1  **Separating overlapping bat calls with a bi-directional long short-term**

2  **memory network**

3  **Using deep neural network to separate overlapping bat calls**

4  Kangkang Zhang[1], Tong Liu[1], Shengjing Song[1], Xin Zhao[1], Shijun Sun[2], Walter

5  Metzner[3], Jiang Feng[1,4*], Ying Liu[1*]

6  [1]Jilin Provincial Key Laboratory of Animal Resource Conservation and Utilization,

7  Northeast Normal University, Changchun, China. [2]School of Environment, Northeast

8  Normal University, Changchun, China. [3]Department of Integrative Biology and

9  Physiology, University of California, Los Angeles, California, USA. [4]Collage of

10  Animal Science and Technology, Jilin Agricultural University, Changchun, China.

11

12  [*]Correspondence

13  E-mail: liuy252@nenu.edu.cn (YL)

14  E-mail: fengj@nenu.edu.cn (JF)

15

16  **Abstract**

17      Acquiring clear and usable audio recordings is critical for acoustic analysis of

18  animal vocalizations. Bioacoustics studies commonly face the problem of overlapping

19  signals, but the issue is often ignored, as there is currently no satisfactory solution.

20  This study presents a bi-directional long short-term memory (BLSTM) network to

21  separate overlapping bat calls and reconstruct waveform audio sounds. The separation

22    quality was evaluated using seven temporal-spectrum parameters. The applicability of

23    this method for bat calls was assessed using six different species. In addition,

24    clustering analysis was conducted with separated echolocation calls from each

25    population. Results showed that all syllables in the overlapping calls were separated

26    with high robustness across species. A comparison between the seven

27    temporal-spectrum parameters showed no significant difference and negligible

28    deviation between the extracted and original calls, indicating high separation quality.

29    Clustering analysis of the separated echolocation calls also produced an accuracy of

30    93.8%, suggesting the reconstructed waveform sounds could be reliably used. These

31    results suggest the proposed technique is a convenient and automated approach for

32    separating overlapping calls using a BLSTM network. This powerful deep neural

33    network approach has the potential to solve complex problems in bioacoustics.

34    **Author summary**

35    In recent years, the development of recording techniques and devices in animal

36    acoustic experiment and population monitoring has led to a sharp increase in the

37    volume of sound data. However, the collected sound would be overlapped because of

38    the existence of multiple individuals, which laid restrictions on taking full advantage

39    of experiment data. Besides, more convenient and automatic methods are needed to

40    cope with the large datasets in animal acoustics. The echolocation calls and

41    communication calls of bats are variable and often overlapped with each other both in

42    the recordings from field and laboratory, which provides an excellent template for

43    research on animal sound separation. Here, we firstly solved the problem of

2

44    overlapping calls in bats successfully based on deep neural network. We built a

45    network to separate the overlapping calls of six bat species. All the syllables in

46    overlapping calls were separated and we found no significant difference between the

47    separated syllables with non-overlapping syllables. We also demonstrated an instance

48    of applying our method on species classification. Our study provides a useful and

49    efficient model for sound data processing in acoustic research and the proposed

50    method has the potential to be generalized to other animal species.

51    **Introduction**

52    The structural identification of vocal units is essential in animal acoustic studies for

53    sound feature analysis, sound emitter recognition, and species identification and

54    monitoring. However, wild animal monitoring, both in the field and in the laboratory,

55    often involves problems caused by the overlapping of different vocal units in time and

56    frequency space, which prevents the components from being suitable for parameter

57    analysis. As a result, the separation of overlapping sounds is an important task in

58    bioacoustic signal processing. However, existing analysis software often struggles to

59    process overlapping calls and previous research on the acoustic identification of

60    animals primarily focuses on extracting target signals from background noise for

61    species classification or population monitoring [1-4]. The process of separating

62    overlapping calls from mixed sounds has received little attention to date and

63    researchers conventionally abandon sounds that overlap in both time and frequency,

64    requiring an extension of the experimental period to obtain sufficient non-overlapping

65    recordings [5, 6]. As such, an effective method for successfully and automatically

66    separating overlapping calls would be of significant interest and benefit to animal

67    researchers.

68    Previous studies using deep neural networks have produced promising results for

69    automated sound recognition in complex acoustic environments for animal species

70    recognition and classification [6-8]. However, in this study, we consider the more

71    difficult task of separating different types of syllables from overlapping calls and

72    reconstructing sound waves from these separated signals. Existing techniques used for

73    animal sound separation often require prohibitive quantities of labelled data. For

74    example, multiple-instance machine learning (MIML) algorithms were proposed for

75    use in sound feature extraction and species identification in birds [1]. However, this

76    technique requires a cropped mask of a signal segment (without overlap) in order to

77    extract each syllable.

