

1 **Title: Inferring the Neural Basis of Binaural Detection Using Deep Learning**

2

3 **Authors:** Samuel S. Smith*¹, Joseph Sollini¹, Michael A. Akeroyd¹

4

5 1- Hearing Sciences, Division of Clinical Neuroscience, School of Medicine,
6 University of Nottingham, Nottingham, NG7 2RD, UK

7

8 ***Correspondence:** samuel.smith@nottingham.ac.uk

9

SUMMARY

10

11 **Neural activity from animals is often used as a proxy for the human brain.**
12 **However, due to distinct environmental pressures, the relevance of perceptual systems**
13 **described in animal models can be unclear. This problem is accentuated when animal**
14 **physiology and human behaviour are not in complete agreement, as is the case for**
15 **binaural hearing-in-noise. As a means to bridge this gap we reverse-engineered**
16 **artificial neural networks from binaural psychophysics. By comparing *in silico***
17 **“physiology” in neural networks with *in vivo* animal data, we were able to make**
18 **inferences as to the basis of binaural perception in humans. We observed the**
19 **emergence of highly specialized solutions to account for low frequency sound**
20 **detection. Artificial neurons developed a sensitivity to temporal delays that increased**
21 **hierarchically and were widely distributed in preference. Network dynamics were**
22 **consistent with a cross-correlator, comparable to the type reported in animal**
23 **physiology. Our results attest to the likely prominence of this neural mechanism in**
24 **human biology. Moreover, this is a primary demonstration that deep learning can infer**
25 **tangible neural mechanisms underlying auditory perception.**

26

27

28

INTRODUCTION

29

30 Deep neural networks (DNNs) have been used to solve problems in many fields of
31 research and are increasingly proving their worth in the field of neuroscience¹. Recent DNN
32 studies have effectively addressed questions of *why* the auditory system is organized the way
33 it is (typically in the context of task optimality)²⁻⁵. However, when human auditory
34 neurophysiology is itself ambiguous, or unknown, we must first question *what* it is, i.e. discover
35 its underlying dynamics. With a few design changes, could DNNs be better leveraged to learn
36 about the underlying human neurophysiology driving audition? We tested this idea by training
37 a DNN configured specifically to mimic human auditory behaviour and investigated what this
38 might reveal about the underlying neural mechanism(s). One potential stumbling block in
39 answering this question is the black-box nature of DNNs. However, new network architectures
40 that put mechanistic interpretability at their forefront (as have shown promise in the field of
41 physics^{6,7}) could help overcome this limitation.

42

43

44

45

46

47

An ideal context in which to examine these inferential properties of DNNs is one where
it is unclear whether non-human neurophysiology satisfactorily explain human audition.
Binaural detection^{8,9}, where interaural differences enhance the detectability of one sound (a
signal) amongst another (e.g. a background noise) by up to 15 dB, represents one such
instance. There is ongoing debate as to which of a number of theoretical frameworks best
relate to human binaural detection¹⁰⁻¹⁴. For example, animal neural data lend support to a

48 theory of binaural cross-correlation^{15–17}. Whereas, human behaviour appears to be equally
49 well, if not better, described by a noise equalization and subtraction scheme^{12,18,19}. These
50 discrepancies have not been resolved with human imaging data^{20–23}, for which resolution and
51 response variability are key limitations²⁴. Further, binaural detection is a highly specialised
52 auditory function for which deficits have real-world consequences^{25,26}. DNNs may offer the
53 opportunity to bridge this gap between animal and human data, and as yet, the inner workings
54 of DNNs constructed to handle binaural audio have scarcely been considered^{27–29}.

55 Here, we reverse-engineered DNNs that accounted for human-like behaviour in a
56 binaural detection task. To best facilitate a mechanistic understanding, the DNN architecture
57 was configured to decode inputs into low-dimensional latent representations from which
58 decisions were based^{6,7}. Not only did the DNN that best mimicked human behaviour learn to
59 utilize binaural information, but it did so in a way strikingly similar to that described in animal
60 physiology. The work attests to the prominence of binaural cross-correlation as a solution to
61 signal detection at low frequencies, and its likely incidence in humans.

62
63

64 RESULTS

65

66 To augment the availability of data, we trained deep neural networks (DNNs) on data
67 from a simulated binaural detection task. These data were generated by a set of equations
68 recognized as effective in predicting human binaural psychophysics^{12,18,19,30–32}. DNNs were
69 trained to mimic detection of a 500 Hz pure tone amongst a broadband noise, each with
70 interaural time differences (ITDs) that varied trial-to-trial. The range of ITDs was restricted to
71 fall within the human physiological range, i.e. as though they came from randomly chosen
72 azimuthal locations in the real-world (**Fig. 1a**). We found that the DNN configuration that
73 optimally predicted unseen binaural detection data did so with a root mean square error
74 (RMSE) of 2.5% (**Extended Data Fig. 1a**). The dynamics of this optimally performing DNN
75 are the focus of this article (summary analytics across other DNN configurations are shown in
76 **Extended Data Fig. 1**).

77

78 Deep neural network accounts for binaural detection psychophysics.

79 As expected, we found that the DNN's detection thresholds (i.e. 69% correct
80 performance) decreased as ITD difference between tone and noise increased (**Fig. 1c,d**). For
81 example, in diotic noise (noise ITD = 0) where the tone came from the left, the detection
82 thresholds were significantly enhanced by 9 dB (from a maximum of 30.8 dB to 21.8 dB, two-
83 sided unpaired t-test, $p \ll 0.0001$). To allow comparative assessment of the DNN and
84 previously published data we tested the network on stimulus configurations typically employed
85 to study binaural detection. These include tones and noise configurations where they are
86 either in-phase or completely out-of-phase across the ears. In the literature these stimuli are

87 denoted as NoSo, NoS π , N π S π , N π So, where N refers to the noise component, S the pure
88 tone signal, with the successive subscripts denoting interaural phase difference (IPD) in
89 radians (note that for a pure tone IPD is linearly proportional to ITD). These stimuli use ITDs
90 that fall outside of the physiological range. For example, a 500 Hz pure tone with an IPD of π
91 corresponds to an ITD of 1000 μ s, larger than that produced by the head size in the simulated
92 training data (maximum of 655 μ s, calculated with Woodworth's formula³³). This meant the
93 DNN had no prior exposure to this size of ITD and so it was unclear how it would function over
94 this range. We found that when the noise signal had zero IPD, the mean detection thresholds
95 predicted by the DNN for corresponding homophasic (NoSo) and antiphasic (NoS π) tones
96 were 30.9 dB and 20.8 dB respectively. The gain in detection, commonly called the binaural
97 masking level difference (BMLD), was 10.1 dB ($p \ll 0.0001$, two-sided unpaired t-test).
98 Comparatively, when instead the noise signal was interaurally out of phase, the mean
99 detection thresholds predicted for the corresponding homophasic (N π S π) and antiphasic
100 (N π So) stimuli were 26.6 dB and 17.1 dB respectively. Their BMLD was 9.5 dB ($p \ll 0.0001$,
101 two-sided unpaired t-test). These BMLDs are similar to those measured in people³⁴ and with
102 estimates from psychophysical equations (10.7 dB and 10.3 dB respectively; **Fig. 1e**).

