

## **Automated detection of Bornean white-bearded gibbon (*Hylobates albibarbis*)**

### **vocalisations using an open-source framework for deep learning**

A. F. Owens<sup>1</sup>, Kimberley Hockings<sup>2</sup>, Muhammed Ali Imron<sup>3</sup>, Shyam Madhusudhana<sup>4,5</sup>,  
Mariaty Ayudia Niun<sup>6</sup>, Tatang Mitra Setia<sup>7</sup>, Manmohan Sharma<sup>1</sup>, Siti Maimunah Soebagio<sup>6,8</sup>,  
F. J. F. Van Veen<sup>1\*</sup> & Wendy M. Erb<sup>5</sup>

<sup>1</sup>Faculty of Environment, Science and Economy, Department of Earth and Environmental Science, Centre for Geography and Environmental Science, University of Exeter, Penryn, TR10 9FE, UK

<sup>2</sup>Centre for Ecology and Conservation, University of Exeter, Penryn, TR10 9FE, UK

<sup>3</sup>Faculty of Forestry, Universitas Gadjah Mada, Yogyakarta, 55281, Indonesia

<sup>4</sup>Centre for Marine Science and Technology, Curtin University, Perth, WA, 6102, Australia

<sup>5</sup>K. Lisa Yang Centre for Conservation Bioacoustics, Cornell Lab of Ornithology, Cornell University, Ithaca, NY 14850, USA

<sup>6</sup>Fakultas Kehutanan dan Pertanian, Universitas Muhammadiyah Palangka Raya, Palangka Raya, 73111, Indonesia

<sup>7</sup>Department of Biology, Faculty of Biology and Agriculture, Universitas Nasional, Jakarta, 12520, Indonesia

<sup>8</sup>Fakultas Kehutanan, Instiper Yogyakarta, Yogyakarta, 55281, Indonesia

1 **ABSTRACT**

2 Passive acoustic monitoring is a promising tool for monitoring at-risk populations of vocal  
3 species, yet extracting relevant information from large acoustic datasets can be time-  
4 consuming, creating a bottleneck at the point of analysis. We adapted an open-source  
5 framework for deep learning in bioacoustics to automatically detect Bornean white-bearded  
6 gibbon (*Hylobates albibarbis*) “great call” vocalisations in a long-term acoustic dataset from a  
7 rainforest location in Borneo. We describe the steps involved in developing this solution, such  
8 as collecting audio recordings, developing training and testing datasets, training neural-network  
9 models, and evaluating model performance. Our best model performed at a satisfactory level (F  
10 score = 0.87), identifying 98% of the highest-quality calls from 90 hours of manually-annotated  
11 audio recordings. We also found no significant difference in the distribution of great call  
12 detections over time between the manual annotations and the model’s output, and greatly  
13 reduced analysis times when compared to a human observer. Future work should seek to apply  
14 our model to long-term acoustic datasets to understand spatiotemporal variations in *H.*  
15 *albibarbis*’ calling activity. With additional information, such as detection probability over  
16 distance, we demonstrate how our model could be used to monitor gibbon population density  
17 and spatial distribution on an unprecedented scale.

18

19

20

21

22

## 23 I. INTRODUCTION

24 Ever-increasing anthropogenic pressures on the environment, such as habitat loss, have  
25 led to widespread population declines in many animal species (Bender et al. 1998). However,  
26 for many species, data in population trends is often sparse and lacking in standardised  
27 methodology (Jetz et al. 2019), leading to a demand for a greater quantity and scale of wildlife  
28 population monitoring programmes to inform conservation responses (Verma et al. 2016). To  
29 meet this demand, conservation scientists and ecologists have turned to developing  
30 technologies to automate data collection, enabling the rapid accumulation of large volumes of  
31 data (Piel and Wich 2021). While this has allowed for unprecedented insight, it can also make  
32 practical aspects of ecological inference challenging (Borowiec et al. 2022).