78    Deep learning networks have been applied to bioacoustic studies but have primarily

79    been used for classification. For instance, convolutional bidirectional recurrent neural

80    networks (CBRNNs) have been used to identify the presence of bird calls in audio

81    samples [4]. Acoustic features were learned by the network (a classifier) and the

82    presence or absence of a bird call was output as an indicator. Convolutional neural

83    networks (CNNs) have been used to predict the presence of a search-phase bat

84    echolocation call in spectrograms. This binary classification problem was used to

85    detect the presence of bats [2]. To our knowledge, the use of deep learning techniques

86    to separate animal calls that overlap in both time and frequency space has yet to be

87    reported.

88    Multiple studies have been conducted using deep learning-based supervised speech

89    separation with humans. Early systems included shallow models that performed a

90    linear transformation of given mixture features during the prediction time interval.

91    This has included Gaussian mixture models [9], support vector machines [10], and

92    non-negative matrix factorization [11]. However, in real-world scenarios, the mapping

93    relationship between mixture signals and sources is typically a nonlinear

94    transformation. Nonlinear models, such as deep neural networks (DNNs), are

95    therefore highly applicable because of their ability to identify nonlinear structures in

96    audio signals [12-14]. Additionally, recurrent neural networks (RNNs) that exhibit the

97    temporal behavior of a time sequence can be trained to predict time-frequency masks

98    for target signals and separate sources from a mixed waveform [15]. Specifically, long

99    short-term memory (LSTM) networks, a variation of RNN models that exhibit strong

100   learning capabilities and simple construction, have been widely used for word and

101   continuous speech recognition [16-18]. By concatenating two separate LSTM

102   networks, bidirectional LSTMs (BLSTMs) can predict each element of a sequence

103   based on past and future context and can naturally account for the temporal dynamics

104   of speech. These models are typically faster and more accurate than standard RNNs in

105   frame-by-frame phoneme classification [19]. In addition, the BLSTM network can

106   compensate for exploding and vanishing gradient issues that can occur during the

107   training of standard RNN models [20]. At present, BLSTMs have achieved

108   state-of-the-art performance for speech recognition [14, 21], natural language

109   processing [22, 23], and speaker-independent speech separation [24]. As such, a

110    BLSTM model was selected in this study for overlapping bat call separation.

111    Echolocating bats have two vocal repertoires, stereotypical echolocation calls for

112    orientation and a variety of communication calls for social activities [25-27].

113    Recordings from both field and laboratory studies indicate that utterances from

114    individual bats often overlap in both time and frequency, which provides an excellent

115    template for research on overlapping sound separation in animals. The primary

116    objective of this study is to develop a technique for separating two target signals

117    (echolocation and socialization calls) from mixtures of acoustic sounds. Although

118    deep leaning has been employed in the acoustic classification of multiple species,

119    including nonhuman primates [28], birds [4], whales [5], and bats [2, 3], the goal of

120    the present study is distinct from these previous cases in which deep neural networks

121    were primarily used as classifiers.

122    Both overlapping and non-overlapping calls (of both echolocation and

123    communication types) were recorded from each of the collected bat species studied in

124    our previous work. We developed a BLSTM network and used the recorded

125    non-overlapping calls to train the model. Recorded overlapping calls were input to the

126    trained model and separated. Independent sound files were then reconstructed for each

127    separated signal. The correctness of these separated signals was measured by

128    comparing the temporal-spectrum parameters between separated calls and the initially

129    recorded (non-overlapping) calls from each species. Finally, clustering analysis was

130    conducted to classify the bats using separated echolocation calls, which provided a

131    practical application of the proposed technique.

## Results

132

133    The proposed algorithm performed well and achieved high accuracy in separating

134    overlapping calls for each of the six species. The BLSTM model was iteratively

135    trained until the training and validation losses reached a minimum. Loss is a

136    summation of errors made with each sample in the training or validation sets and

137    measures how well the model adapts during optimization. Training loss for this model

138    decreased significantly in the first epoch. The validation loss function tended toward

139    an asymptotic value, indicating the training algorithm had converged (S2 Fig). The

140    BLSTM model converged slightly faster when training with CF bat samples (as

141    opposed to FM samples).

142    All echolocation and communication calls in the overlapping signals were

143    correctly extracted during the separation procedure, regardless of their pulse duration

144    or energy characteristics (see Table 1 and Fig 1). In addition, low-intensity FM

145    components in echolocation pulses were successfully extracted from three CF bat

146    species (Figs 1d, 1e, and 1f).