103

104 **Time delay tuning emerges early in the network.**

105 The DNN was able to account for key aspects of human binaural detection behaviour.
106 Given this, we next sought to understand the means by which the DNN derived this behaviour,
107 i.e. the mechanism(s). A common property of animal binaural systems is ITD tuning (**Fig. 2a**).
108 We found that ITD tuning emerged hierarchically within the lower layers of the network. To
109 demonstrate this we characterized "noise delay" functions in DNN nodes, i.e. their response
110 to noises presented with varying ITDs as typically measured in physiology studies³⁵. For nodes
111 in the DNN's first layer, we observed significant ITD tuning in 63 out of 100 nodes (**Fig. 2b**).
112 The noise delay responses exhibited in these nodes were well described by a Gabor
113 function¹⁵, the combination of a cosine windowed by a Gaussian (overlaid in **Fig. 2d**). By the
114 DNN's second layer, significant ITD tuning had emerged in all 100 nodes (**Fig. 2e**). Estimates
115 of each nodes' best ITD (bITD) were inferred from the Gabor fits (in order to account for nodes
116 that were oscillatory in their noise delay responses, a form of phase-ambiguity noted in
117 physiology³⁶, bITD was attributed to the most central tuning peak). In both the first and second
118 layers of the DNN, we observed a wide distribution of bITDs, both within the simulated head-
119 range, and beyond it.

120

121 **Network dynamics match those of a cross-correlation mechanism.**

122 We went on to measure responses to stimuli commonly presented in physiology work
123 to specifically probe binaural detection (i.e. NoSo, NoS π , N π S π , N π So). We found that in the
124 DNN's second layer, node responses varied by bITD and the interaural phase of the noise
125 presented (two-way ANOVA, $F[98,202]=12.5$ for main effect of bITD, $F[1,202]=31.2$ for main

126 effect of noise phase, $F[98,202]=2.9$ for their interaction, $p \ll 0.0001$ for all, **Fig. 3a**). When a
127 signal was presented amongst an in-phase noise (No), responses were largest for nodes with
128 bITDs near $0 \mu\text{s}$ and decreased as bITDs were increasingly non-zero. Conversely, amongst
129 an out-of-phase noise ($N\pi$), responses were lowest for nodes with bITDs near $0 \mu\text{s}$ and
130 increased as bITDs deviated away from this. The effects of tone phase on node dynamics
131 were more subtle, although these dynamics were also in accordance with a nodes' tuning
132 properties (**Fig. 3b,c**). Nodes tuned to smaller ITDs responded most to in-phase tones (So)
133 and least to out-of-phase tones ($S\pi$), and vice-versa for nodes tuned to larger ITDs.

134 These responses are commensurate with a binaural cross-correlation mechanism. The
135 concept of binaural cross-correlation is predicated on the existence of coincidence detectors
136 that encode temporally offset signals (similar to the dynamics already described in layers 1
137 and 2)³⁷. Computationally, a binaural cross-correlation can be calculated by summing the
138 point-by-point product of two temporally offset signals. Comparative outputs from a simple
139 binaural cross-correlation algorithm (namely for signals in noise passed through narrow-band
140 filters centered at 500 Hz) are shown in **Figure 3d-f**, and were found to resemble responses
141 across the DNN's layer 2 nodes (Pearson's $r=0.36$, $p \ll 0.0001$). A number of physiology
142 studies have reported neural responses consistent with a cross-correlation mechanism^{15,16,38,39}
143 (**Fig. 3g-i**).

144

145 **Early network ablation is detrimental to binaural detection.**

146 The functional importance of early layers in the DNN (i.e. the "decoder" portion of the
147 network) was further interrogated by inflicting targeted damage and observing knock-on
148 effects to BMLDs. We set to zero the weights of a fixed proportion of nodes in a specific layer
149 of the network, i.e. ablating them. Our ablation range varied from 0% (none) to 50% (lots). We
150 then observed the corresponding NoSo-No $S\pi$ BMLDs (**Extended Data Fig. 2**). We found that
151 ablation to nodes in the DNN's first layer were most detrimental to BMLDs to such an extent
152 that with as little as 5% ablation the DNN failed to predict a significant BMLD (criterion set at
153 $p < 0.05$ with Bonferonni correction, Student's two-tailed t-test). When ablating the second
154 layer, the DNN initially failed to predict a BMLD when 35% of nodes were ablated. Layers later
155 in the network were more robust to the effects of ablation. BMLDs withstood ablations up to
156 and including 50% (maximum tested) of the nodes in fourth layer, and significant BMLDs were
157 not exhibited when 40% or more of the nodes in the fifth layer were ablated.

158

159 **Low-dimensional representations in DNN imitate neural signature of population-** 160 **level cortical activity.**

161 Responses in the DNN's second layer combine together to form the low-dimensional
162 representations encoded in the DNN's central nodes (i.e. layer 3). This bottleneck architecture
163 has proven successful in extracting key conceptual variables in other fields⁷, and our ambition
164 was for something analogous for binaural detection. In the DNN, six central nodes were

165 deemed operational (**Extended Data Fig. 1b**), and we found noteworthy similarities between
166 these nodes and population level cortical responses in the guinea pig¹⁶ (comparable
167 observations have also been reported in other animals^{15,39} and in the guinea pig inferior
168 colliculus³⁸). The central nodes could take the value of any real number, positive or negative
169 (a necessary limitation imposed by the network architecture). Although not essential for our
170 main conclusions, node responses are presented polarity-corrected to best correspond with
171 noise ITD tuning in cortex (**Fig. 4a,c**; guinea pig auditory cortex).

