 syllables) recorded from 3 female adult bats in a previous study<sup>37</sup>. All raw annotated recordings (wav files)  
 315 can be found here<sup>38</sup>. The original recording were performed in insulated anechoic chambers in small  
 316 groups of <10 bats in order to assure high quality recordings with little background noise. The pre-  
 317 processing of the recordings included selecting sequences where the emitter and context are clear and  
 318 without loud background noise (see<sup>37</sup>). We used the segmentation into syllables provided in the original  
 319 paper. Each syllable was then transformed into an amplitude spectrogram using the STFT function (with  
 320 a window length of 0.007 sec). Spectrograms were trimmed or zero-padded if necessary to create 256x640  
 321 images (representing 0.5 second segments with a frequency resolution of ~140 Hz). These were used as  
 322 the input for a Conditional Variational Autoencoder neural-network (CVAE, see next paragraph). All  
 323 analyses were performed with Python. Neural network analyses were done using Python Keras<sup>39</sup> and  
 324 HMMs were fit using the Pomegranate and HMMlearn Python packages.

325 **Encoding:** The CVAE neural network was composed of seven convolutional layers (in the encoder) and  
 326 another eight in the decoder (see STAR table for a link to the full code). We only used high Signal-to-  
 327 Noise-Ratio syllables to train the CVAE. To this end, we added a 0.05V threshold relative to the noise in  
 328 order to remove weak syllables. This additional processing removed 57% of the syllables. This procedure  
 329 was only relevant for the training of the CVAE, while (unless stated otherwise) all analyses were performed  
 330 on all syllables. The CVAE beta parameter was gradually increased following the KL-annealing procedure  
 331 from 0.1 to 1 (see <https://arxiv.org/abs/1903.10145>). A CVAE network learns a probabilistic mapping  
 332 between a syllable represented by a (256\*640) spectrogram and a latent 512 feature space vector  
 333 (referred to as the embedding) while accounting for the emitter of each vocalization (the Condition). We  
 334 used 80% of the spectrograms for training and 20% of them for testing the network.

335 **PCA:** We used a PCA analysis in order to reduce the 512 feature space to a 40 dimensional space that  
 336 accounted for 42% of the variance. In order to explain the variance encapsulated by our PC's, we chose  
 337 random real syllables and moved along each of the first five leading-PC directions to illustrate their effect.  
 338 We used the CVAE autoencoder to decode the equivalent 512 embedding-vectors back to spectrograms.  
 339 Specifically, the autoencoder enables converting encoding vectors to syllables and vice versa. Thus, given  
 340 a 40-dimension vector, we can convert it to a 512-dimension encoding using the PCs and then convert it  
 341 into a syllable using the autoencoder.

342 **Comparison with acoustic features:** In order to estimate the effect of these leading PCs on the acoustics  
 343 of the vocalizations, we estimated the correlation between changing the PC and the seven following  
 344 acoustic features (each of them estimated for the entire manipulated syllable). Unlike the vocalization  
 345 systems of some animals (e.g., song-birds, mice and some insectivorous bats), fruit-bat vocalizations are  
 346 what we usually term 'noisy' and thus their fundamental frequency (or pitch) is not easy to estimate. For  
 347 the same reason, it is difficult to talk about frequency modulation.

- 348 1) Spectral contrast<sup>40</sup> – the difference between the mean energy in the top quantile (peak energy)  
 349 of the spectrum to that of the bottom quantile (valley energy).  
 350 2) Temporal centroid<sup>41</sup> defined as Eq. 1:

$$\text{Centroid} = \frac{\sum_{n=0}^{N-1} t(n)x(n)}{\sum_{n=0}^{N-1} x(n)}$$

351 where  $x(n)$  represents the magnitude of bin  $n$ , and  $t(n)$  represents the time of that bin  
 352  
 353

- 354  
 355 3) Spectral centroid<sup>41</sup> defined as Eq. 2:  
 356

$$\text{Centroid} = \frac{\sum_{n=0}^{N-1} f(n)x(n)}{\sum_{n=0}^{N-1} x(n)}$$

357  
 358 where  $x(n)$  represents the magnitude of bin  $n$ , and  $f(n)$  represents the center frequency of that  
 359 bin.  
 360  
 361

- 362 4) The spectral rolloff<sup>42</sup> is defined as the center frequency of a spectrogram bin such that at least  
 363 0.85 of the energy of the spectrum in this frame is contained in this bin and the lower frequencies.  
 364

365 5) The temporal rolloff is defined as the center time of a time bin such that at least 0.85 of the  
 366 temporal energy in this frame is contained in this bin and the in earlier times. This feature is a  
 367 good approximation of the duration of the syllable.