147    **Table 1. Separation results.**

| Species | Call type | Number of syllable types | Number of syllables in overlapping calls | Number of overlapping syllables | Number of separated syllables |
|---|---|---|---|---|---|
| *Rhinolophus ferrumequinum* | Echolocation | 1 | 14 | 14 | 14 |
| | Communication | 4 | 8 | 8 | 8 |
| *Vespertilio sinensis* | Echolocation | 1 | 21 | 13 | 13 |
| | Communication | 4 | 8 | 8 | 8 |
| *Hipposideros armiger* | Echolocation | 1 | 28 | 19 | 19 |
| | Communication | 6 | 13 | 13 | 13 |
| *Myotis* | Echolocation | 1 | 54 | 36 | 36 |

| *macrodactylus* | Communication | 6 | 15 | 15 | 15 |
|---|---|---|---|---|---|
| *Rhinolophus* | Echolocation | 1 | 42 | 30 | 30 |
| *pusillus* | Communication | 6 | 10 | 10 | 10 |
| *Ia io* | Echolocation | 1 | 26 | 16 | 16 |
| | Communication | 4 | 11 | 11 | 11 |

148

149  **Fig 1. Spectrograms from original recordings of overlapping calls and calls separated by**

150  **the BLSTM network**. The first graph represents each line of the original overlapping calls

151  and the second and third graphs show the separated echolocation and communication calls,

152  respectively.

153

154  A comparison of seven temporal-spectrum parameters from the separated calls

155  and the original recorded non-overlapping calls showed no significant differences (Fig

156  2 and S3 Table). In addition, parameter deviations in separated calls and original

157  non-overlapping calls showed minimal RMSE values for both echolocation and

158  communication signals (Fig 3 and Fig 4). Clustering analysis performed with

159  separated echolocation calls produced an accuracy of 93.8% across species (Fig 5).

160

161  **Fig 2. Comparisons between the separated and original calls.** Two principle

162  components extracted from seven temporal-spectral parameters were used in the study.

163  Results for echolocation and communication calls are shown in (A-F) and (G-L),

164  respectively.

165  **Fig 3. A comparison of deviations for separated and original echolocation calls.**

166  The RMSE value is shown under each plot. The vertical axis represents values for

167    each parameter and the horizontal axis represents the number of syllables measured.

168    The red triangles represent separated calls and the blue dots represent original calls.

169    Abbreviations include duration (duration), Fstart (starting frequency), Fend (ending

170    frequency), Fpeak (peak frequency), Fmin (minimum frequency), Fmax (maximum

171    frequency), and bandw (bandwidth).

172    **Fig 4. A comparison of deviations in separated and original communication calls.**

173    The RMSE value is shown under each plot. The vertical axes and abbreviations are

174    the same in Fig 3.

175    **Fig 5. Clustering analysis for six bat species based on their separated**

176    **echolocation calls.** Overlapping echolocation signals cannot be used for species

177    identification until after separation.

178

179    **Discussion**

180        The BLSTM network used in the present study achieved high accuracy in

181    separating overlapping echolocation and communication calls from bats. The training

182    and validation loss for the model also exhibited fast convergence and high robustness

183    for bat vocalizations. In particular, the separated calls extracted by the proposed

184    algorithm were reconstructed as waveform files with nearly the same quality as the

185    non-overlapping calls, suggesting BLSTM networks to be useful tools for separating

186    signals in future bioacoustic research, such as sound analysis, acoustic identification,

187    species classification, and wild animal monitoring.

188        It was difficult to compare the performance of this algorithm with that of previous

189    studies, primarily because of differences in the experimental procedure. However, a

190    comparison of temporal-spectrum parameters between separated calls and

191    non-overlapping calls was included as an evaluation metric. The seven parameters

192    used in this study are commonly used in bat studies to describe the temporal-spectral

193    features of syllables [26, 29]. Statistical results for this comparison showed no

194    significant differences and small deviations in parameters between separated calls and

195    original recordings, indicating the system was able to separate calls without affecting

196    syllable quality. In addition, clustering analysis conducted with reconstructed

197    echolocation calls was highly accurate (93.8%) for species classification, indicating

198    that calls separated from overlapping signals could be used to synthesize initial data.

199        The BLSTM network exhibited good performance across all six bat species using

200    both narrow and broad time-frequency calls. It also successfully separated different

201    syllable types from both overlapping echolocation and communication calls (Table 1,

202    Fig 1). No species-specific *a priori* knowledge or particular acoustic sensor was

203    directly encoded into the system, making it generalizable to other animal populations

204    with additional training data. Although the dichotomy between communication and

205    echolocation calls is relatively drastic, the proposed separation system has potential

206    applications for other species, as such mixtures are very common in bats. In the future,

207    more complex emitter-independent separation could be conducted using the proposed

208    system, such as combinations of echolocation or social calls from other animals.