33 Manual extraction of relevant information from large datasets can be time consuming,  
34 resulting in a bottleneck at the point of analysis (Norouzzadeh et al. 2018). This bottleneck is  
35 evident in data generated as part of passive acoustic monitoring (PAM) programmes, which  
36 involve the use of acoustic sensors to autonomously collect sound recordings in the field  
37 (Acevedo and Villanueva-rivera 2010). Advancements in recording device design, greater  
38 availability of lower-cost equipment, and improved data storage options have made the task of  
39 capturing many hours of acoustic data relatively straightforward (Morgan and Braasch 2021).  
40 Data must then be browsed to identify relevant signals of interest, such as species-specific  
41 vocalisations, often by manually listening to each recording in full or visually inspecting the data  
42 in spectrogram form (a time-frequency pictorial representation of an audio signal), or both.

43 PAM can provide a step-change in standardised population monitoring of vocal species at high  
44 temporal resolution and simultaneous large spatial scales, that would be impossible to achieve  
45 with 'traditional' methods relying on manual data collection (Sugai et al. 2019). However, such  
46 PAM programmes often capture datasets so large that they cannot be studied manually in full  
47 in a reasonable timeframe, so automating this limiting data-processing step is critical (Morgan  
48 and Braasch 2021; Clink et al. 2023).

49 Machine learning has proven to be an effective solution for fast and accurate analysis of  
50 acoustic data, including the automated detection of signals of interest (Stowell 2022; Miller et  
51 al. 2023). There are many options available for this task, including artificial neural networks  
52 (ANNs) (Mielke and Zuberbühler 2013), Gaussian mixture models (GMMs) (Heinicke et al. 2015)  
53 and support vector machines (SVMs) (Noda et al. 2016), among others. These each have  
54 associated advantages, and due to the diversity of potential signal types and acoustic  
55 environments, no single method has been shown to be optimal in all situations (Clink et al.  
56 2023). However, it is worth noting that ANNs demonstrate comparatively strong adaptability  
57 and proficiency in understanding complex patterns in data (Haykin 2009; Bengio et al. 2016).  
58 Early work applying ANNs to animal sound made use of the multi-layer perception (MLP)  
59 architecture, with manually selected summary features such as syllable duration, peak  
60 frequency, etc., used to inform the network's predictions (Stowell 2022). While MLPs have been  
61 effective in classifying a wide variety of terrestrial and marine animal calls, the structure of non-  
62 speech acoustic events can be highly variable (Kong et al. 2017), and reducing the data to a  
63 series of manually assigned summary features can restrict the wealth of information available  
64 to train a network, potentially limiting its effectiveness (Stowell 2022).

65 An alternative is to use convolutional neural networks (CNNs) which, like other deep-  
66 learning ANN architectures, rely on feature sets that are not manually selected, but instead  
67 learned during the training process (Morgan and Braasch 2021). CNNs can be applied to visual  
68 representations of audio, such as spectrograms, leveraging their ability to learn patterns that  
69 occur both spatially and temporally in data. This allows CNNs to learn local features regardless  
70 of their spatial position within an image (Knight et al. 2017). CNNs are therefore ideal  
71 candidates for the automated detection of signals within bioacoustic data, where instances of  
72 relevant features within a spectrogram are not predefined or readily identifiable.

73 The suitability of CNNs is making them a highly popular choice of model to process  
74 acoustic data (Stowell 2022). They have been used to analyse vocalisations from a variety of  
75 taxa, including insects (Hibino et al. 2021), fish (Guyot et al. 2021), anurans (Colonna et al.  
76 2016), birds (Narasimhan et al. 2017), bats (Mac Aodha et al. 2018), marine mammals (Miller et  
77 al. 2023) and terrestrial mammals (Bjorck et al. 2019), including primates (Dufourq et al. 2021).  
78 Their potential is far from fully realised however (Rammer and Seidl 2019), and there are  
79 relatively few examples of CNNs being used to answer well-defined research questions in  
80 ecology, as is so with other deep learning approaches (Dufourq et al. 2021). Additionally, there  
81 are few guidelines on how to approach key steps such as model tuning and performance  
82 assessments (Josh Patterson and Adam Gibson 2017; Knight et al. 2017; Stowell 2022). Further  
83 case studies reporting successful applications will advance the development of best practices  
84 for overcoming these challenges (Dufourq et al. 2021).

