# eNeuro

# Natural statistics as inference principles of auditory tuning in biological and artificial midbrain networks

*This Early Release article has been peer-reviewed and accepted, but has not been through the composition and copyediting processes. The final version may differ slightly in style or formatting and will contain links to any extended data.*

**Alerts:** Sign up at www.eneuro.org/alerts to receive customized email alerts when the fully formatted version of this article is published.

# Manuscript Title Page

**1. Manuscript Title (50 word maximum)**

Natural statistics as inference principles of auditory tuning in biological and artificial midbrain networks

**2. Abbreviated Title (50 character maximum)**

Auditory tuning in biological and artificial brain

**3. List all Author Names and Affiliations in order as they would appear in the published article**

Sangwook Park, Department of Electrical and Computer Engineering, Johns Hopkins University

Angeles Salles, Department of Psychological and Brain Sciences, Johns Hopkins University

Kathryne Allen, Department of Psychological and Brain Sciences, Johns Hopkins University

Cynthia F. Moss, Department of Psychological and Brain Sciences, Johns Hopkins University

Mounya Elhilali, Department of Electrical and Computer Engineering, Johns Hopkins University

**4. Author Contributions:**

Designed research SP, ME. Performed research: SP, AS, KA. Analyzed data: SP. Wrote original draft: SP. Revised manuscript: SP, AS, KA, CFM, ME. Funding: CM, ME.

**5. Correspondence should be addressed to (include email address)**

Mounya Elhilali (E-mail: mounya@jhu.edu)

**6. Number of Figures**: 14          **7. Number of Tables**: 1

**8. Number of Multimedia**: 0          **9. Number of words for Abstract**: 250

**10. Number of words for Significance Statement**: 120     **11. Number of words for Introduction**: 763

**12. Number of words for Discussion**: 1121

**14. Conflict of Interest**

Authors report no conflict of interest

39

40    Natural statistics as inference principles of auditory tuning in biological and artificial midbrain networks

41

42    **Abstract**

43    Bats provide a powerful mammalian model to explore the neural representation of complex sounds, as
44    they rely on hearing to survive in their environment. The inferior colliculus (IC) is a central hub of the
45    auditory system that receives converging projections from the ascending pathway and descending
46    inputs from auditory cortex. In this work, we build an artificial neural network to replicate auditory
47    characteristics in IC neurons of the big brown bat. We first test the hypothesis that spectro-temporal
48    tuning of IC neurons is optimized to represent the natural statistics of conspecific vocalizations. We
49    estimate spectro-temporal receptive fields (STRF) of IC neurons and compare tuning characteristics to
50    statistics of bat calls. The results indicate that the FM tuning of IC neurons is matched with the statistics.
51    Then, we investigate this hypothesis on the network optimized to represent natural sound statistics and
52    to compare its output with biological responses.  We also estimate biomimetic STRF's from the artificial
53    network and correlate their characteristics to those of biological neurons. Tuning properties of both
54    biological and artificial neurons reveal strong agreement along both spectral and temporal dimensions,
55    and suggest the presence of nonlinearity, sparsity and complexity constraints that underlie the neural
56    representation in the auditory midbrain. Additionally, the artificial neurons replicate IC neural activities
57    in discrimination of social calls, and provide simulated results for a noise robust discrimination. In this
58    way, the biomimetic network allows us to infer the neural mechanisms by which the bat's IC processes
59    natural sounds used to construct the auditory scene.

60

61    **Significance Statement**

62    Recent advances in machine learning have led to powerful mathematical mappings of complex data.
63    Applied to brain structures, artificial neural networks can be configured to explore principles underlying
64    neural encoding of complex stimuli. Bats use a rich repertoire of calls to communicate and navigate their
65    world, and the statistics underlying the calls appear to align with tuning selectivity of neurons. We show
66    that artificial neural network with a nonlinear, sparse and deep architecture trained on the statistics of
67    bat communication and echolocation calls results in a close match to neurons from bat's inferior
68    colliculus. This tuning optimized to yield an effective representation of spectro-temporal statistics of bat
69    calls appears to underlie strong selectivity and noise invariance in the inferior colliculus.

70

71    **1.  Introduction**

72    Biological neural circuits are believed to provide an efficient code of the sensory world, which allow us
73    to process complex and dynamic stimulus information from our surroundings. Perception of an auditory
74    scene is created by neural activity filtered through several stages of feed-forward and feedback sensory
75    processing. Sound pressure of an acoustic signal is first transduced into a bio-electrical signal in the
76    cochlea. Subsequently, the bio-electrical signal is relayed through the auditory pathway. The inferior
77    colliculus (IC) is an auditory hub that receives ascending inputs from brainstem nuclei and sends

78    information through the thalamus to the auditory cortex, while it also receives descending inputs from
79    auditory cortex (Casseday et al., 2002). The IC encodes complex auditory features such as frequency
80    sweep rate (Williams and Fuzessery, 2010) and patterning (Gordon and O'Neill, 1998) that are necessary
81    for identification of complex auditory objects and therefore plays a key role in representing these
82    objects in a natural listening environment.

83    Echolocating bats build a representation of their surroundings by emitting ultrasonic vocalizations and
84    processing the features of returning echoes to compute the location and features of targets and
85    obstacles in the environment. Bats must rapidly process sonar echoes while concurrently parsing
86    environmental noise and calls emitted by conspecifics. In this complex and rapidly changing auditory
87    scene, the bat's brain efficiently encodes acoustic stimuli and allows the animal to accurately track prey,
88    avoid obstacles, and communicate with conspecifics while dynamically navigating a three-dimensional
89    environment. Humans and other animals face similar challenges in the course of their natural acoustic
90    behaviors. With the goal of elucidating principles underlying auditory scene analysis in the midbrain, we
91    examine the relationship between statistics of the rich acoustic repertoire of bat calls and neural
92    response patterns in the bat's IC to explore artificial networks tuned to map natural statistics in these
93    calls and identify emergent properties that match responses in the IC.

94    Here, we test the hypothesis that the bat's auditory midbrain is optimized to accurately represent the
95    natural statistics in the sounds and echoes that exist in the bat's environment (particularly social and
96    echolocation calls). Past research has suggested that the IC plays a major role in the representation and
97    mapping of communication sounds that give rise to specialized encoding of natural sounds along the
98    ascending auditory system (Aitkin et al., 1994; Suta et al., 2003). An earlier study in the Mexican free
99    tailed bat suggested a possible correspondence between tuning characteristics of individual IC neurons
100   and properties of natural calls from conspecific sounds (Andoni et al., 2007; Brimijoin and O'Neill 2005).
101   In the current study, we corroborate this relationship in a different species and further probe constraints
102   and implications of such optimal encoding of natural sounds on auditory signal processing in a complex
103   scene.

104   We recorded vocalizations from socially housed bats and analyzed the spectro-temporal statistics of
105   natural sounds (e.g. frequency modulation (FM) velocity, directionality). Using the database of collected
106   statistics, we built an artificial network, which projects sounds onto a latent space that efficiently
107   represents statistics of these natural sounds in a strategy of signal reconstruction (Smith and Lewicki,
108   2006). This computational model offers a biomimetic architecture whose main operation is to capture
109   the statistics of natural bat calls, without information about the function of biological neurons. We then
110   ask: Does the emergent tuning of this artificial network match properties of biological neurons in the big
111   brown bat inferior colliculus? To answer this question, we also recorded responses to sound stimuli,
112   spectro-temporal ripples, from individual neurons in the IC of big brown bats.

113   It is known that the spectro-temporal receptive fields (STRF's) suggest a reasonable linear-
114   approximation of neural responses as a transfer function from acoustic stimuli and those are usually
115   used to explore auditory characteristics of the neurons (Andoni et al., 2007; Depireux et al., 2001; Elhilali
116   et al., 2013). We extracted STRF's from IC and artificial neurons and calculated auditory characteristics
117   from these neural response functions (Kowalski et al., 1996; Poon and Yu, 2000). The spectro-temporal
118   tuning characteristics of biological neurons were then compared to both the statistics of natural calls as
119   well as emergent tuning of artificial neurons. By varying the configuration of the artificial network, we

4

120  employed the theoretical network as springboard to examine possible constraints on the configuration
121  of midbrain networks, and gauge the validity of the hypothesis linking biological encoding in the
122  mammalian midbrain to efficient representation of natural sound statistics. While various artificial
123  neural networks can be optimized to reconstruct an input sound from compressed feature on latent
124  space, finding an architecture that closely emulates the biological network provides insights into the
125  underlying functional role of certain brain nuclei. Here, we examine the relationship between the
126  optimal encoding of natural statistics in bat calls and its role in facilitating robust selectivity across sound
127  classes in the repertoire. The graphical abstract in Fig. 1 shows an overview of the approach taken in this
128  work.

129

130  **2.  Materials and Methods**
131  **2.1. Collection of Bat's vocalization**
132  2.1.1.  Animals

133  Big brown bats (*Eptesicus fuscus*) were collected from an exclusion site under a state permit. All
134  experimental procedures were carried out in accordance with a protocol approved by an Institutional
135  Animal Care and Use Committee. A total of approximately 100 bats were housed in our Lab and used for
136  vocal data recordings, and four (2 male, 2 female) bats were used for neurophysiological data collection.

137

138  2.1.2.  Audio recordings for training the biomimetic network

139  A bat call library was built from audio recordings of bats housed in a vivarium room where the
140  temperature is kept at 70-80 $^o$F, and humidity is kept at 30-70 %. This room holds approximately 100
141  bats in groups of 1-6 separated in mesh cages. The recordings were made for two days using an Avisoft
142  CM16/CMPA ultrasonic microphone and the Avisoft-RECORDER software. Mono audio was recorded at a
143  sampling rate of 300 kHz.

144   Natural call recordings from big brown bats were processed to extract meaningful segments. An
145  energy-based signal activity detection was performed on the entire database to remove the silences
146  between calls and to split the recordings into segments containing bat calls (Park et al., 2014). As a
147  result, we constructed species specific databases containing 17,713 calls (about 10 min) for big brown
148  bats. This call database was used for training artificial networks. The data was divided into a training set
149  (15,000 randomly selected calls) to learn network parameters and test set (remaining 2,713 calls) for
150  verifying the network.