172 An interesting feature of binaural processing can be seen when comparing detection
173 behaviour and physiological data. Similar improvements in behavioural performance can be
174 attributed to completely different alterations in network dynamics. For example, BMLDs across
175 NoSo-No π conditions are similar to those found across N π S π -N π So conditions (**Fig. 1e**).
176 However, guinea pig neural data¹⁶ suggests different neural dynamics underlie these similar
177 BMLDs¹⁶. In cortical recordings, population spike counts dropped amongst an No signal, as a
178 pure tone went from So to S π (left panel in **Fig. 4b**). Conversely, amongst an N π signal, as a
179 pure tone transitioned from S π to So, population spike counts increased (right panel in **Fig.**
180 **4b**). These opposing dynamics therefore represent a unique signature of processing

181 The majority of layer 3 nodes displayed the same opposing dynamics in response to
182 homophasic/antiphase stimuli as observed in the guinea pig auditory cortex. Threshold
183 responses in four latent nodes (n_1 , n_3 , n_4 and n_5 , in left panels of **Fig. 4d**) were found to be
184 significantly lower in response to NoS π relative to NoSo ($p < 0.05$, two-sided unpaired t-test).
185 Conversely, threshold responses in the same nodes to N π So at a threshold level were
186 significantly higher in comparison to N π S π ($p < 0.05$, two-sided unpaired t-test, right panels in
187 **Fig. 4d**). One latent node was qualitatively different to the others (node 6 of **Fig. 4d**), and
188 seemed to encode offsets related to the interaural phase of noise (No v. N π , $p \ll 0.0001$, two-
189 sided unpaired t-test).

190
191

192 DISCUSSION

193

194 We set out to discover the efficacy of DNNs as a means of exploring the underlying
195 mechanisms involved in hearing, specifically binaural detection. To do this, we trained DNNs
196 to exhibit binaural detection resembling human behaviour and then examined their internal
197 dynamics as model organisms⁴⁰. The application of DNNs in this way is a promising method
198 in systems neuroscience¹. However, the capacity for DNNs to offer mechanistic understanding
199 beyond broader analogies with auditory processing²⁻⁵ has yet to be established. This work
200 demonstrated not only a number of key similarities with non-human binaural systems but,
201 critically, the method implies that the human auditory system may use alike mechanisms.

202 Perhaps easy to overlook, the DNN was able to successfully utilize binaural
203 discrepancies in auditory stimuli, as opposed to seeking an alternative strategy⁴¹ and/or failing

204 to exhibit binaural detection behaviour. ITD tuning, a well-known characteristic of binaural
205 neurons and normally described in the context of sound localization³⁵, emerged early in the
206 DNN. Although this is a notable finding when considering the potential of DNN models, ITD
207 tuning is axiomatic in most explanations of binaural detection¹⁰. This ambiguity was better
208 clarified by second layer nodes whose dynamics, in response to tones presented in broadband
209 noise, resembled a binaural cross-correlation mechanism¹³. This mechanism was not hard-
210 coded into the network, but inferred. Latent nodes, central within the DNN's architecture, were
211 also compatible with the downstream dynamics of a cross-correlation mechanism and
212 resembled guinea pig cortical neural recordings¹⁶. These results help reinforce and unite
213 findings supportive of binaural cross-correlation as the mechanism underlying binaural
214 detection in people, as opposed to other explanations^{12,18,19}.

215 We also experimented with a technique analogous to neural ablation⁴², observing
216 knock-on effects to detection performance, finding that manipulations early in the system were
217 most detrimental. This result is consistent with atypically small BMLDs associated with
218 peripheral tumors in the human auditory system, but not central lesions⁴³. The potential
219 insights attainable with other experimental DNN techniques is an exciting prospect (e.g.
220 techniques analogous to neural stimulation⁶ in tandem with optogenetics data). Yet,
221 comparison between DNNs and neural biology come accompanied by an asterisk. We make
222 no claims of creating a general-purpose implementation of the human binaural system. The
223 network was not constructed with the goal of accurately mimicking neuronal biophysics or
224 hierarchical complexity, but instead a trade-off was made to favor mechanistic interpretation
225 and optimization performance. The inclusion of additional structural priors (e.g. hemisphericity)
226 and biologically inspired processes (e.g. spiking neural networks^{2,44}) could have merit, but any
227 impact on interpretability should be carefully weighed.

228

229 **Conclusion**

230 In conclusion, our results indicate that an artificial neural network seeks to implement
231 a specialized binaural mechanism to explain human binaural detection. This mechanism,
232 (namely, temporal delay tuning followed by a cross-correlator) corroborates observations
233 made in animal physiology. The work demonstrates the potential for deep learning, in unison
234 with experimental data, to clarify human auditory perception.

235

METHODS

236

237 **Training stimuli.** Stimuli parameters were selected to maximize comparative opportunities
238 with published experimental data. Pure tones were produced at a frequency of 500 Hz, and
239 presented at levels between 0 and 50 dB. Pure tones were 20 ms long (10 periods) and
240 produced with a sample rate of 20 kHz. Pure tones were masked by randomly distributed
241 broadband noise (50-5000 Hz, limited by 6th order Butterworth bandpass filter) with an overall
242 level of 60 dB. The tone and noise were gated simultaneously. Horizontal perception of space
243 at low frequencies is largely based upon ITDs⁴⁵. Given this, tones and noises were simulated
244 with ITDs mapped from two independent angles in the azimuth between -90° (far left) and 90°
245 (far right). ITDs were derived from Woodworth's equation³³, assuming a head radius of 0.0875
246 m.

247

248 **Binaural detection rates and thresholds.** The theory of equalization and cancellation¹²
249 has wide human psychophysical support, successful in predicting BMLD data^{12,46} and binaural
250 pitch phenomena^{12,30,31}, underpinning other models of binaural hearing¹⁹, and proven
251 psychophysically favourable relative to other prominent theories¹⁸. Detection thresholds were
252 calculated from phenomenological equations derived from this theory^{12,32}:

$$253 \quad D(\tau_S, \tau_N) = 31 - 10 \log_{10} \max \left\{ \frac{k - \cos(\omega_0 \tau_S - \omega_0 \tau_N)}{k - \gamma(\tau_N - \tau_0)}, 1 \right\} \text{ dB} \quad (1)$$

254 where τ_S and τ_N are the interaural time lags of the signal and noise, ω_0 is the angular
255 frequency of the pure tone signal, $k = (1 + \sigma_\epsilon^2)e^{\omega_0^2 \sigma_\delta^2}$ where σ_ϵ^2 and σ_δ^2 are jitter (internal
256 noise) parameters, γ is the normalized envelope of the autocorrelation of the narrow-band
257 noise output of a filter centred at the target tone frequency, and τ_0 is an optimal time
258 equalization parameter. The parameters were chosen according to Durlach's original
259 formulation, e.g. γ assumes a filter with triangular gain characteristics. Psychometric functions
260 were derived under the assumption that detection thresholds represented a d' of 1 in a yes-
261 no experiment^{47,48}:

$$262 \quad R(a, D) = \Phi(0.501611 \times 10^{0.1(a-D)}) \quad (2)$$

263 where a is pure tone amplitude and D the detection threshold (**Equation 1**).