209    While deep learning models generally perform better when provided with more data,

210    training with bat calls requires fewer samples than human speech separation, in which

211    available training sets can exceed hundreds of hours [13]. One possible reason for this

212    may be the high signal-to-noise ratio (SNR) of bat sounds recorded with high-quality

213    ultrasound devices. Previous studies have indicated that a high SNR can improve

214    separation accuracy [30] and our results suggest this model was suitable for use with

215    small, high-quality datasets. Although the sound data in this study were sampled in

216    controlled lab conditions, producing recordings that were essentially free of

217    background noise, acoustic analysis software could potentially optimize the separation

218    further by excluding any background noise that was present in the signal.

219         Future studies will also assess the performance of this network for other animal

220    species. Stereotypical patterns and clearly classifiable syllables have been observed in

221    the vocalizations of birds, non-human primates, whales, dolphins, and several other

222    species [31-33]. Features used in the proposed BLSTM were log spectral magnitudes,

223    which can be acquired from any vocal sound. This could potentially lead to robust

224    software that is not specific to a certain species or task. The model could also be

225    generalized to other animals, though limitations may exist. In addition to the quality

226    and quantity of training samples, hyper-parameters must be tuned in accordance with

227    the data [34, 35].

228    **Conclusion**

229         A sound separation model was proposed for extracting bat calls, achieving

230    excellent results. This is the first experimental evidence that the BLSTM model is

231    suitable for separating overlapping bioacoustic signals. These results provided a new

232    source for sound data analysis in animal acoustics research, which may contribute to

233  sample sizes and improve efficiency. This study also demonstrates the potential of

234  deep neural networks for applications to animal vocalization research, including

235  species classification and speech separation.

236  **Materials and Methods**

237  **Sound recording and data preparation**

238  **Species selection and sound sources.** Echolocation calls from bats are primarily

239  composed of constant frequency (CF) components and frequency modulated (FM)

240  components. Social calls are composed of CF, FM, and noise-burst (NB) components.

241  FM calls have short pulse durations and wide bandwidths. As such, they overlap with

242  social calls less in time but more in frequency. In contrast, CF calls have long pulse

243  durations and narrow bandwidths. They overlap with social calls more in time but less

244  in frequency. In consideration of the varied overlapping patterns found in bat calls, we

245  selected both CF bats (*Rhinolophus ferrumequinum*, *Hipposideros armiger*, and

246  *Rhinolophus pusillus*) and FM bats (*Vespertilio sinensis*, *Myotis macrodactylus*, and

247  *Ia io*) to test the separation capabilities of the proposed network, including six

248  different species to test method generalizability.

249      Source sound files from *V. sinensis*, *M. macrodactyllus*, *R. ferrumequinum*, *R.*

250  *pusillus*, and *H. armiger* were collected from previous studies in our lab (S1 Table).

251  Sound files for *Ia io* were selected from unpublished data as follows. Bats captured

252  from the field were housed in a husbandry room with abundant food and fresh water.

253  During each sound recording experiment, 4–5 bats were transferred to a temporary

254  cage. Sound recordings were collected using the Avisoft UltraSoundGate 116H

255    (Avisoft Bioacoustics, Berlin, Germany) and a condenser ultrasound microphone

256    (CM16/CMPA, Avisoft Bioacoustics). The sampling frequency was set to 375 kHz at

257    16 bits. The recording experiment lasted five days in order to acquire a sufficient

258    number of recordings, beginning at 18:00 and finishing at 6:00 the following morning.

259    S1 Table shows sample numbers and locations for the bats, as well as the total

260    duration of sound files selected for the study. All experimental procedures complied

261    with the ABS/ASAB guidelines for the Use of Animals in Research and were

262    approved by the Committee on the Use and Care of Animals at the Northeast Normal

263    University (approval number: NENU-W-2010–101).

264    **Sound analysis.** The total duration of recorded sound files (i.e., original recording

265    files) used for each bat species is shown in S1 Table. We employed Avisoft-SASLab

266    Pro (Version 5.2.12, Avisoft Bioacoustics, Berlin, Germany) to identify

267    non-overlapping and overlapping syllables in echolocation and communication calls.

268    These syllables and calls were described and classified following the nomenclature

269    developed by Kanwal, Matsumura (36) and Ma, Kobayasi (37). The recorded

270    non-overlapping calls were used for preparing training files of each call type and the

271    recorded overlapping calls were used for separation.

272    **Data preparation.** Supervised machine learning algorithms use training samples to

273    "learn" the steps required for completing a task. The training phase in this study

274    involved preparing clear and non-overlapping echolocation and communication calls,

275    selected from original recording sounds. In this process, the BLSTM network learned

276    features found in both call types.