151

152  2.1.3.  Social calls for natural sound representation

153  To investigate discriminability in the artificial network, we used a social call database that includes 26
154  audio clips for 8 different types of bat calls (Fig. 2). These types include six calls, as defined in (Wright et
155  al., 2013), specifically, Echolocation (Echo), Frequency Modulated Bout (FMB), Upward Frequency
156  Modulated (UFM), Long Frequency Modulated (LFM), Short Frequency Modulated (SFM), and Chevron-
157  Shaped (CS); in addition to two additional calls types, Long-Wave and Hook, which resemble a hook in
158  time-frequency space. All audio clips were up-sampled from 250 kHz to 300 kHz.

5

159

### 2.2. Neurophysiological Inferior Colliculus data

161 Recordings of neural responses from IC neurons were used to perform two separate analyses: (1)
162 characterize receptive field tuning of IC neurons; and (2) examine discriminability of IC neurons to
163 different con-specific calls. Methods for receptive field analysis are described next in section 2.2.4, while
164 data used for discriminability analysis are described in section 2.2.5.

165

166     2.2.1.   Receptive field recordings

167 A head-post was adhered to the skull of bats for head fixation as described in (Macias et al., 2018). The
168 inferior colliculus was located using skull and brain landmarks and a surgical drill was used to make a ≤ 1
169 mm diameter craniotomy preserving dura. The neurophysiological recordings were performed in a
170 sound-attenuating and electrically shielded chamber (Industrial Acoustics Company, Inc.). Each bat was
171 restrained individually in a custom-made foam mold and the head was fixed by the head-post. Recording
172 sessions were carried out over 3 to 5 consecutive days, each one lasting no more than 4 hours. Water
173 was offered to the bats every 2 hours. No drugs were administered during recordings. During recordings
174 a silver wire for grounding was placed in between muscle and skull about 5mm rostral to the craniotomy
175 site. The 16-channel recording probe (Neuronexus A1x16-5mm-50-177-A16) was inserted into the brain
176 using a micromanipulator. The surface of the brain was registered as 0 μm for depth reference and the
177 probe was advanced in 10 μm steps using a hydraulic microdrive (Stoelting Co.). Recordings were taken
178 at least 100 μm apart. An OmniPlex D Neural Data Acquisition System recording system (Plexon, Inc.)
179 was used to obtain neural responses with 16-bit precision and 40 kHz sampling rate. A transistor-
180 transistor-logic (TTL) pulse for each stimulus presentation was generated with the National Instrument
181 card used for stimulus presentation and was recorded on channel 17 of the analog channels of the
182 acquisition system for synchronization of acoustic stimuli and neural recordings. The stimuli were
183 recorded on channel 18 of the acquisition system to corroborate synchronization.

184

185     2.2.2.   Moving ripple stimuli

186 A set of ripple stimuli was generated to estimate STRFs of IC neurons (Kowalski et al., 1996; Depireux et
187 al. 2001; Andoni et al., 2007). Ripples are modulated noise stimuli that are dynamic both in time and
188 frequency. Each ripple can be described as

189 $$S(t,x) = 1 + \Delta A \times \sin(2\pi(\omega t + \Omega x) + \phi) \qquad (1)$$

190 where $t$ and $x$ are indices for time and octave scaled frequency. $\Delta A$ and $\phi$ are amplitude and a phase,
191 respectively. And $\omega$ and $\Omega$ represent modulation rates along temporal (Hz) and spectral (cyc/oct) axes.
192 The temporal and spectral modulation parameters were varied from -176 - 176 Hz in steps of 32 Hz and
193 0.0 - 1.5 cyc/oct in steps of 0.15 cyc/oct spectrally (Fig. 3A). Each ripple spanned 6.66 octaves from 1.2
194 kHz to 121 kHz and was 300 ms in duration.

195

196     2.2.3.   Audio playbacks for neural recordings

197 Extracellular recordings from the inferior colliculus of awake animals were taken while they passively
198 listened to broadcast of either ripple stimuli, or pure tones at 70 dB. All stimuli were generated at a
199 sampling rate of 250 kHz using a National Instruments card (PXIe 6358) and transmitted with a
200 calibrated custom-made electrostatic ultrasonic loudspeaker connected to an audio amplifier (Krohn-
201 Hite 7500). The loudspeaker was placed at 60 cm (for all ripple and pure tones stimuli)} from the bat's
202 ear. The frequency response of the loudspeaker was compensated by digitally filtering the playback
203 stimuli with the inverse impulse response of the system as described in (Luo and Moss, 2017).

204 Frequency tuning curves were built by recording neural responses to pure tones of 5 ms duration (with
205 0.5 ms ramping rise and fall). The tones ranged between 20 and 90 kHz (in 5 kHz steps) and the sound
206 pressure levels ranged from 20 to 70 SPL (10 dB steps). At each recording site first, we played 20
207 repetitions of the randomized ripple stimulus and then 15 repetitions of each of the randomized pure
208 tones at a different SPL.

209

210     2.2.4.   Analysis of neuronal responses

211 For the analysis of auditory tuning in response to ripple and pure tone stimuli, responses were sorted
212 offline, then single units were detected using the program 'Wave_clus' (Quiroga et al., 2004). Each
213 individual waveform was inspected and the acceptance threshold for clusters was less than 10% of
214 spikes with < 3 ms inter-spike interval, consistent with the neuronal refractory period. Any sites that
215 showed no response to ripple stimuli were excluded from the spike sorting and further analysis in line
216 with procedures used in other studies (Poon and Yu, 2000; Escabi and Schreiner, 2002; Andoni et al.,
217 2007). After spike sorting, the Euclidian distance error between the mean and variance of number of
218 spikes across trials was computed. Units whose error is less than 1.0 were selected for further analysis,
219 following a Poisson model of spike representation (Corrado et al., 2005; Schwartz et al., 2006). This
220 analysis resulted in 108 single units used for the current study.

221 **Neurophysiological STRFs:** At each recording site, ripple stimuli were repeated 10-20 times in a
222 randomized order for each repetition. A PST histogram was calculated from the spike time sequence of
223 each ripple; then histograms were folded into 32-point periods. The strength and phase of the response
224 to each ripple were estimated directly from the fundamental component obtained by applying a 32-
225 point Fast Fourier Transform (FFT) to the period histogram. Magnitude and phase responses to each
226 ripple were combined together into a magnitude matrix $M(\Omega, \omega)$ and a phase matrix $\Phi(\Omega, \omega)$,
227 respectively. To derive a Ripple Transfer Function (RTF), which is a representation of a STRF in the
228 modulation domain, $M(\Omega, \omega)$ and $\Phi(\Omega, \omega)$ were expanded to four quadrants in the modulation domain
229 spanning from -176 Hz and -1.5 cyc/oct to 176 Hz and 1.5 cyc/oct as $M_e(\Omega, \omega) = M_e^*(-\Omega, -\omega) =$
230 $M(\Omega, \omega)$ and $\Phi_e(\Omega, \omega) = \Phi_e^*(-\Omega, -\omega) = \Phi(\Omega, \omega)$ based on a symmetric property around the origin
231 (Depireux et al., 2001; Andoni et al., 2007). As a result, the RTF was formulated as

232
$$T(\Omega, \omega) = M_e(\Omega, \omega)e^{j\Phi_e(\Omega,\omega)} \qquad (2)$$

233 where $j = \sqrt{-1}$. Finally, a STRF was obtained by performing 2D inverse FFT on the RTF as

234
$$STRF(x, t) = F_{t,-x}^{-1}[T(\Omega, \omega)] \qquad (3)$$

235 where $F^{-1}$ designates the 2D inverse FFT along each axis in the modulation domain.

236

237     2.2.5.    Neural discriminability of con-specific calls

238    In order to examine selectivity of IC neurons to calls from the bat's natural repertoire, we re-used neural
239    data previously collected in an earlier study (Salles et al., 2020), where we collected neuronal responses
240    to Echolocation calls (Echo) versus Frequency-Modulated Bout (FMB) social calls. The study followed the
241    same methodology for data collection as described here. 'Wave_clus' was used to detect and classify
242    single units from the recordings. The spikes responding to either FMB or Echo were counted in windows
243    of 25 ms duration, starting 5ms after stimulus onset. Some units with an average of less than five spikes
244    over 20 times recordings were excluded because they were considered as a non-responsive unit to the
245    stimulus. Multi-unit activity was determined from inter-spike intervals with < 3ms that were inconsistent
246    with neuronal refractory period; and units with greater than 10% of spikes with < 3ms inter-spike
247    interval were excluded from analysis. As a result, total 575 units were finally obtained and their
248    responses are used in the present work to contrast neural discriminability between Echo and FMB calls
249    with artificial neurons.

250

251    **2.3. Responses in artificial neurons**
252     2.3.1.    Artificial network front-end processing

253    To develop a biomimetic architecture, a biologically-inspired auditory spectrogram is used as input for
254    the network (Shamma, 1985a; Shamma, 1985b; Yang et al., 1992; Wang and Shamma, 1994). The
255    auditory spectrogram incorporates four processing stages that emulate peripheral processing in the
256    mammalian system: cochlear filtering, auditory-nerve transduction, hair cell responses, and lateral
257    inhibition (Chi et al., 2005). Briefly, an incoming acoustic waveform is analyzed along a bank of constant-
258    Q filters spanning a logarithmic scale. Then, each frequency channel undergoes a high-pass, nonlinear
259    compression and low-pass filtering followed by lateral inhibition across frequency, following the
260    implementation available in the NSL toolbox (Chi and Shamma, 2005) with the following settings: The
261    frame length was set to 0.2 ms without overlap, and each octave was represented with 24 channels (i.e.
262    128 channels over 5.33 octaves). Octave-scaled center-frequencies were represented as $f_c = 440 \times$
263    $2^{((c-32)/24+\gamma)}$ where $f_c$ is a center frequency of the $c^{th}$ channel, and $\gamma$ is a constant factor of octave
264    shift ($\gamma = 4.38$). Inputs to the artificial network were sampled as square patches of the spectrogram
265    spanning 128 frequency channels (i.e. 5.33 octaves) and 160 time-samples (i.e. 32 ms).