264

265 **Deep neural network.** We trained DNNs^{6,7} to predict the detection rates of tones presented
266 amongst noise, with varying ITDs. Networks took 800 input values comprised of 400 samples
267 from the left-ear waveform and 400 samples from the right-ear waveform. These inputs were
268 passed through two 100 neuron exponential linear unit layers (ELU), referred to as the
269 "decoder" portion of the network. This was followed by a layer of 10 latent Gaussian nodes
270 (>> than the parameters varied in the generation of training stimuli) with minimal uncorrelated
271 representations, constrained by a parameter β . This was followed by another two 100 neuron
272 exponential linear unit layers, referred to as the "decision" portion of the network. All layers

273 were fully connected and feed-forward. DNNs were trained and validated (95%/5% split
274 respectively) on 10^6 instances of a random phase tone at a random level (0-50 dB) in randomly
275 generated white noise, both presented with ITDs mapped from random angles in the azimuth,
276 and the corresponding detection rates (**Equation 2**).

277 The Adam optimization algorithm was used to minimize the cost function:

$$278 \quad C_{\beta}(\hat{x}, x, \sigma, \mu) = \|\hat{x} - x\|_2^2 - \frac{\beta}{2} \sum_i \log(\sigma_i^2) - \mu_i^2 - \sigma_i^2 \quad (3)$$

279 where \hat{x} and x are predicted and true detection rates respectively, and σ and μ are the
280 standard deviation and mean of latent Gaussian nodes respectively. Batch size (number of
281 training instances employed in each iterative update of network parameters) was set to 256.
282 The learning rate (training hyperparameter) was set to 5×10^{-4} for 1000 epochs (total passes
283 of entire training dataset). Ten DNNs were trained for each value of β , namely 0, 10^{-6} , 10^{-5} , 10^{-4} ,
284 10^{-3} , and 10^{-2} (60 in total). The DNN with the least RMSE, between predicted detection rates
285 and ground truth, for the validation dataset, was selected for further analysis. Central nodes
286 were considered operational if the Kullback–Leibler divergence between their individual
287 responses and a unit Gaussian was larger than 0.1 bits (**Extended Data Fig. 1b**).

288

289 **Network detection thresholds.** For a given stimulus configuration, DNN detection
290 thresholds were obtained by calculating the mean of 10 detection rates across tone levels set
291 between 0 and 50 dB in 2.5 dB steps and regressing a psychometric curve (**Equation 2**). This
292 was repeated 10 times for a given stimulus configuration. Stimuli for which detection
293 thresholds were derived included:

- 294 ▪ random phase tones amongst randomly generated broadband noise with ITDs each
295 mapped from fixed azimuthal locations,
- 296 ▪ and random phase tones and randomly generated broadband noise each either in or
297 out of phase (i.e. NoSo, NoS π , N π S π , and N π So).

298 Detection thresholds were also derived following ablations, where a set percentage of a given
299 layers nodes were randomly nullified.

300

301 **Node representations.** Node activations were measured as a function of ITD for broadband
302 noise (50-5000 Hz, 60 dB). ITDs ranged from -2000 μ s to 2000 μ s in steps of 100 μ s. Node
303 activations were also measured in response to So (in-phase) or S π (out-of-phase) signals
304 masked by either No (in-phase) or N π (out-of-phase) broadband noise. Activations in layer 2
305 nodes were measured in response to a tone level of 35 dB amongst a noise level of 60 dB (in
306 **Figure 3a**, activations were displayed with +1 added to their value because the minimum value
307 of the ELU activation function is -1)⁴. In latent Gaussian nodes (in the central layer 3), masked
308 rate-level functions to NoSo, NoS π , N π S π , and N π So were measured amongst pure tone
309 levels varied between 0 and 50 dB in 2.5 dB steps and 60 dB broadband noise. For all
310 response measurements, stimuli were 20 ms long with a sample rate of 20 kHz. For a given
311 stimulus configuration, activations were measured in response to 5000 random generations.

312

313 **Binaural cross-correlation algorithm.** For comparative purposes, outputs from a
314 binaural cross-correlation algorithm were calculated⁴⁹. The stimuli NoSo, NoS π , N π S π , and
315 N π So were generated for a 35 dB tone and a 60 dB randomly distributed broadband noise.
316 Stimuli were sampled at 20 kHz and were 1 s in duration. Signals were passed through
317 gammatone filters centered at 500 Hz and passed through a model of neural transduction⁵⁰.
318 The outputs were then delayed relative to one another, and the cross-products calculated and
319 summated.

320

321 **Statistical analysis.** ITD tuning was quantified by fitting a Gabor function¹⁵ to noise delay
322 responses. The parametric expression for a Gabor function is:

$$323 \quad G = Ae^{-(ITD-bITD)^2/2s^2} \cos(2\pi F(ITD - bITD)) + C \quad (4)$$

324 in which we characterized a nodes' best ITD as the parameter $bITD$, F is the tuning curve
325 frequency, A is a scaling factor (constrained to be positive), C is a constant offset, and s is a
326 decay constant. These parameters were fit with the non-linear least squares algorithm
327 *curve_fit* (a SciPy function⁵¹). An F-test was used to assess whether a Gabor function was a
328 significantly better fit to noise delay responses than a linear function of ITD.

329 We performed Student's two-tailed t-tests (assuming unequal variance) to assess
330 BMLDs and differences in node responses at threshold tone levels. We also used Student's
331 two-tailed t-tests (assuming unequal variance) to assess BMLDs following network ablations,
332 for which a Bonferonni correction was applied to offset the impact of testing multiple ablation
333 rates. Pearson product-moment correlation was calculated between the average responses
334 of nodes to NoSo, NoS π , N π S π , and N π So, and the delay matched outputs of a binaural
335 cross-correlation algorithm. A two-way ANOVA was also run for these nodes responses, with
336 main effects of best ITD and noise phase. For the outlined statistical analyses, the criterion for
337 significance (following multiple comparison corrections, when applied) was set to $p=0.05$. Error
338 bars and lightly shaded underlays in figures are 95% confidence intervals.

339

340 **Resource availability.** Code generated during this study is available at
341 https://github.com/Hearing-Sciences/BinauralDetection_DNN. Further information and
342 requests for resources should be directed to the Lead Contact, Samuel Smith
343 (samuel.smith@nottingham.ac.uk).

344

ACKNOWLEDGMENTS

345

346 S.S.S. and M.A.A. are supported by the Medical Research Council (grant number
347 MR/S002898/1). J.S. is funded by a Nottingham Research Fellowship from the University of
348 Nottingham. We are grateful for access to the University of Nottingham's Augusta HPC
349 service. Thanks to Alan Palmer for assistance in sharing previously published physiology data.