85 Gibbons (family Hylobatidae) are ideal candidates for the automated detection of  
86 species-specific vocalisations. As they reside exclusively in tropical forests, which are often

87 visually challenging and inaccessible, studying their populations using visual methods, such as  
88 line transect and camera trap surveys, is typically very difficult (Vu and Tran 2019). Gibbons  
89 engage in loud, highly stereotyped song bouts which are largely confined to a specific daily  
90 temporal period (Cheyne et al. 2008). During a particular calling bout, they usually emit  
91 multiple calls, allowing for ample training data (Clink et al. 2023). Here, we apply a variation of a  
92 pre-defined CNN architecture, DenseNet (Huang et al. 2016), to identify female great call  
93 vocalisations of the endangered Bornean white-bearded gibbon (*Hylobates albibarbis*) from a  
94 long-term acoustic dataset. Since the great call is performed largely by mated females, it is  
95 often used as an indicator of a gibbon family group, allowing for group density and spatial  
96 distribution estimates to be derived from great call densities (Cheyne et al. 2016). We train a  
97 detector with high precision that minimises false-positive rates (i.e., the rate of detections  
98 incorrectly labelled as the positive class) for application to large acoustic datasets recorded by  
99 PAM arrays. This can then be used to facilitate accurate population monitoring of wild gibbons  
100 on an ever-greater spatiotemporal scale and applied as a case study for developing PAM  
101 frameworks for other endangered loud-call species.

## 102 II. METHODS

### 103 A. Data collection

104 The long-term acoustic dataset used in this study derives from the Mungku Baru  
105 Education Forest (MBEF), a ~50 km<sup>2</sup> area of tropical rainforest in Central Kalimantan, Indonesia.  
106 The MBEF lies in the centre of the wider Rungan Forest Landscape, which covers 1,500 km<sup>2</sup>  
107 between the Kahayan and Rungan rivers north of the provincial capital of Palangka Raya. This

108 represents the largest area of unprotected lowland rainforest remaining on the island of  
109 Borneo (Purnama and Afifah 2021). There is an estimated population of 4,000 white bearded  
110 gibbons in the Rungan Landscape and an estimated density of 2.79 groups per km<sup>2</sup> in the MBEF,  
111 making it of significant importance for the conservation of the species (Buckley et al. 2018).

112         Eight autonomous recording units (ARUs) (Song Meter SM4, Wildlife Acoustics,  
113 Maynard, Massachusetts) were deployed here by WME in July 2018. These were placed on  
114 trees, 5 meters high, in a dispersed grid with approximately 1,200 meters between each device.  
115 The MBEF contains a mosaic of different forest types, and the location of the ARUs was  
116 intended to reflect this heterogeneity, with three situated in ‘kerangas’ heath forest, three in  
117 ‘low pole’ peat swamp forest, and two in ‘mixed swamp forest’, with the latter representing a  
118 transition habitat between the former two (Buckley et al. 2018).

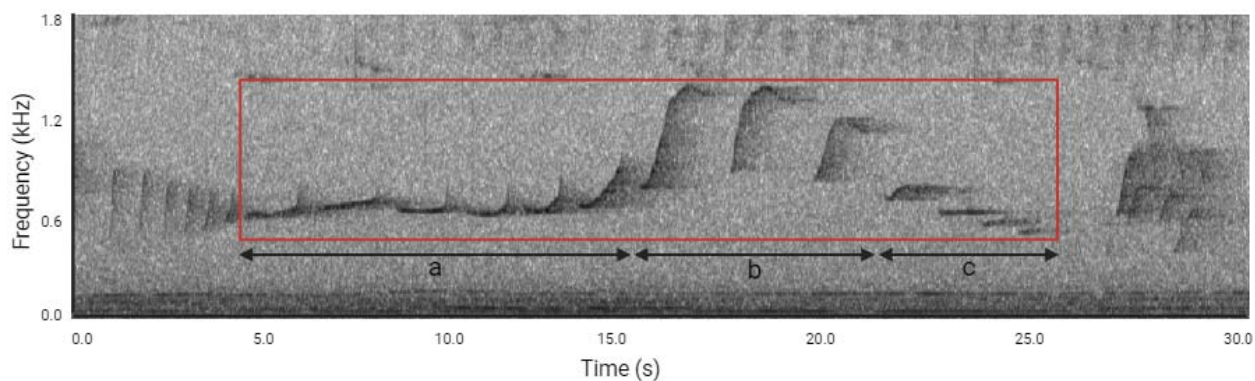
119         The ARUs were set to record continuously from 4 am to 6 pm local time (WIB/UTC +7)  
120 daily between 18<sup>th</sup> July 2018 to 5<sup>th</sup> January 2020. These were set to default settings (see  
121 Appendix [i]) and recorded on two channels with a sampling rate of 24 kHz. Audio was captured  
122 in 16-bit Waveform Audio File Format (WAV) and saved as one-hour files. Memory cards and  
123 batteries were changed every two weeks.