266

267     2.3.2.    Structure of artificial network

268    An artificial neuron, i.e. node mimicking a biological neuron, is mathematically modelled by a linear
269    combination of pre-node outputs and a non-linear activation function. An artificial network is
270    constructed by connecting a large number of nodes to each other. Using nonlinear activation functions
271    enables the network to perform nonlinear computations on feedforward propagation. For this study, we
272    favored a generative architecture using an autoencoder composed of an encoder, which compresses
273    original data into a compact code; and a decoder, which reconstructs the original signal from that code
274    (Baldi, 2012; Doersch, 2016). The intuition is to directly test our hypothesis that the network would infer

275    a statistical model of the training dataset of natural calls, and if successful should allow a faithful
276    reconstruction of the inputs.

277    The proposed architecture is shown in Fig. 4. First an encoder stage **E** is composed of convolutional
278    layers, pooling layers, and a fully connected layer. A latent vector represents compressed features
279    learned from the input data. A decoder stage **D** composed of reverse operations using transposed
280    convolutions, reconstructs the input features from a latent vector. A sampling stage, interposed
281    between the encoder and decoder, emulates neural activity yielding sparse binary activations.

282    Using the same general building block composed of convolution and pooling layers, this study
283    investigates various configurations of the network by varying: (1) *depth*, which is the number of blocks.
284    In Fig. 4, the black-flow shows a double stacking structure as an example. A deeper network can be
285    constructed by stacking more blocks, on the other hand, a shallow network can be created by removing
286    a block; (2) *nonlinearity*, by varying the slope of nonlinear activation function employed; and (3) *sparsity*,
287    by controlling the density of sampling in the latent space.

288    The encoder architecture **E** follows a convolutional neural network (CNN) framework in order to reduce
289    the number of trainable parameters, hence controlling for over-fitting issues and generalizability to
290    unseen data (Dietterich, 1995). The convolutional layers compute output feature maps using 2D
291    convolutions between input feature maps and several filters as

292    $$I_o^l[f,t,k] = \sum_{\xi,\tau,m} I_i^l[\xi,\tau,m] f^l[\xi-f,\tau-t,m,k] \qquad (4)$$

293    where $f, t, l, k$ and $m$ are indices for spectral, temporal, layer, channel of output feature map, and
294    channel of input feature map respectively. $I_i$, $I_o$, and $f^l$ are feature maps for input and output, and
295    convolutional filter applied in the $l^{th}$ layer, respectively. Multi-scale filters are employed in each
296    convolutional layer to balance broad span (in time and frequency) vs. localized analyses. Then, output
297    feature maps concatenate filter outputs using multi-sized filters (Fig. 5A) (Szegedy et al., 2015). Specifics
298    of both filter composition and dimensions of intermediate feature maps are summarized in Table 1.
299    Neural activation by an acoustic feature is emulated by applying a nonlinear function after convolution
300    as

301    $$I_a^l[f,t,k] = \max(I_o^l[f,t,k], \ \alpha \times I_o^l[f,t,k]) \qquad (5)$$

302    where $\alpha$ is a constant within an interval [0, 1] (Maas et al., 2013). Next, pooling layers compress the
303    output from the previous convolutional layer by extracting a maximum among some values enclosed by
304    a non-overlapping window (i.e. max-pooling) $I^l$ (Scherer et al., 2010). As a result, the width and height
305    of the output are reduced by half. At the top of the encoder, a fully connected layer is applied for
306    mapping into a latent space, which involves natural statistics requiring to reconstruct original input, as
307    $v_c = W_l^T \times flatten(I^L)$ where $I^L$ is a feature map in the last pooling layer, $W$ is weight matrix in the
308    fully connected layer, and $flatten(.)$ is a reshape function from a 3D tensor to a vector.

309    In the middle stage, a binary code vector $v_b$ is generated by performing a Bernoulli sampling process. A
310    sigmoid function is applied to the latent vector to calculate prior probabilities. Thus, the output of the
311    middle stage is represented as $v_b = Bernoulli(\sigma(v_c))$ where $\sigma(.)$ is a sigmoid function.

312    The decoder **D** is composed of a fully connected layer and transposed convolution layers. In the fully
313    connected layer, a latent vector is expanded into an initial space as $\hat{v} = W_l \times v_b$, and the vector $\hat{v}$ is

9

314 reshaped to a 3D tensor as a set of initial feature maps as $\hat{I}^l = reshape(\hat{v})$. From initial feature maps, a
315 transposed convolution using multi-scale filters is sequentially performed until the output has the same
316 dimensions as the input patch (Shelhamer et al., 2017; Radford et al., 2015). Convolutional filters used in
317 the encoder are applied for transposed convolution after transposing input channel from output
318 channel dimension as $\dot{f}^l[f, t, k, m]$. A transposed convolution using multi-scale filters is performed in
319 three steps (Fig. 5B). First, the input feature map $\hat{I}^l$ is split into submaps, $[\hat{I}_1^l, \hat{I}_2^l, ..., \hat{I}_N^l]$, as many as the
320 number of filters. Second, transposed convolution is individually performed for each pair of submap and
321 filter. Finally, a set of output feature maps is obtained by averaging the results of the second step.

322

323     2.3.3.   Training artificial network

324 The network was trained using the cost function:

325 $$L = \frac{1}{2}\sum_n[(x_n - D(E(x_n)))^2 + \lambda(\rho - \sum_i \sigma(v_{c_i}))^2]\qquad(6)$$

326 where $x_n$ is an input patch with respect to the $n^{th}$ index, $E(.)$ represents an encoder function while
327 $D(.)$ is for a decoder, and $\rho$ means the average number of active nodes. The first term represents the
328 mean square error between an input patch and its reconstruction by the autoencoder. The sparse
329 constraint prevents overfitting as well as emulates sparsity of active neurons in the brain. Let $Y$ be a
330 random variable representing the number of active nodes by the Bernoulli process. Then, the
331 distribution known as the Poisson binomial distribution is denoted as

332 $$\Pr(Y = \rho) = \sum_A[\prod_{i \in A} \sigma(v_{c_i}) \prod_{j \in A^c}(1 - \sigma(v_{c_j}))]\qquad(7)$$

333 where $A$ is a set whose elements are possible combination for choosing $\rho$ nodes from $N$ nodes. This
334 distribution can be approximated by $Binomial(N, \mu/N)$ where $\mu = \sum_i \sigma(v_{c_i})$ (Choi and Xia, 2002). The
335 network training was implemented using TensorFlow (Abadi et al., 2016). AdamOptimization was
336 applied for an optimizer with $1.0e - 4$ learning rate. And, $\lambda$ was set to $1.0e - 4$. For more details,
337 readers can find the implementation on http://www.github.com/JHU-LCAP/BioSonar-IC-model /.

338   Comparisons between the biological neurons and artificial neurons were performed to infer the
339 network configuration that best matches the characteristics of IC neurons (as explained next). The best
340 configuration composed of a triple stacking network, a parameter of nonlinearity $\alpha = 0.2$ in (5), and 10%
341 sparsity constraint in (6).

342

343     2.3.4.   Biomimetic STRFs

344 Once trained, the network was interrogated following the same procedure as biological neurons. The
345 same ripple stimuli were given as input to the network and activity of the nodes before applying the
346 sigmoid activation and the Bernoulli sampling, $v_c$ in Fig. 4 was characterized. Each ripple was
347 transformed into an auditory spectrogram (as described earlier). A sequence of input patches for each
348 ripple were then composed by applying a sliding window (window length: 160 frames) in every 2 ms
349 (sliding step: 10 frames) (Fig. 3B). Input patches in the sequence were consecutively fed into the pre-
350 trained encoder, then a latent vector $v_c$ was obtained every 2 ms. The same procedure for extracting

351  biological STRF's was followed (See. Section 2.2.4). To find the magnitude $m$ and phase $\phi$ of the
352  responses, we performed a 32-point Fast Fourier Transform (FFT) and derived the magnitude and
353  unwrapped phase of the fundamental component (Fig. 3C). By repeating this procedure for all ripples,
354  the magnitude and phase were collected in a matrix $M(\Omega, \omega)$ and a $\Phi(\Omega, \omega)$, respectively (Fig. 3D).
355  These modulation responses were then converted into time-frequency STRF profiles by performing a 2D
356  inverse FFT on the RTF (Fig. 6). Note that, in this study, all network architectures employed a total 100
357  artificial neurons (spanning a 100-dimensional latent space) so that 100-biomimetic STRFs were used for
358  analysis.

359

360  **2.4. Analysis of auditory characteristics**
361  2.4.1.  Natural statistics and Auditory characteristics

362  *Frequency Modulation (FM) velocity (statistics of bat calls)*: To characterize conspecific vocalizations, we
363  calculated FM velocities of each call segment in our database. Since moving ripples were used as bases
364  components of the Fourier modulation domain (Singh and Theunissen, 2003), we derived auditory
365  spectrograms of each call, then performed a 2D FFT after mean subtraction to remove constant
366  components. $T_c(\Omega, \omega) = F_{f,t}[S(f,t) - \bar{S}]$ where $F_{f,t}[.]$ is the 2D FFT, $S$ is an auditory spectrogram of a
367  bat call, and $\bar{S}$ is its mean over the time and frequency axes. A velocity line was estimated by performing
368  a line fitting on the magnitude of 2D FFT result. Finally, the FM velocity of a bat call was acquired by
369  calculating the slope of the velocity line.

370  *Best Velocity (BV)*: We defined a best velocity as the center of mass with respect to response power in a
371  magnitude plot. To estimate the center of mass, we performed a Gaussian surface fitting on the 1st
372  quadrant of magnitude plot. After normalization as $\bar{M}_e = M_e/[\sum_{\Omega,\omega} M_e \Delta\Omega\Delta\omega]$ where $\Delta\omega$ and $\Delta\Omega$ are
373  respectively step size of temporal and spectral modulation rate, the fitting was performed to estimate
374  mean vector and covariance matrix, by minimizing a square mean error function as
375  $Err = \frac{1}{2}\sum_{\Omega,\omega}(ln(\bar{M}_e) - ln(G_{\mu,\Sigma}))^2$, where $G_{\mu,\Sigma}$ is a Gaussian distribution with mean vector $\mu$ and
376  covariance matrix $\Sigma$. By performing the Least Square Error (LSE) estimator iteratively (Kay, 1993), we
377  derived the Gaussian mean vector and covariance matrix. Best Velocity was defined as the slope of the
378  mean vector (Fig. 7).