350

351

AUTHOR CONTRIBUTIONS

352

353
354 Conceptualization, S.S.S. and M.A.A.; Methodology, S.S.S.; Investigation, S.S.S. and J.S.;
355 Writing – Original Draft, S.S.S.; Writing – Review & Editing, S.S.S., J.S. and M.A.A.; Funding
356 Acquisition, M.A.A.; Resources, M.A.A.; Supervision, M.A.A.

357

358

COMPETING INTERESTS

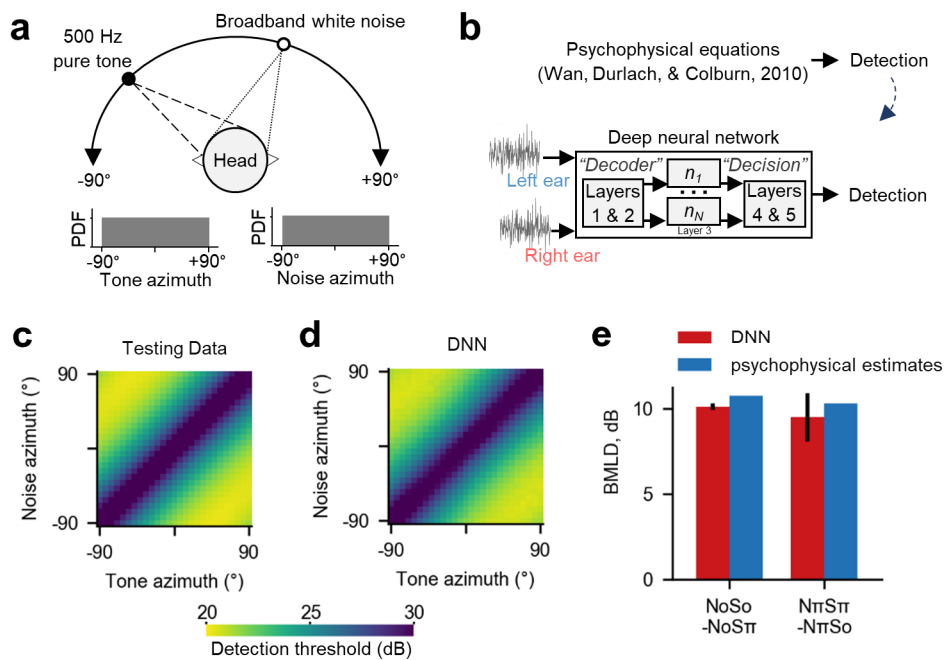
359

360

361 The authors declare no competing interests.

362

FIGURE LEGENDS



363 **Figure 1. Deep neural network accounts for binaural detection psychophysics.**

364 Data from a frontal field binaural detection task (**a**, generated using psychophysical equations)

365 were used to train DNNs (**b**) to detect a pure tone (black circle in **a**) in broadband noise (empty

366 circle in **a**). Locations (and hence ITDs) of the tone and noise were chosen at random on each

367 trial and were equally likely to come from each location (bottom panels of **a**). The DNN was a

368 5-layer network with a low-dimensional central layer, i.e. layer 3, designed to promote

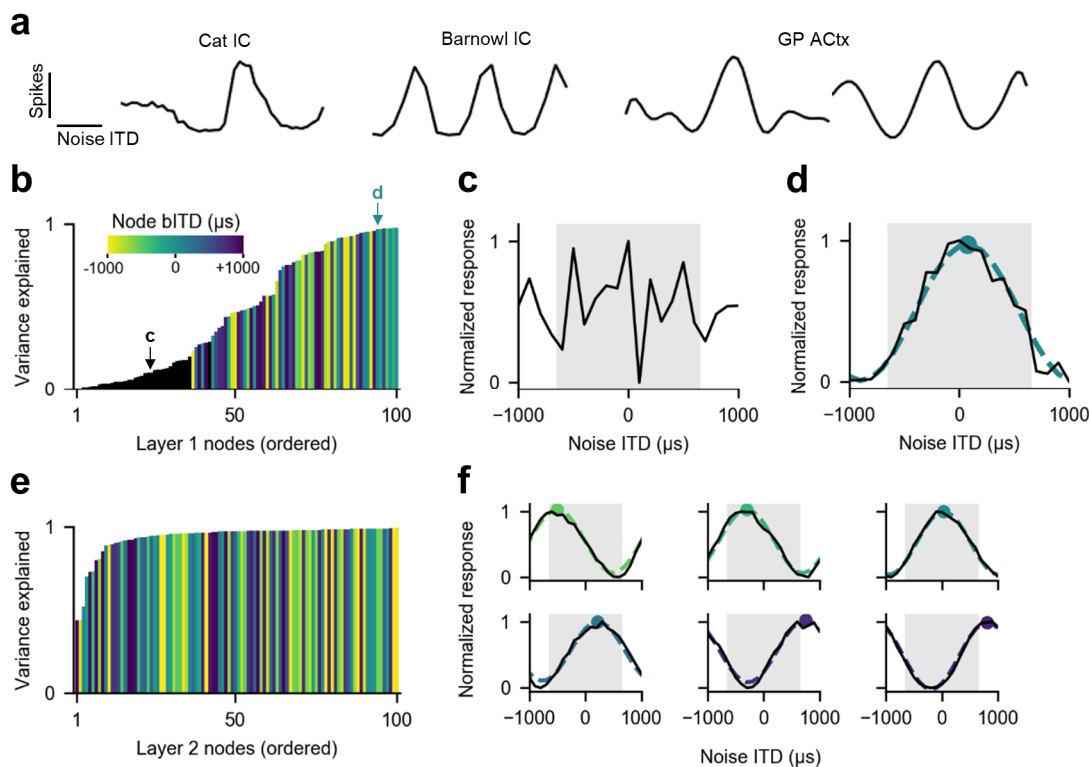
369 interpretation of the internal workings of the network (**b**). The DNN performance (**d**) for unseen

370 testing data (**c**) was found to be comparable. In addition, binaural masking level differences

371 (BMLDs) were derived for experimental stimulus configurations (NoSo/NoS π , N π S π /N π So,