368  
 369 6) The spectral bandwidth<sup>41</sup> is defined as:  
 370

$$\left( \sum_k S(k)(f(k) - f_c)^p \right)^{\frac{1}{p}}$$

371  
 372 where  $S(k)$  is the spectral magnitude at frequency bin  $k$ ,  $f(k)$  is the frequency at bin  $k$ , and  $f_c$  is the  
 373 spectral centroid. We used  $p=2$ , and thus this is equivalent to a weighted standard deviation.  
 374

375  
 376 7) The Spectral flatness<sup>43</sup>, also known as Wiener entropy, which quantifies how tone-like a sound is,  
 377 as opposed to how noise-like.  
 378

379 To determine which acoustic features contribute most to the variance, we computed the Pearson  
 380 correlation of each PC and the above acoustic features; that is, for 100 syllables, we varied the syllables  
 381 by moving along each PC and computed the respective value of the acoustic feature. We then selected  
 382 the features with the lowest Pearson P-values.

383 **Examining context using HMMs:** Using the trained CVAE, we encoded the syllables (without filtering weak  
 384 syllables) into sequences of  $N*512$  (where  $N$  is the number of syllables in the acoustic sequence). Each  
 385 sequence of syllables was then translated into a sequence of PC-weights (where each syllable is encoded  
 386 by 40 PC weights). Here, we only used sequences annotated for three contexts – feeding aggression,  
 387 general fighting and mating aggression, as provided in ref<sup>38</sup> comprising a total of 326 sequences. We  
 388 extracted all ( $n= 1-7$ ) n-grams from the sentences using a sliding window (resulting in a total of 12,900 n-  
 389 grams, but we also controlled for this step by running the entire procedure on the original data only). We  
 390 trained a 3-hidden state multivariate Gaussian HMM, using a supervised approach. That is, we trained the  
 391 HMM such that each hidden state is equivalent to one of the three annotated contexts (*feeding, fighting*  
 392 *or mating*). We evaluated the accuracy of this model on the test set and estimated the performance for  
 393 every n-gram separately. We performed an 8-fold cross-validation procedure, each time randomly  
 394 selecting 87.5% of the data for training.

395 To examine the **compositional syntax** hypothesis we ran the above-noted trained context-HMMs on each  
 396 syllable in the sequences separately. We then examined (using a binomial test) whether the probability  
 397 of a syllable being classified as belonging to a context of the respective sequence was higher than expected  
 398 by chance (0.33). For example, we tested whether the syllables in mating sequences were also classified  
 399 as mating syllables above chance.

400 **Statistics:** To test the effect of the number of syllables in a sequence on context recognition accuracy, we  
 401 used generalized linear models (GLMs) with the accuracy of classification set as the explained variable and  
 402 the number of syllables, the context and their interaction set as fixed factors. We used a logistic link  
 403 function because the explained variable is a proportion. This analysis was also used for the different  
 404 permutation controls.

405 **Acknowledgement:** We thank I. Arnon and Y. prat for reading on commenting on the manuscript. We  
 406 thanks M. Taub for her assistance with the graphics. This project was partially funded by the ERC project  
 407 BehaviorIsland

408 **Author Contribution:** Conceptualization, Y.Y and Y. A; Methodology, Y.A; Software, Y.A.; Writing –  
409 Original Draft, Y.Y. and Y.A.; Writing – Review & Editing, Y.Y. and Y.A.; Project Administration, Y.Y