379  *Orientation (Ori)*: To characterize velocity selectivity, we defined orientation as the angle between a
380  line connecting the origin to the center of mass and a dominant eigenvector of the Gaussian covariance
381  matrix. Note that the dominant eigenvector indicates the dominant direction of magnitude spread at
382  the center of mass (Fig. 7) (Andoni et al., 2007).

383  *Inseparability (Ins)*: Singular value decomposition (SVD) is applied to each STRF for calculating
384  inseparability (Depireux et al., 2001). This approach decomposes the STRF into a linear combination of
385  rank-1 matrices; in other words, $STRF = \sum_i \lambda_i u_i v_i^H$, where $u_i$ and $v_i$ are respectively left- and right-
386  eigenvectors (column vector) corresponding to a singular value $\lambda_i$, and $H$ means Hermitian transpose
387  (Strang, 2009). Based on this definition, a STRF is called separable if the STRF can be approximated by
388  summation of just a few matrices otherwise it is inseparable. We measured inseparability of a STRF
389  calculated as $Ins = 1 - \lambda_1^2/\sum_i \lambda_i^2$, where $\lambda_1 \geq \lambda_2 \geq \lambda_3 \geq \cdots$. Note that the inseparability is bounded
390  within the interval [0,1] where the $Ins$ is equal to 0 for separable STRFs otherwise it goes to 1.

11

391    *Direction Selectivity Index (DSI):* To investigate direction selectivity of STRF's, we compared total power
392    in the 1st and the 2nd quadrant of the RTF. If a STRF favors downward-moving ripples, total power in
393    the 1st quadrant of magnitude plot is larger than the other since the 1st quadrant is composed of
394    responses evoked by downward-sweeping ripples. From this perspective, a DSI was defined as
395    $DSI = (P_2 - P_1)/(P_2 + P_1)$ where $P_i$ is a power in the $i^{th}$ quadrant of RTF, and it is calculated by
396    $P_i = \sum_{(\Omega,\omega)\in Q_i} |T(\Omega,\omega)|$ where $Q_i$ means the $i^{th}$ quadrant. Since the power on each quadrant is a non-
397    negative value, the DSI is bounded within the interval [-1,1] where downward/upward selectivity is
398    represented to negative/positive DSI while 0 represents no selectivity in the direction. DSI for natural
399    vocalizations was derived using the Fourier representation described to derive FM velocity.

400    *Best Frequency (BF)*: To investigate frequency selectivity of STRFs, we defined a BF as the frequency of
401    the maximum peak of absolute STRF, $|STRF|$ over the entire time and frequency spans. Best frequency
402    (spectral peak) of natural calls was computed by finding the peak frequency of the average spectrum.

403

404        2.4.2.    A bootstrap for statistical comparison}

405    We performed a bootstrap analysis to evaluate similarity between distributions of characteristics (e.g.
406    FM velocity, BV) comparing natural calls, IC neurons, and artificial neurons. The procedure selects
407    random 30 samples from natural calls in each iteration with replacement. For IC neurons, random
408    samples from each of the 4 bats are used in each iteration to maintain a balanced representation across
409    bats. In case of artificial neurons, we trained 10 independent-networks (using different initialization
410    procedures) and combined the neurons from each network into a complete set that was then sampled
411    during the bootstrap procedure. For each comparison and each bootstrap repetition, the distance
412    between means was noted. 1000 repetitions were used to generate a distribution of mean distances
413    $d_{(\mu,\sigma)}$ where $\mu$ and $\sigma$ are the mean and standard deviation. The *p-value* for accepting null hypothesis
414    was calculated as $p = 1 - 2\int_0^{|\varepsilon|} d_{(0,\sigma)}(x)dx$ where $d_{(0,\sigma)}$ is a zero-mean Gaussian distribution with
415    same variance $\sigma$, and $\varepsilon$ was a real number satisfying $d_{(0,\sigma)}(\varepsilon) = d_{(\mu,\sigma)}(\varepsilon)$.

416

417        2.5. **Natural sound representation with artificial neurons**
418        2.5.1.    Analysis of response selectivity in artificial neurons

419    We explored response selectivity to bat calls in biomimetic neurons. To replicate the study performed
420    on IC neurons (Salles et al., 2020), FMB and Echo calls in the sound database were used to measure
421    responses on artificial neurons. Each audio clip was fed into the network after converting to auditory
422    spectrogram (See, Section 2.3.1), then we obtained activation probabilities for 100-nodes due to the
423    stimulus as $\sigma(v_c)$ in Fig. 4. We averaged the activation probabilities over the same type of calls, and
424    placed the results for 100-nodes onto a 2D scatter plot. Since the IC neurons are categorized into 3-
425    groups such FMB selective, Echo selective, and Non-selective (Salles et al., 2020), we performed k-
426    means clustering (k=3) on the principal axis by the principal component analysis (PCA).

427

428        2.5.2.    Social call representation with artificial neurons

429  We explored bat's call representation with the biomimetic network. In order to perform stochastic
430  analysis, we made 10-copies for each audio clip in the natural sound database (See. Section 2.1.3) by a
431  data augmentation based on temporal shift so that 260 audio clips were ready for the response analysis
432  on artificial network. After converting the audio clips for 8 types of bat call to auditory spectrogram (See,
433  Section 2.3.1), the spectrograms were fed into the network to obtain the network's responses, the $v_c$ in
434  Fig. 4. Then, we estimated the Gaussian distributions for the responses to each call type, and measured
435  a distance between two distributions by using the Jensen-Shannon divergence (JSD) as $JSD(P\|Q) =$
436  $(KLD(P\|M) + KLD(Q\|M))/2$ where $P$ and $Q$ represent two target distributions, $KLD$ is the Kullback-
437  Leibler divergence (KLD), and $M = (P + Q)/2$ (Endres and Schindelin, 2003). Unlike the KLD, the JSD is
438  bounded within the interval [0,1] where 0 means that two distributions are equal. Finally, we quantified
439  a discriminability across the classes by averaging JSDs of all cases choosing 2 of 8. In evaluation, we
440  calculated the averaging JSDs with 10-models for each configuration that were trained on different
441  initial values and summarized the mean and standard deviation of the 10 results. Additionally, we
442  explored the noise effect on the sound representation with simulated audios produced by adding
443  Gaussian random noise to each of the 260 audio clips depending on signal to noise ratio (SNR).

444   To compare with neural data, we performed this analysis between FMB and Echo responses. In the
445  same manner, we calculated JSD based on the networks. We adopted neural data used in the previous
446  study (Salles et al., 2020). Among 575 neurons, we chose 351 neurons which were recorded with same
447  version of stimuli, and constructed 351-dimensional vector to represent response pattern across the
448  neurons by concatenating the number of spikes on each neuron. Once the vector is projected onto 100
449  dimensional space based on PCA, we estimated Gaussian distributions for FMB and Echo responses.
450  Then, we calculated JSD between the distributions.

451

452  **3.  Results**
453  **3.1. Database of natural big-brown bat calls**

454  Acoustic recordings of bat calls emitted while socially housed in the laboratory yielded a data set of
455  natural calls containing a wide range of vocalization types. Fig. 8 shows the time-frequency
456  representation of several types of vocalizations in the database. The bat vocalizations include isolated
457  (non-overlapping) calls representing communication (Fig. 8A-C) or echolocating (Fig. 8D-G) sounds as
458  well as overlapping calls from two distinct bats (Fig. 8G-H). Best Velocity (BV) values reflect the broad
459  range of FM energies in these social communication calls (BV = 18 oct/s, 43 oct/s and 140 oct/s for Fig.
460  8A, B and C respectively). Echolocation calls show even higher FM energies with shorter signals (BV =
461  274 oct/s and 333 oct/s for Fig. 8D and E respectively). In Fig. 8G and H, we note presence of multiple
462  calls though the statistics derived from that segment are largely influenced by the dominant call (BV =
463  381 oct/s and 410 oct/s for Fig. 8G and H respectively). The natural complexity in the animal's auditory
464  environment was maintained in this study and no supervised curation of this data set was performed
465  beyond removal of silence segments (see *Methods*). We also note presence of ambient background in all
466  recordings as a result of the cage environment and recording setup used to collect this data.

467

468  **3.2. Auditory characteristics of biological STRF's**

469   To explore auditory characteristics of big brown bat midbrain, we calculated STRF's from neural
470   recordings of IC neurons. Fig. 9A highlights examples from 6 neurons, revealing a downward sweep
471   selectivity, with excitation and inhibition represented as red and blue areas, respectively. The best
472   frequency (BF) is also shown as red dashed line indicating the maximum peak, positive or negative, of
473   the STRF. We evaluated auditory characteristics across all neural recordings with respect to BV, DSI
474   (direction selectivity), orientation, and inseparability (Fig. 10, yellow histograms). Using a bootstrap
475   procedure, we compared the auditory characteristics of IC neurons to properties of natural calls (Fig. 10,
476   gray background regions for standard deviation). The analysis revealed that the distribution of BVs in IC
477   neurons is statistically equivalent to that of natural calls ($\mu$=-1.63, $\sigma$=17.13, *p-value*=0.9622, Fig. 10A). A
478   match was also observed for direction selectivity $\mu$ =-0.01, $\sigma$=0.03, *p-value*=0.8789, Fig. 10B). This result
479   is consistent with the hypothesis that IC neurons have consistent tuning to the statistics of conspecific
480   vocalizations (Andoni et al., 2007). We noted that the majority of IC neurons (93.6%) favored downward
481   sweeps (Fig. 10B) (Gittelman et al., 2009) while their orientation is centered around 0 deg. Most IC
482   neurons yield higher than rank-one STRF's (average inseparability index 0.49 $\pm$ 0.09).