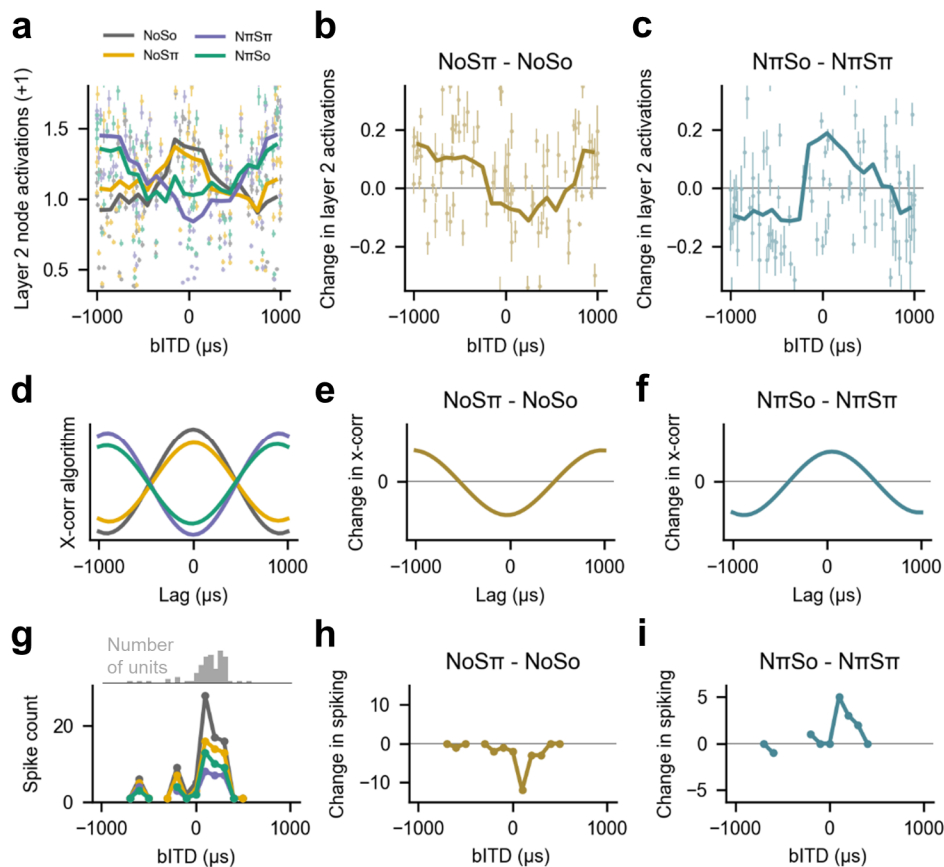
372 **e**, Note: π is beyond the DNNs trained range). Error bars for DNN represent 95% confidence

373 intervals.



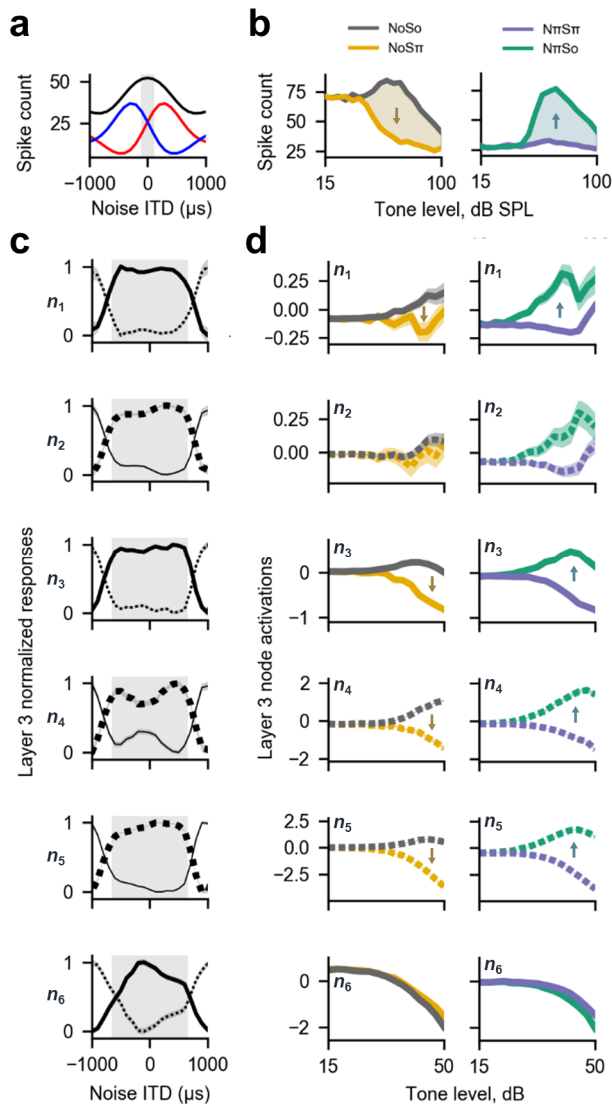
374 **Figure 2. Time delay tuning emerges early in the network.**

375 Neural ITD tuning curves have been observed in a number of animal species (a), including
 376 cat³⁵ (inferior colliculus, IC), barnowl³⁹ (IC), and guinea pig¹⁶ (auditory cortex, GP ACtx). ITD
 377 tuning emerged as a property of nodes within the early layers of the DNN and increased
 378 hierarchically between layer 1 (b) and layer 2 (e). ITD tuning was defined as the proportion of
 379 variance explained (R^2) by fits (Gabor functions) regressed to noise delay responses for nodes
 380 in the DNN's layers. Bars are color-coded by the nodes' best ITD (bITD; black indicates the
 381 Gabor fit was not significantly better than a linear fit). Individual examples of ITD tuning within
 382 a subselection of nodes in layer 1 (c and d) and layer 2 (f). The gray box underlays represent
 383 the ITD-limit for our simulation (modelled on the human head).



384 **Figure 3. Network dynamics match those of a cross-correlation mechanism of**
 385 **the type suggested in animal physiology.**

386 The internal mechanism of layer 2 of the DNN was probed by considering activation of nodes
 387 with different ITD tuning (a) for a set of typically employed binaural detection stimuli (NoSo,
 388 NoST π , N π ST π , N π So; color-coded). Single node data in light colors where error bars represent
 389 95% confidence intervals. Moving averages are overlaid and color-coded. DNN activation (a)
 390 was found to closely match a simple cross-correlation model (d). It was also comparable to
 391 animal physiological data¹⁶ (g) at the bITDs sampled (g, bottom panel shows the total spike
 392 counts of guinea pig auditory cortical neurons tuned to a given bITD, top panel shows neural
 393 count for each bITD). The activation change of nodes to paired stimulus configurations (NoSo
 394 vs NoST π and N π ST π vs N π So, b and c respectively) produced a profile across bITD (Error
 395 bars represent 95% confidence intervals. Moving averages are overlaid and color-coded). This
 396 profile matched that of a simple cross-correlation model (e and f, change in cross-correlation
 397 at different lags). In addition, a similar profile has been observed in the guinea pig animal
 398 model (h and i, change in spike count at different bITDs).