277    Training samples consisted of randomly selected non-overlapping syllables in

278    echolocation and communication calls from each bat species (in the original

279    recordings), with signal-to-noise ratios (SNRs) above −20 dB. The echolocation

280    training files contained 1,300–6,240 pulses and the communication training files

281    contained 780–1,800 syllables (S1 Table). Although the quantity of selected syllables

282    varied between studies, the data was sufficient for model training. Efforts were made

283    to include roughly equivalent quantities of each syllable type. Time intervals between

284    syllables in the training files were consistent with those of the original recordings. The

285    lengths of training files for echolocation and communication calls were the same for

286    each bat species (S1 Table).

287    **Model training and call separation**

288    **Model structure and training stage.** We developed a network with four BLSTM

289    layers, followed by one feedforward layer (Fig 6). Each BLSTM layer included one

290    forward and one backward basic LSTM layer, both of which were added with dropout

291    functions (tensorflow.nn.rnn_cell.DropoutWrapper). Each BLSTM layer contained

292    300 hidden cells and the feedforward layer corresponded to the embedding dimension

293    (i.e., a 3D matrix with depth N=40 in this experiment). Stochastic gradient descent

294    with a momentum of 0.9 and a fixed learning rate of $10^{-3}$ was used for training. The

295    tanh activation function and the Adam optimizer were adopted to support adaptive

296    learning rates and faster convergence. The structure and hyper-parameters for the

297    model were designed based on the work of Hershey, Chen (21).

298    **Fig 6**. **The BLSTM model architecture and workflow graph.**

14

299     The model was trained using the files for one bat species in each trial.

300     Echolocation and communication call training files were loaded using the librosa

301     (version 0.6.2) Python package. Frames from the two sound files were read and added

302     together to create sound mixtures. Sound features used for training (log spectral

303     magnitudes) were extracted from this mixture. The extraction process was completed

304     using a short-time Fourier transform (STFT) with a Hamming window (length of 512

305     and shift of 256).

306     The mixture from each bat species was then segmented into 100-frame samples,

307     all of which were divided into a training set and a validation set using a ratio of 2:1

308     (see S1 Table for detailed sample quantities). The training set, validation set, and

309     indicator labels were combined and input to the model. The validation set was used to

310     optimize tuning parameters and evaluate call separation performance. Indicator labels

311     were set to 0 or 1, representing the two types of calls in the mixture. Ideal binary

312     masks were used to train the network and gradients were calculated using shuffled

313     mini-batches (batch size of 128) from larger segments.

314     The output of this model was a set of embeddings that included learned features

315     for both echolocation and communication calls. In this framework, the deep network

316     assigned embedding vectors to each time-frequency bin in the spectrogram. The

317     network then minimized the distance between embeddings dominated by the same

318     call type in each bin while maximizing the distance between embeddings dominated

319     by different call types. The output was then compared with the validation set and

320     indicator labels to calculate loss, which was back propagated from the output to the

321 input through each layer. Model weights and parameters were then updated based on

322 the calculated loss and training was completed after sufficient iteration epochs.

323 **Separation stage.** In this stage, overlapping echolocation and communication calls

324 were randomly selected from the original recordings to create a sound file of test sets,

325 used for separation. The log spectral magnitudes of the overlapping calls were then

326 extracted, combined into samples, and input to the trained model. The phases of calls

327 extracted from the sound files were also saved for use in sound reconstruction. The

328 trained model then output embeddings for each segment (100 frames) in a process

329 similar to the training stage. Embeddings were clustered using the k-means method

330 from Scikit-learn (Version 0.20.0) to produce time-frequency masks. The number of

331 clusters corresponded to the number of call types in the mixture (2 - echolocation and

332 communication). These masks and the clustering method were then used to determine

333 which parts of each segment in the overlapped calls would be preserved or neglected

334 based on their correspondence to each call type. For example, if the maximum

335 magnitudes were more likely to belong to echolocation calls, the related mask values

336 were set to 1 and the others were set to 0, allowing the echolocation calls to be separated

337 correctly. Finally, output calls were reconstructed using the inverse fast Fourier

338 transform (IFFT) function numpy.fft.ifft in NumPy (Version 1.15.1). The IFFT

339 transformed the magnitude into a wave using phase information saved at the beginning

340 of the separation stage. The model produced two waveform files, each containing one

341 call type. Additional detail concerning the sound separation algorithms can be found in

342 the work of Hershey (2016).