410

411 **Declaration of interests:**

412 The authors declare no competing interests.

413

414 **Inclusion and diversity:**

415 We support inclusive, diverse and equitable conduct of research.

416

417

418

419

420

421 **Bibliography**

422

- 423 1. Seyfarth, R.M., Cheney, D.L., Bergman, T., Fischer, J., Zuberbühler, K., and Hammerschmidt, K.  
424 (2010). The central importance of information in studies of animal communication. *Animal*  
425 *Behaviour* *80*, 3–8. 10.1016/j.anbehav.2010.04.012.
- 426 2. Suzuki, T.N. (2016). Semantic communication in birds: evidence from field research over the past  
427 two decades. *Ecological Research* *31*, 307–319. 10.1007/S11284-016-1339-X.
- 428 3. Bloomfield, T.C., Gentner, T.Q., and Margoliash, D. (2011). What birds have to say about  
429 language. *Nature Neuroscience* *14*, 947–948. 10.1038/nn.2884.
- 430 4. Gentner, T.Q., Fenn, K.M., Margoliash, D., and Nusbaum, H.C. (2006). Recursive syntactic pattern  
431 learning by songbirds. *Nature* *440*, 1204. 10.1038/NATURE04675.
- 432 5. Suzuki, T.N., Wheatcroft, D., and Griesser, M. (2020). The syntax–semantics interface in animal  
433 vocal communication. *Philosophical Transactions of the Royal Society B* *375*.  
434 10.1098/RSTB.2018.0405.
- 435 6. Griesser, M., Wheatcroft, D., and Suzuki, T.N. (2018). From bird calls to human language:  
436 exploring the evolutionary drivers of compositional syntax. *Current Opinion in Behavioral*  
437 *Sciences* *21*, 6–12. 10.1016/J.COBEHA.2017.11.002.
- 438 7. Zuberbühler, K. (2018). Combinatorial capacities in primates. *Current Opinion in Behavioral*  
439 *Sciences* *21*, 161–169. 10.1016/J.COBEHA.2018.03.015.

- 440 8. Arnold, K., and Zuberbühler, K. (2006). Semantic combinations in primate calls. *Nature* 2006  
441 441:7091 441, 303–303. 10.1038/441303a.
- 442 9. Ten Cate, C., and B Slater, P.J. (1991). Song learning in zebra finches: how are elements from two  
443 tutors integrated? *Anim. Behav* 42, 150–152.
- 444 10. Lachlan, R.F., Verhagen, L., Peters, S., and ten Cate, C. (2010). Are There Species-Universal  
445 Categories in Bird Song Phonology and Syntax? A Comparative Study of Chaffinches (*Fringilla*  
446 *coelebs*), Zebra Finches (*Taenopygia guttata*), and Swamp Sparrows (*Melospiza georgiana*).  
447 *Journal of Comparative Psychology* 124, 92–108. 10.1037/A0016996.
- 448 11. Berwick, R.C., Okanoya, K., Beckers, G.J.L., and Bolhuis, J.J. Songs to syntax: the linguistics of  
449 birdsong. 10.1016/j.tics.2011.01.002.