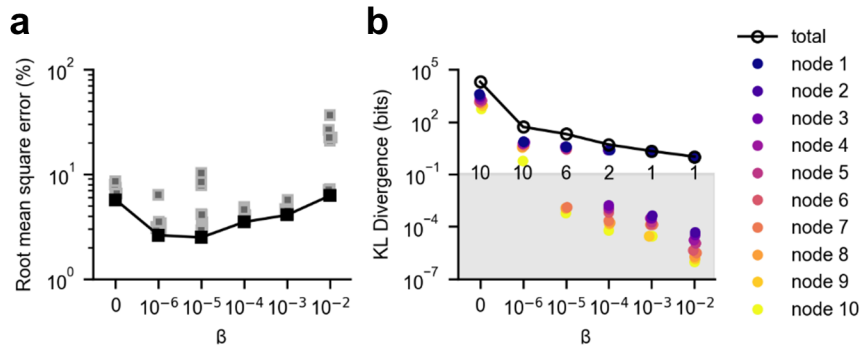


399 **Figure 4. Low-dimensional representations in DNN imitate neural signature of**
 400 **population-level cortical activity.**

401 A noise delay function from a representative neuron in right guinea pig auditory cortex¹⁶ (a,
 402 red line) alongside its reflection representative of left cortex (a, blue line) and their sum (black
 403 line). Population masked rate-level functions recorded from guinea pig auditory cortex¹⁶, in
 404 response to experimental binaural detection stimuli (NoSo, NoS π , N π S π , N π So) are shown
 405 amongst arrows indicating changes as stimuli become more easily detectable (b). Noise delay
 406 functions (c) from operational nodes in layer 3 (n_1 - n_6) are alongside masked rate-level
 407 functions (d) in response to the same stimuli configurations as in (b), i.e NoSo, NoS π , N π S π ,
 408 N π So. Dashed lines represent polarity corrected responses (such that noise delay functions
 409 in (c) are peaked as in (a)). The gray box underlays represent the ITD-limit for the guinea pig
 410 (in a) or an average human head (in c). Lightly shaded regions represent 95% confidence
 411 intervals.

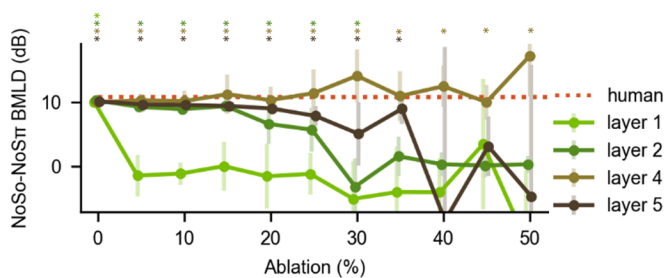
412

EXTENDED DATA



413 **Figure 1. Meta-parameter search**

414 Prediction error for 60 (10 for each value of β [see Methods]) DNNs tested on the validation
415 dataset (a). The DNNs with the minimum error for each value of β are represented with large
416 black squares. Information transmission was investigated for the most accurate DNNs for each
417 value of β (b). Displayed is the total Kullback–Leibler (KL) divergence between latent (layer 3)
418 nodes in layer 3 and an isotropic Gaussian (empty black circles). The KL divergence between
419 each individual node and a unit Gaussian is also shown (color coded in order). The gray region
420 represents nodes deemed to be suppressed during training. The number of nodes with KL
421 divergences above this region (out of 10) are typed on the upper border of this region.



422 **Figure 2. Early network ablation is detrimental to binaural detection.**

423 BMLDs (NoSo-NoSt π) were predicted following varying levels of random ablation (setting to
424 zero) to nodes in the DNN's layers. Error bars represent 95% confidence interval. Color-coded
425 asterisks indicate significant BMLDs ($p < 0.05$ following Bonferonni correction).

REFERENCES

1. Cichy, R. M. & Kaiser, D. Deep Neural Networks as Scientific Models. *Trends in Cognitive Sciences* vol. 23 305–317 (2019).
2. Khatami, F. & Escabí, M. A. Spiking network optimized for word recognition in noise predicts auditory system hierarchy. *PLOS Comput. Biol.* **16**, e1007558 (2020).
3. Kell, A. J. E., Yamins, D. L. K., Shook, E. N., Norman-Haignere, S. V. & McDermott, J. H. A Task-Optimized Neural Network Replicates Human Auditory Behavior, Predicts Brain Responses, and Reveals a Cortical Processing Hierarchy. *Neuron* **98**, 630–644.e16 (2018).
4. Koumura, T., Terashima, H. & Furukawa, S. Cascaded tuning to amplitude modulation for natural sound recognition. *J. Neurosci.* **39**, 5517–5533 (2019).
5. Furukawa, S., Terashima, H., Koumura, T. & Tsukano, H. Data-driven approaches for unveiling the neurophysiological functions of the auditory system. *Acoust. Sci. Technol.* **41**, 63–66 (2020).
6. Higgins, I. *et al.* B-VAE: Learning basic visual concepts with a constrained variational framework. in *5th International Conference on Learning Representations, ICLR 2017 - Conference Track Proceedings* (2017).
7. Iten, R., Metger, T., Wilming, H., del Rio, L. L. & Renner, R. Discovering physical concepts with neural networks. *Phys. Rev. Lett.* **124**, 010508 (2018).
8. Hish, I. J. Binaural summation and interaural inhibition as a function of the level of masking noise. *Am. J. Psychol.* **61**, 205–213 (1948).
9. Hirsh, I. J. The Influence of Interaural Phase on Interaural Summation and Inhibition. *J. Acoust. Soc. Am.* **20**, 536–544 (1948).
10. Domnitz, R. H. & Colburn, H. S. Analysis of binaural detection models for dependence on interaural target parameters. *J. Acoust. Soc. Am.* **59**, 598–601 (1976).
11. Jeffress, L. A. Binaural signal detection- Vector theory(Vector correlation theory and neural mechanisms of binaural signal detection in human auditory system). *Found. Mod. Audit. theory.* **2**, 351–368 (1972).
12. Durlach, N. I. Equalization and Cancellation Theory of Binaural Masking-Level Differences. *J. Acoust. Soc. Am.* **35**, 1206–1218 (1963).
13. Colburn, H. S. Theory of binaural interaction based on auditory-nerve data. II. Detection of tones in noise. *Cit. J. Acoust. Soc. Am.* **61**, 525 (1977).
14. Dietz, M. *et al.* A framework for testing and comparing binaural models. *Hear. Res.* **360**, 92–106 (2018).
15. Lane, C. C. & Delgutte, B. Neural correlates and mechanisms of spatial release from masking: Single-unit and population responses in the inferior colliculus. *J. Neurophysiol.* **94**, 1180–1198 (2005).
16. Gilbert, H. J., Shackleton, T. M., Krumbholz, K. & Palmer, A. R. The neural substrate