483    The distribution of frequency tuning (BF) of IC neurons tended to fall between 10 and 30 kHz.
484   Particularly, BF's of 87% of neurons are below 30 kHz (Fig. 11B). In contrast, spectral peaks observed in
485   the vocalization database revealed a higher spectral peak (37.17 $\pm$ 5.62) as shown in Fig. 11A. This
486   profile is likely driven by the strength of the first harmonic component in vocalization which tends to be
487   stronger than other components. As seen from the examples in Fig. 8, most vocalizations contain
488   multiple harmonic peaks with higher energy in the first component resulting in a difference between the
489   BF of IC neurons and spectral peaks of the calls database ($\mu$=-12.58, $\sigma$=1.82, *p-value*=0).

490

### 3.3. Auditory characteristics of artificial STRF's

492   Using natural calls, an artificial network was trained to best represent the statistics of the vocalization.
493   Characteristics of model neurons were analyzed in the same way as biological neurons using spectro-
494   temporal receptive fields. The distribution of model characteristics is shown in Fig. 10, overlaid in blue.
495   Compared to natural calls, model neurons reveal a statistically matching distribution with respect to BV
496   (bootstrap $\mu$=-2.62, $\sigma$=19.85, *p-value*=0.9473) and DSI (bootstrap $\mu$=-0.005, $\sigma$=0.01, *p-value*=0.9382).
497   Model neurons also match the spectral peak of natural calls (bootstrap $\mu$=-0.49, $\sigma$=4.75, *p-value*=0.9592,
498   Fig. 11C). These results are not surprising given that the model was trained to mimic the statistics of
499   these calls. Still, the model was not specifically configured to match specific directionality or velocity
500   patterns but rather represent the time-frequency profile of the calls as a whole.

501    In parallel, the comparison between model and biological neurons reveal remarkable agreement. A
502   bootstrap procedure was performed to compare all auditory characteristics of these STRF's, and results
503   are shown in inset panels in Fig. 10. We note that characteristics of biomimetic neurons match the
504   properties of IC neurons including BVs ($\mu$=2.92, $\sigma$=16.61, *p-value*=0.9300), DSI ($\mu$=0.01, $\sigma$=0.07, *p-
505   value*=0.9358), orientation ($\mu$=1.34, $\sigma$=3.26, *p-value*=0.8370), and inseparability ($\mu$=-0.02, $\sigma$=0.04, *p-
506   value*=0.8079). The BFs of artificial neurons are statistically different from IC neurons (bootstrap $\mu$=12.10,
507   $\sigma$=4.44, *p-value*=0.1731), even though there is substantial overlap at the range of 0-40 kHz. The BFs of
508   artificial neurons are more broadly distributed over the entire frequency range with about 14% of
509   artificial neurons having high BF (above 60 kHz) (Ferragamo et al., 1997).

510

### 3.4. Architecture of the biomimetic network

512 While results reported so far focus on the 'best' biomimetic network, we also investigated how changing
513 the architecture of the model affects the tuning parameters of artificial neurons. We systematically
514 varied the model in terms of structural complexity (the number of stacking blocks), sparsity of the latent
515 space and non-linearity of the activation function. Fig. 12A shows the mean and standard deviation of
516 characteristics of model neurons across 10-network validations for each pair of complexity and sparsity
517 ($\alpha = 0.2$). The mean FM velocities and orientation in the natural calls database are represented by a
518 black line on each panel; while the grey regions represent 95% confidence intervals for each mean. The
519 results show that a very shallow model (mono-stacking) results in a greatly biased negative orientation,
520 as well slower BV estimates. By increasing the model depth, there is an increased match between the
521 model's spectro-temporal configuration (represented by BV and orientation) and that of natural
522 statistics. Furthermore, extremely low or high sparsity values also result in over or under-estimating
523 statistics of natural calls; with 10% sparsity results in a great match with average statistics of the natural
524 calls.

525 Using the triple stacking network with 10% sparsity, we investigated the effect of the model non-
526 linearity on the same auditory characteristics of model neurons (Fig. 12B). Setting the non-linearity
527 parameter to 1.0 results in a fully linear processing which clearly produces in a mismatch between the
528 model and call characteristics. By increasing the degree of non-linearity (decreasing alpha), we note a
529 closer match between the two.

530 It should be noted that across all the different configurations of the model, all architectures were able
531 to converge (i.e. minimize the reconstruction error between the spectrogram of a given call sound and
532 its reconstruction using the model's latent space). Fig. 12C shows the average reconstruction error over
533 the 10 models for each parameter set. While all models successfully converge to reconstruct natural
534 calls and encode statistics of in the database, only a few configurations result in a reasonable match to
535 the spectro-temporal characteristics of model neurons. As a matter of fact, the model was not
536 constrained to match these properties in its latent space; it is merely trained to represent the call
537 spectrograms as faithfully as possible. This requirement has multiple plausible solutions, and only
538 certain configurations result in a close match with velocity and orientation characteristics of natural calls.

539

### 3.5. Natural call representation with the biomimetic network

541 So far, the results suggest that a deep nonlinear architecture with high sparsity to achieve an optimal
542 representation of the statistics of natural bat vocalizations is capable to replicate auditory characteristics
543 of the bat's midbrain. We next examined the implications of this mapping to facilitate discrimination of
544 the large variety in the call repertoire. A study revealed that tuning characteristics of bat IC neurons
545 differentially encode different sound categories in the bat vocalizations, specifically echolocation calls
546 and food-claiming FMB (frequency-modulated bout) social calls (Salles et al., 2020). We examined
547 whether the artificial network, trained simply to emulate natural statistics in the bat repertoire (without
548 knowledge of different sound classes) also yields distinct activations of these different groups. Fig. 13A
549 top replicates the response selectivity of biological IC neurons, showing a scatter plot of average

550  activation probabilities for each neuron in response to FMB calls (x-axis) vs. Echolocation (echo) calls (y-
551  axis), projected on the principal axis by PCA. The figure inset shows the original neural responses before
552  data projection. Fig. 13B-D depict a similar analysis of call selectivity for the mono, double and triple
553  artificial network, respectively. Note that each panel from B to D was produced by one network of 10-
554  models for example. The top panels show a scatter. Across the 3 network configurations, we note that
555  the mono stacking model induces mostly non-selective activation across its neural population (Fig. 13B)
556  while the double stacking model yields biased responses in favor of Echo calls (Fig. 13C). The triple
557  stacking model reveals a more balanced activation from Echolocation and FMB social call types (Fig. 13D)
558  that closely matches biological selectivity.

559   We extend the analysis of call selectivity in the artificial network to other classes of calls in the bat
560  repertoire. We evaluated discriminability across 8 types of calls using the Jensen-Shannon divergences
561  (JSDs) (Endres and Schindelin, 2003). Fig. 14 shows the results for various network depths, linearity and
562  sparsity for the calls in the database (clean) as well as with additional simulated additive noise with
563  decreasing signal-to-noise (SNR) ratios. The triple stacking model (with high sparsity and nonlinear
564  activation) produces the most discriminable responses, as well as more robust discrimination even in
565  presence of noise. Shallower architectures are clearly affected by presence of resulting in reduced
566  discriminability. Linear activations and low sparsity appear to also affect discriminability and robustness
567  to noise albeit not at the same rate. These results suggest that the optimal representation of call
568  statistics likely plays a role in facilitating the identification of different sound classes even in presence of
569  noise. Similarly, a study with guinea pigs has shown the robust discrimination in the responses to
570  communication sounds (Souffi et al., 2020). Such hypothesis aligns with earlier reports (Chechik et al.,
571  2006) but remains to be validated in the IC of the big brown bat. As reference, we computed JSD for the
572  two echo and FMB call classes (Fig. 14) for both artificial and biological neurons. Both measures reveal a
573  close agreement and high discriminability that far surpasses selectivity from shallower architectures.

574

575  ### 4.  Discussion

576  **The biomimetic artificial network provides a nonlinear response model of neural selectivity**

577  To examine the tuning of auditory neurons, each cell can be considered as a system with a mapping
578  function $F$ that represents a relation between stimulus $s$ and neural response $r$ i.e. $r = F(s)$. While
579  characterizing the full system function may be theoretically or experimentally nearly impossible,
580  linearized models using STRF are often used to build a computational response model as $r(t) =$
581  $\int s(t,f) * h(t,f)df$ where $t$ and $f$ is respectively time and frequency index, $s(t,f)$ is spectro-temporal
582  representation for a stimulus, * is a convolution operator, and $h(t,f)$ represents a spectro-temporal
583  receptive fields (STRF's) (Depireux et al., 2001; Elhilali et al., 2013; Machens et al., 2004; Fritz et al.,
584  2003). This model is often applied with reasonable success to predict neural responses to other sound
585  classes including conspecific vocalizations or other natural sounds. Although the linear model is a
586  reasonable approximation for mimicking neural responses in the brain, it is limited in its ability to inform
587  nonlinear transformations that are usually observed in between stimulus and response (Escabi and
588  Schreiner, 2002; Theunissen et al., 2000). One of the main advantages of including nonlinear activations
589  in a feed-forward propagation in the proposed neural network is that it implicitly incorporates the
590  effects of these nonlinear mappings in the propagation of activity throughout the network. Still, we are
591  able to evaluate the linearized portion of the response (via STRF's of artificial neurons) without explicitly

592 incorporating the nonlinear terms in the STRF model itself. This black-box approach to incorporate
593 complexities of neural mapping via deep neural networks opens the possibility to more intricate
594 readouts of the representation of artificial networks. We anticipate that such biomimetic artificial
595 network can be used to build a system mimicking the bat's ability for object shape recognition using its
596 bio-sonar.