## 124         **B. Manual annotation**

125         Separate manually-annotated training and testing datasets were required to develop  
126 the automated detector. To create these samples, recordings between 4-10am were selected  
127 from a single day every four weeks from a randomly selected device for each habitat. This

128 covers the temporal period in which most *H. albibarbis* great calls occur (Cheyne et al. 2008)  
129 and ensured that a variety of potential sound environments were included as training inputs to  
130 the model, improving its ability to generalise over a wider range of applications. The resultant  
131 subset contained 300 hours of recording, covering 50 days spanning from October 2018 to  
132 December 2019.

133 The selected sound files were then loaded into the sound analysis software Raven Pro  
134 1.6, (K. Lisa Yang Center for Conservation Bioacoustics, Ithaca, NY, USA) and visualised as  
135 spectrograms (see Appendix [ii]). With assistance from a team of undergraduate interns, each  
136 recording was listened to in full and visually scanned to identify great call events. A selection  
137 was created for each event by drawing a box around the call in the spectrogram, providing  
138 information about its time-frequency boundaries (see Figure 1).



139  
140 FIG. 1. A spectrogram image of a female *H. albibarbis* great call, created using Raven Pro 1.6.  
141 This vocalisation comprises of introductory (a), climax (b) and descending (c) notes. The red box  
142 represents the time and frequency boundaries of a manually-annotated selection. The time  
143 boundaries span from the start of the first introductory note to the end of the last descending

144 note. The frequency boundaries span from the lowest frequency descending note up to the  
145 highest frequency climax note.

146 Each selection was annotated for its completeness and quality. A selection was marked  
147 as “clear” when the entirety of the call could be heard and was visually clear in the  
148 spectrogram, “faint” when the whole call could be heard but was not fully shown in the  
149 spectrogram (and vice versa), and very faint when the call was only partially seen and heard  
150 (i.e., part of the great call was not captured in the recording). Each annotation was confirmed  
151 by AFO to reduce the influence of inter-observer variability. In total, 1611 great calls were  
152 annotated.

153 The manually-annotated data was then randomly split, with 70% allocated for training  
154 (210 hours, 1089 calls) and 30% allocated for testing (90 hours, 522 calls). Due to the  
155 ambiguous nature of “very faint” calls, they were removed from the training process and used  
156 solely for testing thereafter. This prevented misleading information, i.e., non-target events,  
157 being fed into the positive class weightings. Models were trained using “clear” and “faint”  
158 instances (729 calls), as well as only “clear” instances (522 calls).

### 159 **C. Automated detection**

160 The development and testing of the automated detector utilised Koogu (version 0.7.2)  
161 (Madhusudhana 2023), an open-source framework for deep learning from bioacoustic data.  
162 Koogu offers a variety of functions for deep learning, including preparing audio for use as inputs  
163 to machine learning models, training models, assessing their performance, and using trained

164 models for the automated analysis of large datasets. For a full workflow describing how the  
165 following steps were implemented within Koogu, see Supplementary Material A.

166 **1. Data preparation**

167 All audio files were first down sampled to 4,500 Hz to reduce the overall file size and  
168 improve the efficiency of downstream computations (Miller et al. 2023). The resulting Nyquist  
169 frequency (2,250 Hz) is above the highest manually-annotated great call frequency (2,077 Hz),  
170 and so no relevant information was lost in this process. The recordings were then split into  
171 consecutive 28 second (s) segments (longer than the longest manually-annotated great call at  
172 27.7s) with a hop size of 1s, leaving an overlap of 27s between clips. The waveform of each  
173 segment was normalised by scaling the amplitudes to occur in the range -1.0 to 1.0.