- 450 12. Clarke, E., Reichard, U.H., and Zuberbühler, K. (2006). The Syntax and Meaning of Wild Gibbon  
451 Songs. *PLOS ONE* 1, e73. 10.1371/JOURNAL.PONE.0000073.
- 452 13. Ouattara, K., Lemasson, A., and Zuberbühler, K. (2009). Campbell's Monkeys Use Affixation to  
453 Alter Call Meaning. *PLOS ONE* 4, e7808. 10.1371/JOURNAL.PONE.0007808.
- 454 14. Mercado, E., Herman, L.M., and Pack, A.A. (2004). Song copying by humpback whales: themes  
455 and variations. *Animal Cognition* 2004 8:2 8, 93–102. 10.1007/S10071-004-0238-7.
- 456 15. Kershenbaum, A., Ilany, A., Blaustein, L., and Geffen, E. (2012). Syntactic structure and  
457 geographical dialects in the songs of male rock hyraxes. *Proceedings of the Royal Society B:*  
458 *Biological Sciences* 279, 2974–2981. 10.1098/RSPB.2012.0322.
- 459 16. Jansen, D.A.W.A.M., Cant, M.A., and Manser, M.B. (2012). Segmental concatenation of individual  
460 signatures and context cues in banded mongoose (*Mungos mungo*) close calls. *BMC biology* 10.  
461 10.1186/1741-7007-10-97.
- 462 17. Esser, K.-H., Condon, C.J., Suga, N., and Kanwal, J.S. (1997). Syntax processing by auditory cortical  
463 neurons in the FM–FM area of the mustached bat *Pteronotus parnellii*. *Proceedings of the*  
464 *National Academy of Sciences* 94, 14019–14024. 10.1073/PNAS.94.25.14019.
- 465 18. Bohn, K., Montiel-Reyes, F., and Salazar, I. (2016). The Complex Songs of Two Molossid Species.  
466 In *Sociality in Bats* (Springer International Publishing), pp. 143–160. 10.1007/978-3-319-38953-  
467 0\_6.
- 468 19. Kanwal, J.S., Matsumura, S., Ohlemiller, K., and Suga, N. (1994). Analysis of acoustic elements and  
469 syntax in communication sounds emitted by mustached bats. *Journal of the Acoustical Society of*  
470 *America* 96, 1229–1254. 10.1121/1.410273.
- 471 20. Wilkinson, G.S. (2003). Social and vocal complexity in bats.
- 472 21. Chaverri, G., Gillam, E.H., and Vonhof, M.J. (2010). Social calls used by a leaf-roosting bat to  
473 signal location. *Biology letters* 6, 441–444. 10.1098/rsbl.2009.0964.
- 474 22. Schoeman, M.C., and Goodman, S.M. (2012). Vocalizations in the Malagasy Cave-Dwelling Fruit  
475 Bat, *Eidolon dupreanum* : Possible Evidence of Incipient Echolocation? *Acta Chiropterologica* 14,  
476 409–416. 10.3161/150811012X661729.
- 477 23. Bohn, K.M., Smarsh, G.C., and Smotherman, M. (2013). Social context evokes rapid changes in  
478 bat song syntax. *Animal Behaviour* 85, 1485–1491. 10.1016/J.ANBEHAV.2013.04.002.