597

**598 Midbrain responses are optimized to represent the statistics of natural calls in a bat's soundscape**

599 In this study, we explored the hypothesis that the bat's IC neurons are tuned to represent the FM
600 velocity and spectro-temporal structure of conspecific vocalizations. Evidence in support of this Sender-
601 Receiver matching has been previously reported in the pallid bat (Fuzessery, et al., 2006) and Mexican
602 free-tailed bat (Andoni et al., 2007), as well as other species such as zebra finches (Woolley et al., 2005)
603 (also see (Woolley and Portfors, 2013). Here, we report similar findings in the big brown bat, and
604 establish a close correspondence between acoustic characteristics of natural calls and tuning of spectro-
605 temporal receptive fields of IC neurons of the big brown bat. Going beyond this relationship, an artificial
606 network trained independently on these natural calls reveals tuning properties that not only conform
607 with spectro-temporal features of the calls (which they were trained on), but also unveils IC-like tuning
608 structure and complexity (e.g. separability) that the model was not specifically trained on (Fig. 10). This
609 result hints that the midbrain architecture gives rise to tuning configurations that leverage the spectro-
610 temporal richness of the bat's repertoire to not only represent these features with high fidelity but also
611 enable selective responses to discriminate between classes of natural calls.

612   The artificial network used in the current study shows that the neural encoding of an incoming stimulus
613 gives rise to a response across neural populations that enables it to faithfully reconstruct this stimulus,
614 revealing a high fidelity mapping without loss of information. While not explicitly happening in the brain,
615 this stimulus reconstruction from the internal latent space is the basis for training the artificial network
616 which yields emergent tuning that matches the biology. It is important to note that tuning properties of
617 artificial neurons were derived using moving ripples which invoke the principle of signal decomposition
618 by separating each conspecific call into a sum of ripples with different orientations, rates and phases, in
619 line with the Fourier theory of signal representation. While the network was never trained on these
620 ripples, its response to each ripple spectral motion pattern both in terms of magnitude and phase (both
621 needed for STRF reconstruction) suggest a quantitative correspondence with the downward-sweeping
622 signals that are prominent in the bat repertoire. It is also important to note that not all known coding
623 properties of the bat midbrain are represented in STRFs (Brimijoin and O'Neill, 2010) and that future
624 steps to test time varying response properties (such as an adaptation) would further validate the ability
625 of this network to replicate the biology of the bat IC (Lesica and Grothe, 2008; Rabinowitz et al., 2013;
626 Lohse et al., 2020).

627

**628 A deep architecture with sparsity is best suited to model the statistics of natural calls**

629 Varying the architecture of the network led to different latent spaces to represent the characteristics of
630 the database of natural calls. Specifically, changing the complexity of the network (in terms of depth),
631 sparsity and nonlinearity converged on different solutions for representing conspecific sounds. Under all

632 configurations, the networks were able to reconstruct the input spectrogram with minimal error
633 indicating that its latent space is sufficiently informative about the statistics in the training database (Fig.
634 12C). Nonetheless, only a specific configuration with high sparsity, nonlinearity and sufficient depth is
635 able to replicate biological tuning properties, giving insights into coding principles underlying
636 configuration of IC networks probed in this study. Naturally, while this investigation cannot rule out
637 other configurations that would also reveal a strong match to biology, it can eliminate parameters that
638 converge on solutions that are far from the biology (e.g. shallow networks, linear models). It is worth
639 noting that we were unable to train a quadruple stacking network to represent statistics in the database
640 so we are unable to comment on the extend to which an even deeper network may correlate with
641 biological tuning. The output of a 4th block could be missing spectro-temporal features due to over-
642 compression. This is an issue that could explored using large input patches or modifying the pooling step.

643

644 **Tuning to conspecific natural sounds may underlie selective and robust encoding of auditory objects**

645 We note that directional selectivity to FM sweeps in biological and artificial neurons, results in high
646 discriminability between different classes of calls. Specifically, these results support the notion that by
647 having neural sub-population tuned to different subsets of spectro-temporal statistics, the network is
648 able to encode and differentially respond to different vocalizations and social or echolocation calls. This
649 discriminability is enhanced in the triple sparse and nonlinear network that best matches biological
650 tuning and much reduced in other network configurations despite the fact that these other models were
651 also successfully trained to represent the same natural statistics in the bat call repertoire. This variability
652 may stem from correlated behavior across the neural population which was previously shown to play an
653 important role in enhanced discriminability of vocalizations in the auditory midbrain (Schneider and
654 Woolley, 2010). This encoding selectivity remains fairly stable in presence of stationary ambient noise
655 suggesting that the high dimensional mapping encoding incoming natural calls results in a noise
656 invariant representation that is believed to start emerging at the level of the inferior colliculus and
657 further strengthen in auditory cortex (see, Willmore et al., 2014).

658

659

660

661

662

663

664 **References**

665 Abadi M, Agarwal A, Barham P, Brevdo E, Chen Z, Citro C, et al.. TensorFlow: Large-Scale Machine Learning on
666 Heterogeneous Distributed Systems; 2016. Available from: http://arxiv.org/abs/1603.04467.

667 Aitkin L, Tran L, Syka J. The responses of neurons in subdivisions of the inferior colliculus of cats to tonal,
668 noise and vocal stimuli. Experimental Brain Research. 1994;98(1):53–64.

669 Andoni S, Li N, Pollak GD. Spectrotemporal Receptive Fields in the Inferior Colliculus Revealing Selectivity for

670 Spectral Motion in Conspecific Vocalizations. Journal of Neuroscience. 2007;27(18):4882–4893.

671 Baldi P. Autoencoders, Unsupervised Learning, and Deep Architectures. In: ICML workshop on unsupervised
672 and transfer learning; 2012. p. 37–49.

673 Brimijoin WO, O'Neill WE. On the prediction of sweep rate and directional selectivity for FM sounds from two-
674 tone interactions in the inferior colliculus. Hearing Research. 2005;210(1-2):63–79.

675 Brimijoin WO, O'Neill WE. Patterned tone sequences reveal non-linear interactions in auditory
676 spectrotemporal receptive fields in the inferior colliculus. Hearing Research. 2010;267(1-2):96–110.

677 Casseday JH, Fremouw T, Covey E. The Inferior Colliculus: A Hub for the Central Auditory System. In:
678 Integrative Functions in the Mammalian Auditory Pathway. New York: Springer; 2002. p. 238–318.

679 Chechik G, Anderson MJ, Bar-Yosef O, Young ED, Tishby N, Nelken I. Reduction of Information Redundancy in
680 the Ascending Auditory Pathway. Neuron. 2006;51(3):359–368.

681 Chi T, Ru P, Shamma SA. Multiresolution spectrotemporal analysis of complex sounds. Journal of the
682 Acoustical Society of America. 2005;118(2):887–906.

683 Chi T, Shamma S. NSL Matlab Toolbox; 2005. Available from:
684 http://www.isr.umd.edu/∼speech/nsltools.tar.gz.

685 Choi KP, Xia A. Approximating the number of successes in independent trials : Binomial versus Poisson. The
686 annals of applied probability. 2002;12(4):1139–1148.

687 Corrado GS, Sugrue LP, Seung HS, Newsome WT. Linear-Nonlinear-Poisson Models of Primate Choice
688 Dynamics. Journal of the Experimental Analysis of Behavior. 2005;84(3):581–617.

689 Depireux DA, Simon JZ, Klein DJ, Shamma SA. Spectro-temporal response field characterization with dynamic
690 ripples in ferret primary auditory cortex. Journal of neurophysiology. 2001;85(3):1220–1234.

691 Dietterich T. Overfitting and undercomputing in machine learning. ACM Computing Surveys. 1995;27(3):326–
692 327.

693 Doersch C. Tutorial on Variational Autoencoders; 2016. Available from: http://arxiv.org/abs/1606.05908.

694 Elhilali M, Shamma SA, Simon JZ, Fritz JB. A Linear Systems View to the Concept of STRF. Handbook of Modern
695 Techniques in Auditory Cortex. Nova Science Pub Inc; 2013. p. 33–60.

696    Endres DM, Schindelin JE. A new metric for probability distributions. IEEE Transactions on Information
697    Theory. 2003;49(7):1858–1860.

698    Escabi MA, Schreiner CE. Nonlinear Spectrotemporal Sound Analysis by Neurons in the Auditory Midbrain.
699    Journal of Neuroscience. 2002;22(10):4114–4131.

700    Ferragamo MJ, Haresign T, Simmons JA. Frequency tuning, latencies, and responses to frequency-modulated
701    sweeps in the inferior colliculus of the echolocating bat, Eptesicus fuscus. Journal of Comparative Physiology -
702    A Sensory, Neural, and Behavioral Physiology. 1997;182(1):65–79.

703    Fritz J, Shamma S, Elhilali M, Klein D. Rapid task-related plasticity of spectrotemporal receptive fields in
704    primary auditory cortex. Nature Neuroscience. 2003;6(11):1216–1223.

705    Fuzessery ZM, Richardson MD, Coburn MS. Neural Mechanisms Underlying Selectivity for the Rate and
706    Direction of Frequency-Modulated Sweeps in the Inferior Colliculus of the Pallid Bat. Journal of
707    Neurophysiology. 2006;96(3):1320–1336.

708    Gittelman JX, Li N, Pollak GD. Mechanisms underlying directional selectivity for frequency-modulated sweeps
709    in the inferior colliculus revealed by in vivo whole-cell recordings. Journal of Neuroscience.
710    2009;29(41):13030–13041.

711    Gordon M, O'Neill WE. Temporal processing across frequency channels by FM selective auditory neurons can
712    account for FM rate selectivity. Hearing Research. 1998;122(1-2):97–108.

713    Kay SM. Fundamentals of statistical signal processing. Englewood Cliffs, N.J.: Prentice-Hall PTR; 1993.

714    Kowalski N, Depireux DA, Shamma SA. Analysis of dynamic spectra in ferret primary auditory cortex. I.
715    Characteristics of single-unit responses to moving ripple spectra. Journal of Neurophysiology.
716    1996;76(5):3503–3523.

717    Lesica NA, Grothe B. Efficient temporal processing of naturalistic sounds. PLoS ONE. 2008;3(2).

718    Lohse M, Bajo VM, King AJ, Willmore BDB. Neural circuits underlying auditory contrast gain control and their
719    perceptual implications. Nature Communications. 2020;11(1):1–13.

720    Luo J, Moss CF. Echolocating bats rely on audiovocal feedback adapt sonar signal design. Proceedings of the
721    National Academy of Sciences of the United States of America. 2017;114(41):10978–10983.

722    Park J, Kim W, Han DK, Ko H. Voice activity detection in noisy environments based on double-combined
723    Fourier transform and line fitting. The Scientific World Journal. 2014;2014.