### Model evaluation

The quality of reconstructed echolocation and communication calls was assessed by comparing their temporal-spectrum parameters to the non-overlapping calls selected from the original recording files (excluding training data). Avisoft-SASLab Pro was used for automatic parameter measurements of duration, bandwidth, peak frequency, minimum frequency, maximum frequency, starting frequency, and ending frequency. A t-SNE (t-distributed stochastic neighbor embedding - R3.6.1 package) analysis was adopted for dimensionality reduction. Two dimensions were extracted from these seven parameters for original and separated syllables and compared with one-way ANOVA (aov in R3.6.1) or two-sided Wilcoxon signed-rank tests (wilcox.test in R3.6.1), depending on their fit to a normal Gaussian distribution. The significance level was set to 0.05 for all tests. We adopted the root mean square error (RMSE) to measure and avoid obscuring individual variations between reconstructed and original calls. Clustering analysis was conducted using the reconstructed echolocation calls from the six bat species, to assess whether the separated calls could be further used in species classification.

### Acknowledgements

### Author contributions

365  **Conceptualization**: Kangkang Zhang, Walter Metzner.

366  **Data curation**: Tong Liu, Shengjing Song, Xin Zhao.

367  **Formal Analysis**: Kangkang Zhang, Tong Liu.

368  **Methodology**: Kangkang Zhang, Ying Liu, Jiang Feng.

369  **Software**: Kangkang Zhang, Shijun Sun.

370  **Funding acquisition**: Ying Liu, Jiang Feng, Walter Metzner.

371  **Supervision**: Jiang Feng.

372  **Visualization**: Kangkang Zhang, Tong Liu.

373  **Writing – original draft**: Kangkang Zhang, Ying Liu.

374  **Writing – review & editing**: Ying Liu, Walter Metzner.

375  ## Competing interests

376  The authors have declared that no competing interests exist.

## References

1. Briggs F, Lakshminarayanan B, Neal L, Fern XZ, Raich R, Hadley SJK, et al. Acoustic classification of multiple simultaneous bird species: A multi-instance multi-label approach. 2012;131(6):4640-50. doi: 10.1121/1.4707424.

2. Aodha OM, Gibb R, Barlow KE, Browning E, Firman M, Freeman R, et al. Bat detective—Deep learning tools for bat acoustic signal detection. PLOS Computational Biology. 2018;14(3):156869.

3. Walters CL, Freeman R, Collen A, Dietz C, Brock Fenton M, Jones G, et al. A continental-scale tool for acoustic identification of European bats. Journal of Applied Ecology. 2012;49(5):1064-74. doi: https://doi.org/10.1111/j.1365-2664.2012.02182.x.

4. Adavanne S, Drossos K, Cakir E, Virtanen T. Stacked convolutional and recurrent neural networks for bird audio detection. european signal processing conference. 2017:1729-33. doi: 10.23919/EUSIPCO.2017.8081505.

5. Shamir L, Yerby C, Simpson R, Benda-Beckmann AMv, Tyack P, Samarra F, et al. Classification of large acoustic datasets using machine learning and crowdsourcing: Application to whale calls. The Journal of the Acoustical Society of America. 2014;135(2):953-62. doi: https://doi.org/10.1121/1.4861348.

6. Priyadarshani N, Marsland S, Castro I. Automated birdsong recognition in complex acoustic environments: a review. Journal of Avian Biology. 2018;49(5):jav-01447. doi: https://doi.org/10.1111/jav.01447.

7. Redgwell RD, Szewczak JM, Jones G, Parsons S. Classification of echolocation calls from 14 species of bat by support vector machines and ensembles of neural networks. Algorithms. 2009;2(3):907-24.

8. Sprengel E, Jaggi M, Kilcher Y, Hofmann T, editors. Audio based bird species identification using deep learning techniques. LifeCLEF 2016; 2016.

9. Kim G, Lu Y, Hu Y, Loizou PC. An algorithm that improves speech intelligibility in noise for normal-hearing listeners. The Journal of the Acoustical Society of America. 2009;126(3):1486-94. doi: 10.1121/1.3184603.

10. Wang Y, Wang D. Towards Scaling Up Classification-Based Speech Separation. IEEE Transactions on Audio, Speech, and Language Processing. 2013;21(7):1381-90. doi: 10.1109/TASL.2013.2250961.

11. Grais EM, Erdogan H, editors. Spectro-temporal post-smoothing in NMF based single-channel source separation. Signal Processing Conference (EUSIPCO), 2012 Proceedings of the 20th European; 2012: IEEE.

12. Nugraha AA, Liutkus A, Vincent E. Deep Neural Network Based Multichannel Audio Source Separation. Signals Commun Techn. 2018:157-85. doi: 10.1007/978-3-319-73031-8_7. PubMed PMID: WOS:000441018800008.

13. Wang D, Chen J. Supervised Speech Separation Based on Deep Learning: An Overview. IEEE Transactions on Audio, Speech, and Language Processing. 2018;26(10):1702-26. doi: 10.1109/Taslp.2018.2842159. PubMed PMID: WOS:000438718300001.