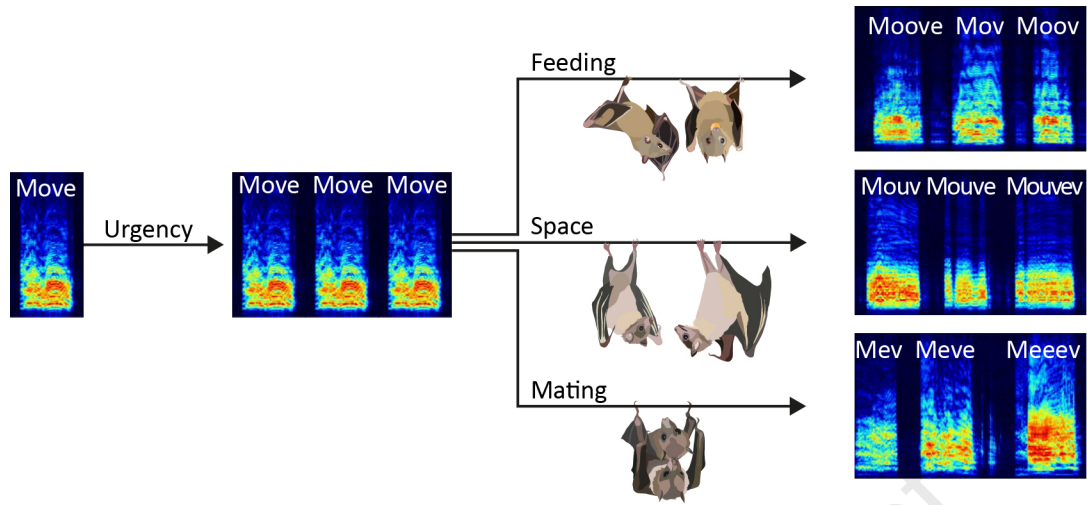


- 479 24. Fernandez, A.A., Burchardt, L.S., Nagy, M., and Knörnschild, M. (2021). Babbling in a vocal  
480 learning bat resembles human infant babbling. *Science* 373, 923–926.  
481 10.1126/SCIENCE.ABF9279.
- 482 25. Prat, Y., Taub, M., and Yovel, Y. (2016). Everyday bat vocalizations contain information about  
483 emitter, addressee, context, and behavior. *Scientific Reports*,.
- 484 26. Prat, Y. (2019). Animals Have No Language, and Humans Are Animals Too:  
485 <https://doi.org/10.1177/1745691619858402> 14, 885–893. 10.1177/1745691619858402.
- 486 27. Kershenbaum, A., Blumstein, D.T., Roch, M.A., Akçay, Ç., Backus, G., Bee, M.A., Bohn, K., Cao, Y.,  
487 Carter, G., Cäsar, C., et al. (2016). Acoustic sequences in non-human animals: a tutorial review  
488 and prospectus. *Biological reviews of the Cambridge Philosophical Society* 91, 13.  
489 10.1111/BRV.12160.
- 490 28. Cate, C. ten, and Okanoya, K. (2012). Revisiting the syntactic abilities of non-human animals:  
491 natural vocalizations and artificial grammar learning. *Philosophical Transactions of the Royal*  
492 *Society B: Biological Sciences* 367, 1984. 10.1098/RSTB.2012.0055.
- 493 29. Knörnschild, M., Eckenweber, M., Fernandez, A.A., and Nagy, M. (2016). Sexually Selected  
494 Vocalizations of Neotropical Bats. In *Sociality in Bats* (Springer International Publishing), pp. 179–  
495 195. 10.1007/978-3-319-38953-0\_8.
- 496 30. Oord, A. van den, Li, Y., and Vinyals, O. (2018). Representation Learning with Contrastive  
497 Predictive Coding.
- 498 31. Jin Park, T., Kanda, N., Dimitriadis, D., Han, K.J., Watanabe, S., and Narayanan, S. A Review of  
499 Speaker Diarization: Recent Advances with Deep Learning.
- 500 32. Orangutans Modify Their Gestural Signaling According to Their Audience's Comprehension -  
501 ScienceDirect <https://www.sciencedirect.com/science/article/pii/S0960982207016405>.
- 502 33. Zuberbühler, K. (2020). Syntax and compositionality in animal communication. *Philosophical*  
503 *Transactions of the Royal Society B* 375. 10.1098/RSTB.2019.0062.
- 504 34. Jorgewich-Cohen, G., Townsend, S.W., Padovese, L.R., Klein, N., Praschag, P., Ferrara, C.R.,  
505 Ettmar, S., Menezes, S., Varani, A.P., Serano, J., et al. (2022). Common evolutionary origin of  
506 acoustic communication in choanate vertebrates. *Nature Communications* 2022 13:1 13, 1–7.  
507 10.1038/s41467-022-33741-8.
- 508 35. Harten, L., Katz, A., Goldshtein, A., Handel, M., and Yovel, Y. (2020). The ontogeny of a  
509 mammalian cognitive map in the real world. *Science* 369, 194–197.
- 510 36. Toledo, S., Shohami, D., Schiffner, I., Lourie, E., Orchan, Y., Bartan, Y., and Nathan, R. (2020).  
511 Cognitive map-based navigation in wild bats revealed by a new high-throughput tracking system.  
512 *Science* 369, 188–193. 10.1126/science.aax6904.
- 513 37. Prat, Y., Taub, M., and Yovel, Y. (2016). Everyday bat vocalizations contain information about  
514 emitter, addressee, context, and behavior. *Scientific Reports* 6, 39419.
- 515 38. Prat, Y., Taub, M., Pratt, E., and Yovel, Y. (2017). Data Descriptor: An annotated dataset of  
516 Egyptian fruit bat vocalizations across varying contexts and during vocal ontogeny. *Scientific Data*  
517 4. 10.1038/sdata.2017.143.