724    Poon PWF, Yu PP. Spectro-temporal receptive fields of midbrain auditory neurons in the rat obtained with
725    frequency modulated stimulation. Neuroscience Letters. 2000;289(1):9–12.

726    Quiroga RQ, Nadasdy Z, Ben-Shaul Y. Unsupervised Spike Detection and Sorting with Wavelets and
727    Superparamagnetic Clustering. Neural Computation. 2004;16(8):1661–1687.

728    Rabinowitz NC, Willmore BDB, King AJ, Schnupp JWH. Constructing Noise-Invariant Representations of Sound
729    in the Auditory Pathway. PLoS Biology. 2013;11(11).

730    Radford A, Metz L, Chintala S. Unsupervised Representation Learning with Deep Convolutional Generative
731    Adversarial Networks; 2015. Available from: http://arxiv.org/abs/1511.06434.

732    Maas AL, Hannun AY, Ng AY. Rectifier nonlinearities improve neural network acoustic models. In: in ICML
733    Workshop on Deep Learning for Audio, Speech and Language Processing. vol. 28; 2013. p. 1–6.

734    Machens CK, Wehr MS, Zador AM. Linearity of cortical receptive fields measured with natural sounds. The
735    Journal of neuroscience. 2004;24(5):1089–1100.

736    Macıas S, Luo J, Moss CF. Natural echolocation sequences evoke echo-delay selectivity in the auditory
737    midbrain of the FM bat, eptesicus fuscus. Journal of Neurophysiology. 2018;120(3):1323–1339.

738    Salles A, Park S, Sundar H, Macias S, Elhilali M, Moss CF. Neural Response Selectivity to Natural Sounds in the
739    Bat Midbrain. Neuroscience. 2020;434:200–211.

740    Scherer D, Muller A, Behnke S. Evaluation of Pooling Operations in Convolutional Architectures for Object
741    Recognition. In: International Conference on Artificial Neural Networks (ICANN); 2010. p. 92–101.

742    Schneider DM, Woolley SMN. Discrimination of Communication Vocalizations by Single Neurons and Groups
743    of Neurons in the Auditory Midbrain. Journal of Neurophysiology. 2010;103(6):3248–3265.

744    Schwartz O, Pillow JW, Rust NC, Simoncelli EP. Spike-triggered neural characterization. Journal of Vision.
745    2006;6(4):484–507.

746    Shamma SA. Speech processing in the auditory system I: The representation of speech sounds in the
747    responses of the auditory nerve. The Journal of the Acoustical Society of America. 1985a;78(5):1612–1621.

748    Shamma SA. Speech processing in the auditory system II: Lateral inhibition and the central processing of
749    speech evoked activity in the auditory nerve. The Journal of the Acoustical Society of America.
750    1985b;78(5):1622–1632.

751    Shelhamer E, Long J, Darrell T. Fully Convolutional Networks for Semantic Segmentation. IEEE Transactions
752    on Pattern Analysis and Machine Intelligence. 2017;39(4):640–651.

753    Singh N, Theunissen F. Modulation spectra of natural sounds and ethological theories of auditory processing.
754    Journal of the Acoustical Society of America. 2003;106:3394–3411.

755    Smith EC, Lewicki MS. Efficient auditory coding. Nature. 2006;439:978–982.

756    Souffi S, Lorenzi C, Varnet L, Huetz C, Edeline JM. Noise-Sensitive but More Precise Subcortical
757    Representations Coexist with Robust Cortical Encoding of Natural Vocalizations. Journal of Neuroscience.
758    2020;40(27):5228–5246.

759    Strang G. Introduction to linear algebra. 4th ed. Wellesley-Cambridge Press; 2009.

760    Suta D, Kvasnak E, Popelar J, Syka J. Representation of Species-Specific Vocalizations in the Inferior Colliculus
761    of the Guinea Pig. Journal of Neurophysiology. 2003;90(6):3794–3808.

762    Szegedy C, Liu W, Jia Y, Sermanet P, Reed S, Anguelov D, et al. Going deeper with convolutions. In: Proceedings
763    of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition. vol. 07-12-June; 2015.
764    p. 1–9.

765    Theunissen FE, Sen K, Doupe AJ. Spectral-temporal receptive fields of nonlinear auditory neurons obtained
766    using natural sounds. Journal of Neuroscience. 2000;20(6):2315–2331.

767 Wang K, Shamma SA. Self-normalization and noise-robustness in early auditory representations. IEEE
768 Transactions on Speech and Audio Process. 1994;2:421–435.

769 Williams AJ, Fuzessery ZM. Facilitatory mechanisms shape selectivity for the rate and direction of FM sweeps
770 in the inferior colliculus of the pallid bat. Journal of Neurophysiology. 2010;104(3):1456–1471.

771 Willmore BDB, Cooke JE, King AJ. Hearing in noisy environments: noise invariance and contrast gain control.
772 The Journal of Physiology. 2014;592(16):3371–3381.

773 Woolley SMN, Fremouw TE, Hsu A, Theunissen FE. Tuning for spectro-temporal modulations as a mechanism
774 for auditory discrimination of natural sounds. Nature Neurosci. 2005;8(10):1371–1379.

775 Woolley SMN, Portfors CV. Conserved mechanisms of vocalization coding in mammalian and songbird
776 auditory midbrain. Hearing Research. 2013;305(1):45–56.

777 Wright GS, Chiu C, Xian W, Wilkinson GS, Moss CF. Social calls of flying big brown bats (Eptesicus fuscus).
778 Frontiers in Physiology. 2013;4 AUG(August):1–9.

779 Yang X, Wang K, Shamma SA. Auditory representations of acoustic signals. IEEE Transactions on Information
780 Theory. 1992;38(2):824–839.

781

782

**Figures and Tables**

Fig. 1 Overview of study foci. A database of natural calls from a colony of big brown bats is collected and analyzed for its auditory characteristics. Shown in the figure is a distribution of FM velocities. Right: Tuning characteristics of biological neurons from the big brown bat inferior colliculus are derived using Spectro-Temporal Receptive Field (STRF) method, and properties of biological neurons are derived (e.g. best velocity, BV). Shown in the figure is a brain slice identifying the location of the IC in the big brown bat (from Salles et al., 2020). Left: Computational models with various configurations are examined and emergent tuning properties of artificial networks are derived to compare against statistics of natural calls as well as biological neurons.

Fig. 2 Example spectrograms of 8 bat calls in social call database. **A**, Echolocation (Echo). **B**, Frequency modulated bout (FMB). **C**, Upward frequency modulated (UFM). **D**, Long frequency modulated (LFM). **E**, Short frequency modulated (SFM). **F**, Chevron shaped (CS). **G**, Hook. **H**, Long-wave.}

Fig. 3 Ripple transfer function extraction. **A**, A subset of ripple stimuli. **B**, Input patch sequence configuration from a ripple stimulus to a code vector for characterizing the network. **C**, Examples of responses on each node (each element of the code vector), and definition of magnitude $m$ and phase $\phi$ in a ripple response. **D**, Magnitude and phase plots on one of nodes.}

Fig. 4 Convolutional layered autoencoder structure for biomimetic network. The flow denoted in black shows a double stacking structure as a standard example. Based on this structure, a deeper structure can be constructed by stacking more modules, on the other hand, a shallow structure is created by removing a module on the top of the standard example.

Fig. 5 Operations using multi-scale filters. **A**, convolution using multi-scale filters. **B**, transposed convolution using multi-scale filters.

Fig. 6 STRF calculation. **A**, expanded magnitude and phase matrices which are matching to Fig. 2D. **B**, 2-Dimensional (above) and 3-Dimensional (bottom) representation of the STRF that is obtained by performing 64 by 64 interpolation and Gaussian smoothing sequentially. In 2D representation, red area represents excitation regions while blue represents inhibition regions.

Fig. 7 Descriptions for best velocity and orientation. **A**, magnitude plot. **B**, the result of Gaussian surface fitting. The red ellipse represents Gaussian mean vector $\mu$ (center) and covariance matrix $\Sigma$ (rotation), the best velocity is defined by a slope of Gaussian mean vector and the orientation error is defined by an angle difference between mean vector and covariance rotation.

23

820  Fig. 8 Example spectrograms of several types of calls monophonic cases for social communication (A-C),
821  echolocation (D-F), and polyphonic cases (G-H). Note the differences in frequency content, duration, and
822  sweep velocity. Note that peak frequency is represented onto each panel as the black dashed line.

823

824  Fig. 9 Examples for biological STRF and biomimetic STRF. **A**, biological STRFs obtained from bat's IC
825  neuron. **B**, biomimetic STRFs obtained from a triple stacking network with 10% sparsity. Note that red
826  and blue area show excitation and inhibition regions, respectively.

827

828  Fig. 10 Histogram of biological and biomimetic STRFs according to auditory characteristics. **A**, Best
829  velocity. **B**, Direction selectivity index. **C**, Orientation (The zero-mean is marked as the red line). **D**,
830  Inseparability.

831

832  Fig. 11 Analysis of best frequency in dataset, IC neurons, and artificial neurons. **A**, histogram of peak
833  frequencies in natural calls, the background grey line represents averaging spectrum envelop of natural
834  calls. **B**, BFs on IC neurons. **C**, BFs on artificial neurons.

835

836  Fig. 12 Best velocity and Orientation of biomimetic STRF's depending on network's configuration. **A**, for
837  the number of stacking modules and sparsity (0.2 lReLu). **B**, for nonlinearity (Triple stacking model with
838  10% sparsity). **C**, average reconstruction error over the 10 networks for each parameter set, blue and
839  red dashed lines are mean of the errors for all configurations at the beginning of training and the end of
840  the training, respectively.