14. Marchi E, Ferroni G, Eyben F, Gabrielli L, Squartini S, Schuller B, editors. Multi-resolution linear prediction based features for audio onset detection with bidirectional LSTM neural networks. 2014 IEEE international conference on acoustics, speech and signal processing (ICASSP); 2014: IEEE.

15. Weninger F, Hershey JR, Le Roux J, Schuller B, editors. Discriminatively trained recurrent neural
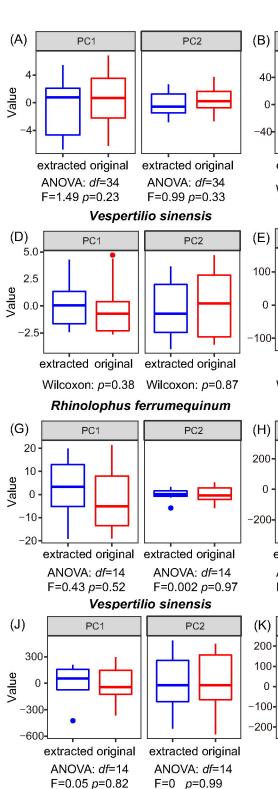
420    networks for single-channel speech separation. Proceedings 2nd IEEE Global Conference on Signal

421    and Information Processing, GlobalSIP, Machine Learning Applications in Speech Processing

422    Symposium, Atlanta, GA, USA; 2014.

423    16.   Eck D, Graves A, Schmidhuber J. A new approach to continuous speech recognition using LSTM

424    recurrent neural networks. Technical Report. 2003.

425    17.   Beringer N, editor Human language acquisition methods in a machine learning task. Eighth

426    International Conference on Spoken Language Processing; 2004.

427    18.   Graves A, Beringer N, Schmidhuber J, editors. A Comparison Between Spiking and Differentiable

428    Recurrent Neural Networks on Spoken Digit Recognition. international conference on modelling

429    identification and control; 2004.

430    19.   Graves A, Schmidhuber J. Framewise phoneme classification with bidirectional LSTM and other

431    neural network architectures. Neural Networks. 2005;18(5):602-10. doi:

432    https://doi.org/10.1016/j.neunet.2005.06.042.

433    20.   Makino S. Audio Source Separation: Springer; 2018.

434    21.   Hershey JR, Chen Z, Roux JL, Watanabe S, editors. Deep clustering: Discriminative embeddings

435    for segmentation and separation. 2016 IEEE International Conference on Acoustics, Speech and Signal

436    Processing (ICASSP); 2016 20-25 March 2016.

437    22.   Wöllmer M, Eyben F, Graves A, Schuller B, Rigoll GJCC. Bidirectional LSTM networks for

438    context-sensitive keyword detection in a cognitive virtual agent framework. 2010;2(3):180-90.

439    23.   Huang Z, Xu W, Yu K. Bidirectional LSTM-CRF models for sequence tagging. arXiv preprint

440    arXiv:150801991. 2015.

441    24.   Li C, Zhu L, Xu S, Gao P, Xu B, editors. CBLDNN-Based Speaker-Independent Speech

442    Separation Via Generative Adversarial Training. international conference on acoustics, speech, and

443    signal processing; 2018.

444    25.   Kunz TH, Fenton MB. Bat ecology: University of Chicago Press; 2005.

445    26.   Gillam E, Fenton MB. Roles of acoustic social communication in the lives of bats. Bat

446    Bioacoustics: Springer; 2016. p. 117-39.

447    27.   Luo B, Huang X, Li Y, Lu G, Zhao J, Zhang K, et al. Social call divergence in bats: a comparative

448    analysis. Behavioral Ecology. 2017;28(2):533-40. doi: 10.1093/beheco/arw184.

449    28.   Pozzi L, Gamba M, Giacoma C. The Use of Artificial Neural Networks to Classify Primate

450    Vocalizations: A Pilot Study on Black Lemurs. Am J Primatol. 2010;72(4):337-48. doi:

451    10.1002/ajp.20786. PubMed PMID: WOS:000276124300007.

452    29.   Jin L, Wang J, Zhang Z, Sun K, Kanwal JS, Feng J. Postnatal development of morphological and

453    vocal features in Asian particolored bat, Vespertilio sinensis. Mammalian Biology. 2012;77(5):339-44.

454    doi: 10.1016/j.mambio.2012.05.001. PubMed PMID: WOS:000309023900005.

455    30.   Weng C, Yu D, Seltzer ML, Droppo J. Deep neural networks for single-channel multi-talker

456    speech recognition. IEEE Transactions on Audio, Speech, and Language Processing.

457    2015;23(10):1670-9. doi: 10.1109/Taslp.2015.2444659. PubMed PMID: WOS:000356518100012.