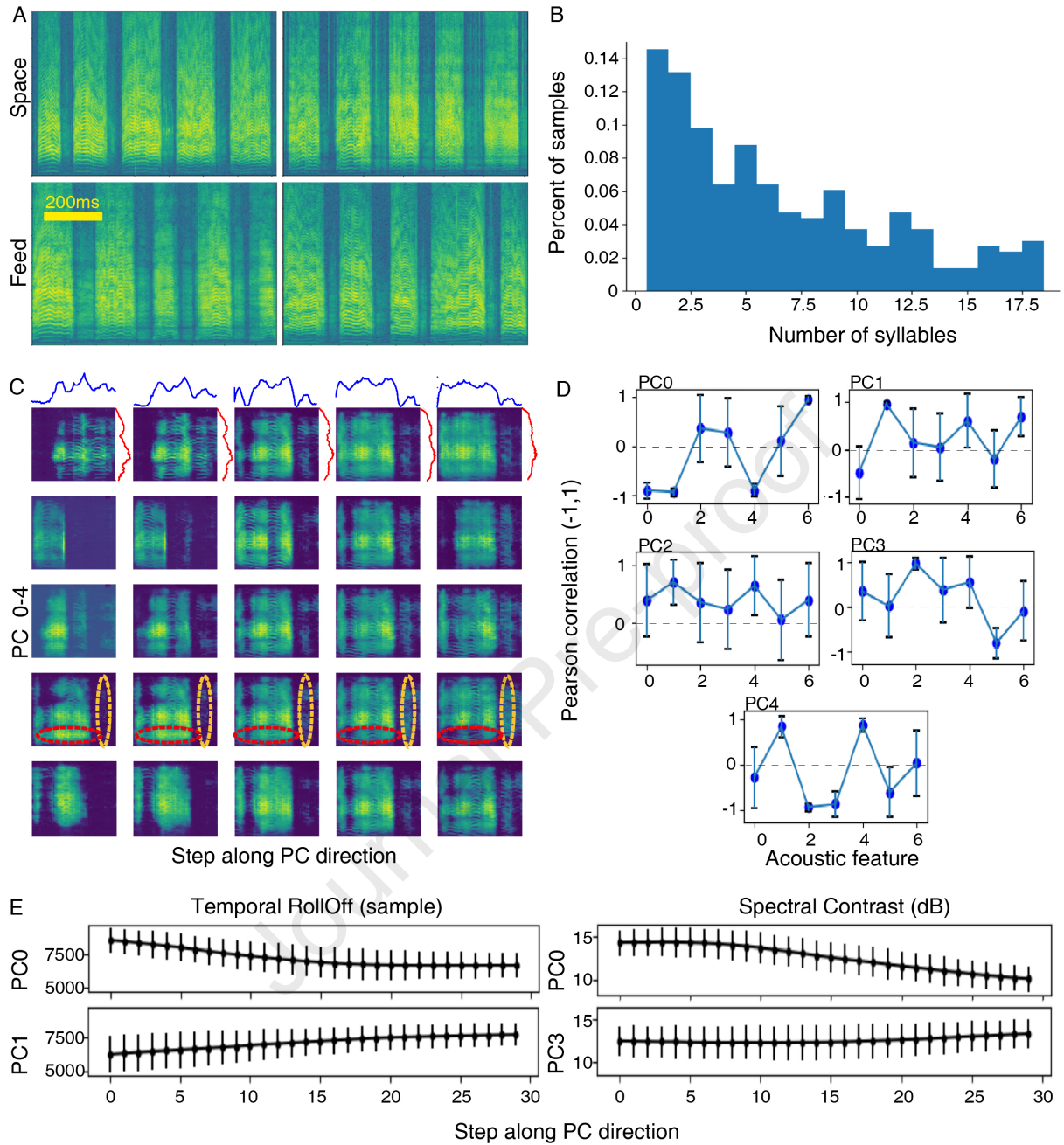
- 518 39. Chollet, F.& others (2015). Keras.
- 519 40. Jiang, D.N., Lu, L., Zhang, H.J., Tao, J.H., and Cai, L.H. (2002). Music type classification by spectral  
520 contrast feature. Proceedings - 2002 IEEE International Conference on Multimedia and Expo,  
521 ICME 2002 1, 113–116. 10.1109/ICME.2002.1035731.
- 522 41. Klapuri, A., and Davy, M. (2006). Signal processing methods for music transcription. Signal  
523 Processing Methods for Music Transcription, 1–440. 10.1007/0-387-32845-9.
- 524 42. (17) A large set of audio features for sound description (similarity and classification) in the  
525 CUIDADO project | Request PDF  
526 [https://www.researchgate.net/publication/200688649\\_A\\_large\\_set\\_of\\_audio\\_features\\_for\\_sou  
527 nd\\_description\\_similarity\\_and\\_classification\\_in\\_the\\_CUIDADO\\_project](https://www.researchgate.net/publication/200688649_A_large_set_of_audio_features_for_sound_description_similarity_and_classification_in_the_CUIDADO_project).
- 528 43. Dubnov, S. (2004). Generalization of spectral flatness measure for non-Gaussian linear processes.  
529 IEEE Signal Processing Letters 11, 698–701. 10.1109/LSP.2004.831663.
- 530
- 531