841

842  Fig. 13 Selectivity to FMB vs. Echolocation call for; **A**, IC neurons (Salles et al., 2020), spike frequency was
843  calculated by dividing the number of spikes by total number of spikes starting 5 ms after stimulus onset.
844  The horizontal axis on the bottom was circularly shifted with zero-centered non-selective neurons. **B**,
845  mono-stacking model. **C**, double-stacking model. **D**, triple-stacking model.}

846

847  Fig. 14 Natural sound representation by biomimetic network in different SNR conditions
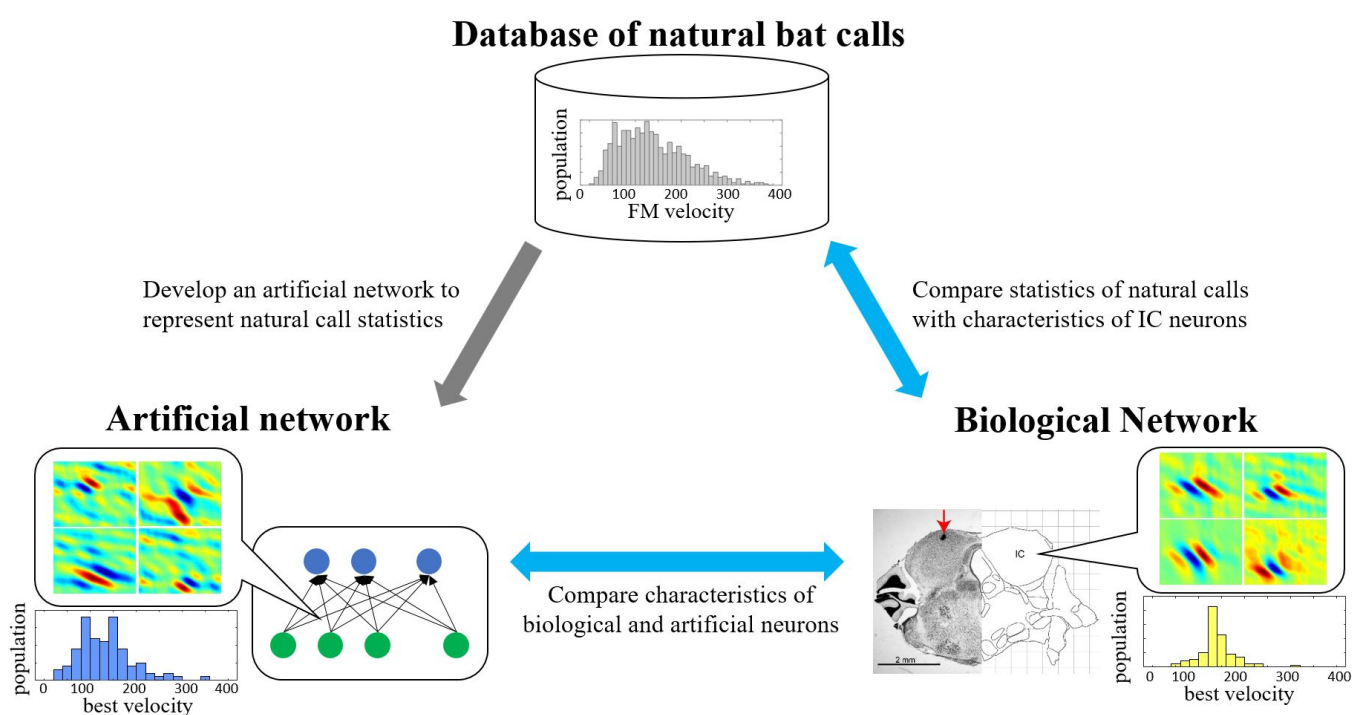
848
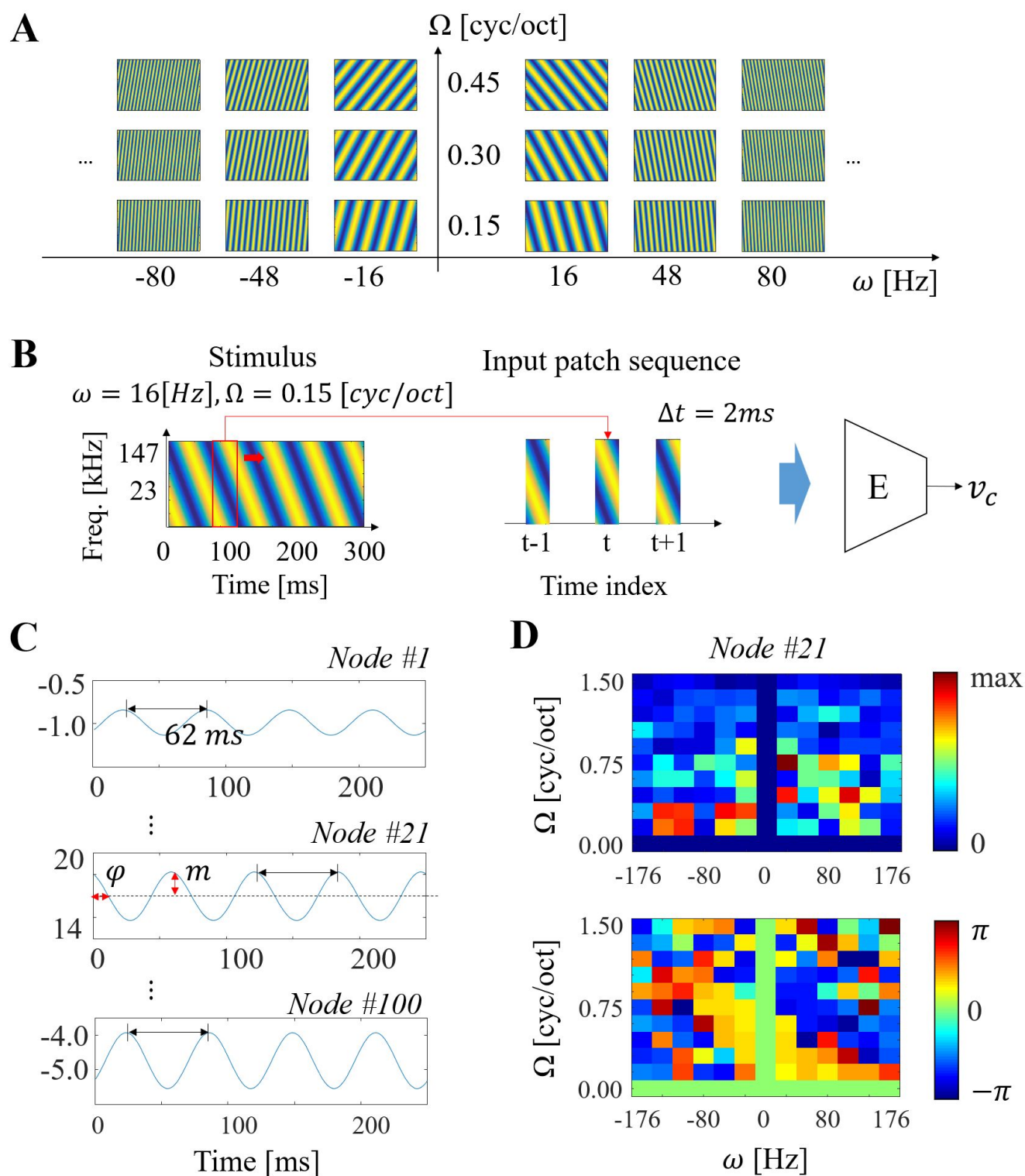
849  Table 1. Description of network parameters, midlevel feature maps, and input.

850

851

852

24

**Database of natural bat calls**

Develop an artificial network to
represent natural call statistics

Compare statistics of natural calls
with characteristics of IC neurons

**Artificial network**

**Biological Network**

Compare characteristics of
biological and artificial neurons

**A**

$$I_o^l$$

Convolution

$$I_i^l = I^{l-1}$$

$$I_o^l$$

Concatenating

$$f_1^l * \quad f_2^l * \quad \cdots \quad * f_k^l$$

$$I_i^l = I^{l-1}$$

**B**

$$\tilde{I}^l$$

$$\tilde{I}^l$$

Transposed_convolution

$$\tilde{I}^{l-1}$$

Split feature maps

$$\dot{f}_1^l * \quad \dot{f}_2^l * \quad \cdots \quad * \dot{f}_k^l$$

Averaging

$$\tilde{I}^{l-1}$$

**A**

$M_e(\omega, \Omega)$



$\Phi_e(\omega, \Omega)$



**B**

$\mathcal{F}_{t,f}^{-1}[T(\omega, \Omega)]$

**A**



**B**

**A** *FM velocity of natural calls*

Density

0.3

0.2

0.1

0 100 200 300 400

BV [oct/s]

Inset: $BV_{net} - BV_{ic}$, y-axis 50, 100, 150, x-axis -60 to 60

**B** *DSI of natural calls*

Density

0.2

0.1

-1.0 -0.5 0 0.5 1.0

Direction Selectivity Index

Inset: $DSI_{net} - DSI_{ic}$, y-axis 60, 120, x-axis -0.3 to 0.3

**C**

Density

0.3

0.2

0.1

-30 -20 -10 0 10 20 30

Orientation [deg.]

Inset: $Ori_{net} - Ori_{ic}$, y-axis 80, 160, x-axis -15 to 15

**D**

Density

0.3

0.2

0.1

0 0.2 0.4 0.6 0.8 1.0

Inseparability

Inset: $Ins_{net} - Ins_{ic}$, y-axis 50, 100, 150, x-axis -0.15 to 0.15

Biomimetic STRF    Biological STRF

**A**



**B**



**C**

**A**



# of spk. to FMB

# of spk. to Echo

(34,14)

Spike rate to Echo

1.0

0.5

(0.71, 0.29)

0    0.5    1.0

Spike rate to FMB

N=351

Population

20

10

0

-40    0    40

Principal axis
(circularly shifted)

**B**

Spike rate to Echo

1.0

0.5

0    0.5    1.0

Spike rate to FMB

1st Principle axis

N=100

Population

12

6

0

-0.3    0    0.3

1st Principal axis

**C**

Spike rate to Echo

1.0

0.5

0    0.5    1.0

Spike rate to FMB

1st Principle axis

N=100

Population

12

6

0

-0.3    0    0.3

1st Principal axis

**D**

Echo selective

Non-selective

Spike rate to Echo

1.0

0.5

FMB selective

0    0.5    1.0

Spike rate to FMB

1st Principle axis

N=100

Population

12

6

0

-0.3    0    0.3

1st Principal axis