174 The resulting start and end times of each segment were then compared to those from  
175 manual annotations. Segments that fully contained the temporal extents of an annotated great  
176 call were considered as positive inputs while segments without temporal overlap were  
177 considered as negative inputs. Segments with partial overlap were excluded from training as  
178 these could resemble non great call events and so lead to uncertainties during the training  
179 process.

180 Spectrograms for both the positive and negative class were then computed (see  
181 Appendix [iii]), resulting in input spectrograms with a shape of 384 x 580 (height x width) pixels.  
182 To address the imbalance between positive and negative classes, the maximum number of  
183 training inputs for each class was reduced to 10,000. This utilised all the positive class inputs  
184 while randomly subsampling inputs from the negative class. Following this, there were 5,763

185 “clear” and 1,253 “faint” positive class spectrograms as well as 10,000 negative class  
186 spectrograms.

## 187 **2. Data augmentation**

188 Despite the manually-annotated calls showing a high level of variance in background  
189 noise, call duration and note length, for example, data augmentation was applied to further  
190 improve input variance. To do this, several pre-defined augmentations supported by Koogu  
191 were applied both on waveforms before conversion into spectrograms, and on the  
192 spectrograms themselves. These were performed at each epoch, i.e., each time the model  
193 passed through the entire training dataset, during the feeding of inputs into the model. This  
194 meant that the same original sample could have different levels of augmentations between  
195 epochs.

196 Firstly, Gaussian noise (Schlüter and Grill 2015) was added to 25% of the training input’s  
197 waveforms at each epoch to simulate varying levels of background noise. The amount of noise  
198 added randomly varied from -20dB to -30dB below the peak dB of the input signal.

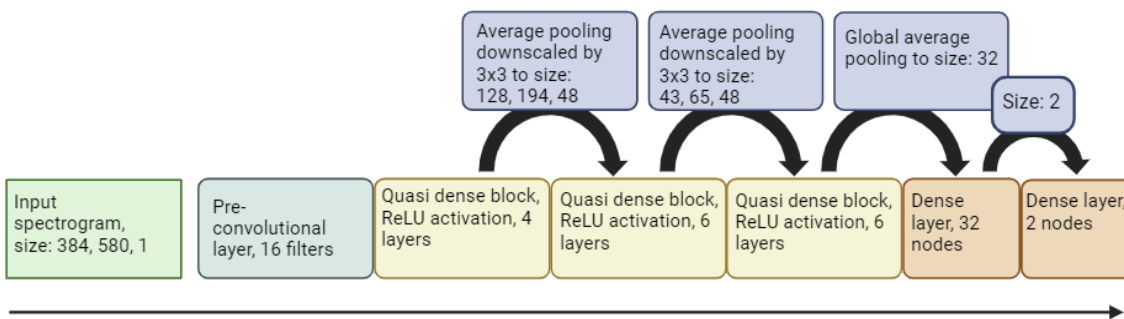
199 The spectrogram was then smeared and squished along the time axis (Madhusudhana  
200 2023). These augmentations were independently added to 20% of the input spectrograms each  
201 epoch, at a magnitude of -1,1 (smearing backwards and forwards by one frame of the  
202 spectrogram) and -2,2 (stretching and squishing over up to 2 frames of the spectrogram). This  
203 process essentially blurred inputs along the time axis, as while the target calls were contained  
204 within a standard frequency range, the duration of the signals was highly variable.

205           Finally, Koogu’s “AlterDistance” augmentation was applied on 25% of the spectrogram  
206 inputs. This aimed to mimic the effect of increasing or reducing the distance between the  
207 calling gibbon and the receiver, by attenuating or amplifying higher frequencies while keeping  
208 lower frequencies relatively unchanged. This was applied by a random factor between -5dB  
209 (attenuation) and 5dB (amplification).