458    31.   Naguib M, Riebel K. Bird song: a key model in animal communication. Encyclopedia for

459    language and linguistics. 2006;2:40-53.

460    32.   Filatova OA, Deecke VB, Ford JKB, Matkin CO, Barrett-Lennard LG, Guzeev MA, et al. Call

461    diversity in the North Pacific killer whale populations: implications for dialect evolution and

462    population history. Animal behaviour. 2012;83(3):595-603. doi:

463    https://doi.org/10.1016/j.anbehav.2011.12.013.

20

464  33.  Herzing DL. Clicks, whistles and pulses: Passive and active signal use in dolphin communication.
465  Acta Astronautica. 2014;105(2):534-7. doi: 10.1016/j.actaastro.2014.07.003.
466  34.  Lecun Y, Bengio Y, Hinton G. Deep learning. Nature. 2015;521(7553):436. doi:
467  10.1038/nature14539. PubMed PMID: 26017442.
468  35.  Goodfellow I, Bengio Y, Courville A. Deep Learning: The MIT Press; 2016.
469  36.  Kanwal JS, Matsumura S, Ohlemiller K, Suga N. Analysis of acoustic elements and syntax in
470  communication sounds emitted by mustached bats. The Journal of the Acoustical Society of America.
471  1994;96(3):1229-54. doi: https://doi.org/10.1121/1.410273.
472  37.  Ma J, Kobayasi K, Zhang S, Metzner W. Vocal communication in adult greater horseshoe bats,
473  Rhinolophus ferrumequinum. Journal of comparative physiology A, Neuroethology, sensory, neural,
474  and behavioral physiology. 2006;192(5):535-50. doi: https://doi.org/10.1007/s00359-006-0094-9.
475  PubMed PMID: 16418857.

476

# Supporting information

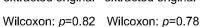**S1 Table. A summary of calls used for model training.**

**S2 Fig. Training loss and validation loss during model training.**

**S3 Table. Statistical comparisons of principle components extracted from seven parameters.** No significant differences were observed between parameters for separated and original syllables. A one-way ANOVA was used to test the normal distributed data and a two-sided Wilcoxon signed-rank test was used to assess the data that did not conform well to a normal distribution.

**V.sinensis**
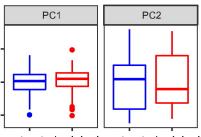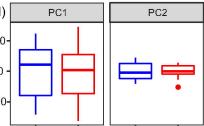
**Ia io**

**M.macrodactylus**

**R.ferrumequinum**

**H.armiger**

**R.pusillus**

(A) *Vespertilio sinensis*
ANOVA: *df*=34 F=1.49 *p*=0.23 (PC1)
ANOVA: *df*=34 F=0.99 *p*=0.33 (PC2)

(B) *Ia io*
Wilcoxon: *p*=0.82 (PC1)
Wilcoxon: *p*=0.78 (PC2)

(C) *Myotis macrodactylus*
Wilcoxon: *p*=0.20 (PC1)
Wilcoxon: *p*=0.94 (PC2)

(D) *Rhinolophus ferrumequinum*
Wilcoxon: *p*=0.38 (PC1)
Wilcoxon: *p*=0.87 (PC2)

(E) *Hipposideros armiger*
Wilcoxon: *p*=0.44 (PC1)
Wilcoxon: *p*=0.67 (PC2)

(F) *Rhinolophus pusillus*
Wilcoxon: *p*=0.81 (PC1)
Wilcoxon: *p*=0.31 (PC2)

(G) *Vespertilio sinensis*
ANOVA: *df*=14 F=0.43 *p*=0.52 (PC1)
ANOVA: *df*=14 F=0.002 *p*=0.97 (PC2)

(H) *Ia io*
ANOVA: *df*=20 F=0.002 *p*=0.97 (PC1)
ANOVA: *df*=20 F=0.002 *p*=0.97 (PC2)

(I) *Myotis macrodactylus*
ANOVA: *df*=28 F=0.05 *p*=0.82 (PC1)
ANOVA: *df*=28 F=0.02 *p*=0.89 (PC2)

(J) *Rhinolophus ferrumequinum*
ANOVA: *df*=14 F=0.05 *p*=0.82 (PC1)
ANOVA: *df*=14 F=0 *p*=0.99 (PC2)

(K) *Hipposideros armiger*
ANOVA: *df*=24 F=0.17 *p*=0.69 (PC1)
Wilcoxon: *p*=0.98 (PC2)

(L) *Rhinolophus pusillus*
ANOVA: *df*=18 F=0.33 *p*=0.58 (PC1)
ANOVA: *df*=18 F=0.01 *p*=0.91 (PC2)