- for binaural masking level differences in the auditory cortex. *J. Neurosci.* **35**, 209–220 (2015).
17. Palmer, A. R. & Shackleton, T. M. The physiological basis of the binaural masking level difference. *Acta Acustica united with Acustica* vol. 88 312–319 (2002).
 18. Culling, J. F. Evidence specifically favoring the equalization-cancellation theory of binaural unmasking. *J. Acoust. Soc. Am.* **122**, 2803 (2007).
 19. Breebaart, J., van de Par, S. & Kohlrausch, A. Binaural processing model based on contralateral inhibition. I. Model structure. *J. Acoust. Soc. Am.* **110**, 1074–1088 (2001).
 20. Wack, D. S. *et al.* Functional Anatomy of the Masking Level Difference, an fMRI Study. *PLoS One* **7**, e41263 (2012).
 21. Wack, D. S., Polak, P., Furuyama, J. & Burkard, R. F. Masking Level Differences – A Diffusion Tensor Imaging and Functional MRI Study. *PLoS One* **9**, e88466 (2014).
 22. Sasaki, T. *et al.* Neuromagnetic evaluation of binaural unmasking. *Neuroimage* **25**, 684–689 (2005).
 23. Fowler, C. G. Electrophysiological evidence for the sources of the masking level difference. *Journal of Speech, Language, and Hearing Research* vol. 60 2364–2374 (2017).
 24. Jiang, D., McAlpine, D. & Palmer, A. R. Detectability index measures of binaural masking level difference across populations of inferior colliculus neurons. *J. Neurosci.* **17**, 9331–9339 (1997).
 25. Schnupp, J. W. H. & Carr, C. E. On hearing with more than one ear: Lessons from evolution. *Nature Neuroscience* vol. 12 692–697 (2009).
 26. Ching, T. Y. C., van Wanrooy, E., Dillon, H. & Carter, L. Spatial release from masking in normal-hearing children and children who use hearing aids. *J. Acoust. Soc. Am.* **129**, 368–375 (2011).
 27. Adavanne, S., Politis, A., Nikunen, J. & Virtanen, T. Sound Event Localization and Detection of Overlapping Sources Using Convolutional Recurrent Neural Networks. *IEEE J. Sel. Top. Signal Process.* **13**, 34–48 (2018).
 28. Vecchiotti, P., Ma, N., Squartini, S. & Brown, G. J. End-to-end Binaural Sound Localisation from the Raw Waveform. in *ICASSP, IEEE International Conference on Acoustics, Speech and Signal Processing - Proceedings* vols 2019-May 451–455 (Institute of Electrical and Electronics Engineers Inc., 2019).
 29. Francl, A. & McDermott, J. Deep neural network models of sound localization reveal how perception is adapted to real-world environments. *bioRxiv* 2020.07.21.214486 (2020) doi:10.1101/2020.07.21.214486.
 30. Klein, M. A. & Hartmann, W. M. Binaural edge pitch. *J. Acoust. Soc. Am.* **70**, 51–61 (1981).
 31. Hartmann, W. M. & McMillon, C. D. Binaural coherence edge pitch. *J. Acoust. Soc. Am.* **109**, 294–305 (2001).

32. Wan, R., Durlach, N. I. & Colburn, H. S. Application of an extended equalization-cancellation model to speech intelligibility with spatially distributed maskers. *J. Acoust. Soc. Am.* **128**, 3678–3690 (2010).
33. Woodworth, R., Barber, B. & Schlosberg, H. Experimental psychology. (1954).
34. Durlach, N. I. & Colburn, H. S. Binaural phenomena. in *Handbook of perception, Vol IV, Hearing* (eds. Carterette, E. C. & Friedman, M. P.) 365–447 (New York: Academic, 1978).
35. Joris, P. & Yin, T. C. T. A matter of time: internal delays in binaural processing. *Trends in Neurosciences* vol. 30 70–78 (2007).
36. Singheiser, M., Gutfreund, Y. & Wagner, H. The representation of sound localization cues in the barn owl's inferior colliculus. *Front. Neural Circuits* **6**, 45 (2012).
37. Jeffress, L. A. A place theory of sound localization. *J. Comp. Physiol. Psychol.* **41**, 35–39 (1948).
38. McAlpine, D., Jiang, D. & Palmer, A. R. Binaural masking level differences in the inferior colliculus of the guinea pig. *J. Acoust. Soc. Am.* **100**, 490–503 (1996).
39. Asadollahi, A., Endler, F., Nelken, I. & Wagner, H. Neural correlates of binaural masking level difference in the inferior colliculus of the barn owl (*Tyto alba*). *Eur. J. Neurosci.* **32**, 606–618 (2010).
40. Scholte, H. S. Fantastic DNimals and where to find them. *NeuroImage* vol. 180 112–113 (2018).
41. Malhotra, G., Evans, B. D. & Bowers, J. S. Hiding a plane with a pixel: examining shape-bias in CNNs and the benefit of building in biological constraints. *Vision Res.* **174**, 57–68 (2020).
42. Meyes, R., Lu, M., de Puisseau, C. W. & Meisen, T. *Ablation studies in artificial neural networks*. *arXiv* (arXiv, 2019).
43. Olsen, W. O., Noffsinger, D. & Carhart, R. Masking level differences encountered in clinical populations. *Int. J. Audiol.* **15**, 287–301 (1976).
44. Goodman, D. F. M. & Brette, R. Learning to localise sounds with spiking neural networks. in *Advances in Neural Information Processing Systems 23: 24th Annual Conference on Neural Information Processing Systems 2010, NIPS 2010* (2010).
45. Rayleigh, Lord. XII. On our perception of sound direction. *London, Edinburgh, Dublin Philos. Mag. J. Sci.* **13**, 214–232 (1907).
46. Domnitz, R. H. & Colburn, H. S. Analysis of binaural detection models for dependence on interaural target parameters. *J. Acoust. Soc. Am.* **59**, 598–601 (1976).
47. Egan, J. P., Lindner, W. A. & Mcfadden, D. *Masking-level differences and the form of the psychometric function*. (1969).
48. Ingleby, J. D. Signal detection theory and psychophysics. *J. Sound Vib.* **5**, 519–521 (1967).
49. Akeroyd, M. A binaural cross-correlogram toolbox for MATLAB. (2017).

50. Meddis, R., Hewitt, M. J. & Shackleton, T. M. Implementation details of a computation model of the inner hair-cell/auditory-nerve synapse. *J. Acoust. Soc. Am.* **87**, 1813–1816 (1990).
51. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nat. Methods* **17**, 261–272 (2020).