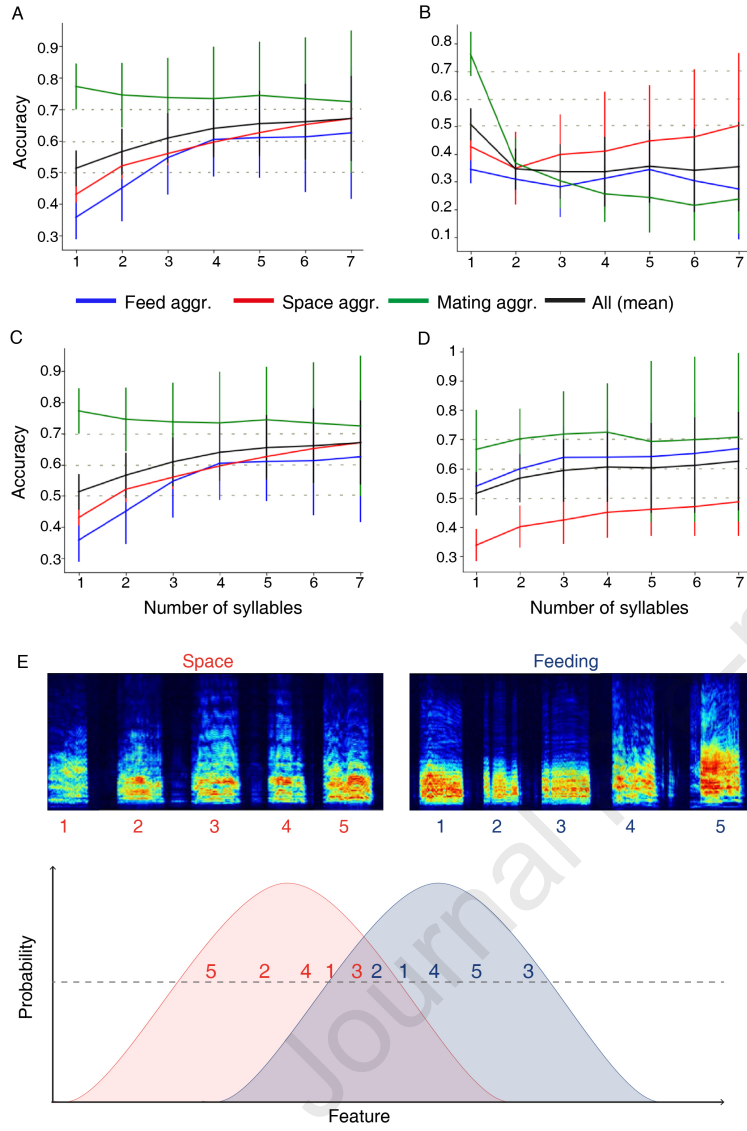
## KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Original wav files	Previous study	Prat et al. 2017, Scientific data
Chemicals, Peptides, and Recombinant Proteins		
Not relevant		
Deposited Data		
Acoustic syllable encodings	Self-recordings	<a href="https://data.mendeley.com/datasets/mjfv43zgtv/3">https://data.mendeley.com/datasets/mjfv43zgtv/3</a>
Experimental Models:		
Three female Egyptian fruit bats ( <i>Rousettus aegyptiacus</i> )	Caught in a cave in central Israel	Taxonomy ID: 9407
Acoustics Recordings		
Microphones + AD converters	Avisoft Bio-acoustics	CM16, SM1612
Quantification and statistical Analysis		
Stats (GLMs) were run in Matlab 2019	The Mathworks	<a href="https://www.mathworks.com/downloads/">https://www.mathworks.com/downloads/</a> ;
All samples were randomized to control for possible biases. Exclusion was based on signal quality. The exact criteria are explained in the methods		
Software and Algorithms		
Self-written code	Self-written code in Python	<a href="https://data.mendeley.com/datasets/mjfv43zgtv/3">https://data.mendeley.com/datasets/mjfv43zgtv/3</a>



Journal Pre-proof







Fruit bats emit sequences of vocalizations while interacting with conspecifics

Artificial neural networks can be used to encode bat vocalizations

Longer sequences of vocalizations convey more information about their context

The order of the syllables in the sequence does not seem to affect information

Journal Pre-proof