### 210           **3. Network parameters and training**

211           The DenseNet architecture (Huang et al. 2016) was chosen as the base CNN architecture for this  
212 study. Early variations of the model suffered from overfitting, occurring when the model learns  
213 noise or random fluctuations in the training data rather than the underlying pattern itself. This  
214 is a symptom of when the model is too complex relative to its intended task. Bearing this in  
215 mind, the standard DenseNet architecture was adapted to a “quasi-DenseNet” architecture  
216 (Madhusudhana et al. 2021). This reduces the number of connections within each dense block,  
217 limiting model size and complexity. To limit the model’s complexity further and improve  
218 computational efficiency, bottleneck layers were also added (Huang et al. 2016). Finally, batch  
219 normalisation was enabled to improve model convergence. For the final model architecture,  
220 see Figure 2.

221



222

223 FIG. 2. Flowchart showing the final model architecture. The final model had a growth rate of 12,  
224 began with a 16-filter pre-convolutional layer, contained 3 quasi-dense blocks with 4, 6 and 6  
225 layers respectively and finished on a 32-node dense layer. Average pooling layers downscaled  
226 the inputs by a factor of 3x3 (height x width) in the transition blocks between quasi-dense  
227 blocks. Global Average pooling was used to reduce the spatial dimensions of outputs of the  
228 final block to the 32-node feature vectors.

229 Training inputs were then divided further, with 15% randomly selected as a validation  
230 set to evaluate the model's performance throughout the training process. Dropout layers were  
231 added (Srivastava et al. 2014) at a rate of 5% to further reduce overfitting and improve  
232 generalisation. The models were then trained over 80 epochs using the Adam optimizer  
233 (Kingma and Ba 2014) with a minibatch size of 24. The learning rate was initially set at 0.01 and  
234 then reduced successively by a factor of 10 at epochs 20, and 40.

#### 235 **4. Testing**

236 Trained models were then applied to the test dataset to provide a preliminary  
237 assessment of model performance and establish a desirable detector threshold value. To do  
238 this, each test segment was assigned a confidence score by the model between 0 and 1

239 indicating how likely it was to contain a great call. For thresholds 0 to 1, with an interval of 0.01,  
240 performance scores were outputted containing the number of true positives (TP), false  
241 positives (FP) and false negatives (FN). If there was a 100% overlap between an input and an  
242 annotated great call, and the confidence score was above the threshold it was marked as a TP.  
243 If there was no- or partial overlap and the score was above the threshold it was marked as an  
244 FP. If there was full overlap but the confidence score was below the threshold, then it would be  
245 marked as a FN.

246         These quantities were used to compute recall, precision, and F score at each threshold.  
247 Recall was defined as:  $TP/(TP+FN)$ , precision as:  $TP/(TP+FP)$ , and F score as:  $2(P \times R)/(P+R)$  (where  
248 P = precision and R = recall). The optimal threshold was then selected using maximum F score,  
249 as it has been shown to be a good indicator of overall model performance (Clink et al., 2023).

250         Once a threshold had been decided, the model was re-run on the test dataset audio files  
251 to produce detections and analyse the models' output. Neighbouring segments for which  
252 scores were above the identified threshold were grouped together to form a single detection.  
253 The score of each detection was then set as the maximum of its component segment scores,  
254 and its start and end times were set to the start time of the first segment and the end time of  
255 the last segment respectively. This considered every segment generated, regardless of the  
256 amount of overlap with any annotation. These detections were then outputted in the format of  
257 Raven Pro selection tables.

#### 258         **D. Post-processing**

259 Initial inspection of the models' outputs showed that they produced clusters of  
260 detections for each great call, overestimating the number of calls within the data. To minimise  
261 the number of repeat detections for each target event, these were grouped together further  
262 (see Supplementary Material B). If any detections overlapped with the first model detection,  
263 then they were assigned to the same group. This process iterated across all the models'  
264 detections, forming groups based off the first detection in each cluster of detections. These  
265 groups were then filtered to retain only the highest scoring detection(s) within each group.  
266 Where the start times of remaining selections were the same, only one detection was retained.  
267 There was no further judgement between remaining detections where their scores were tied, as  
268 this could indicate that two great calls overlapped or occurred close in time.

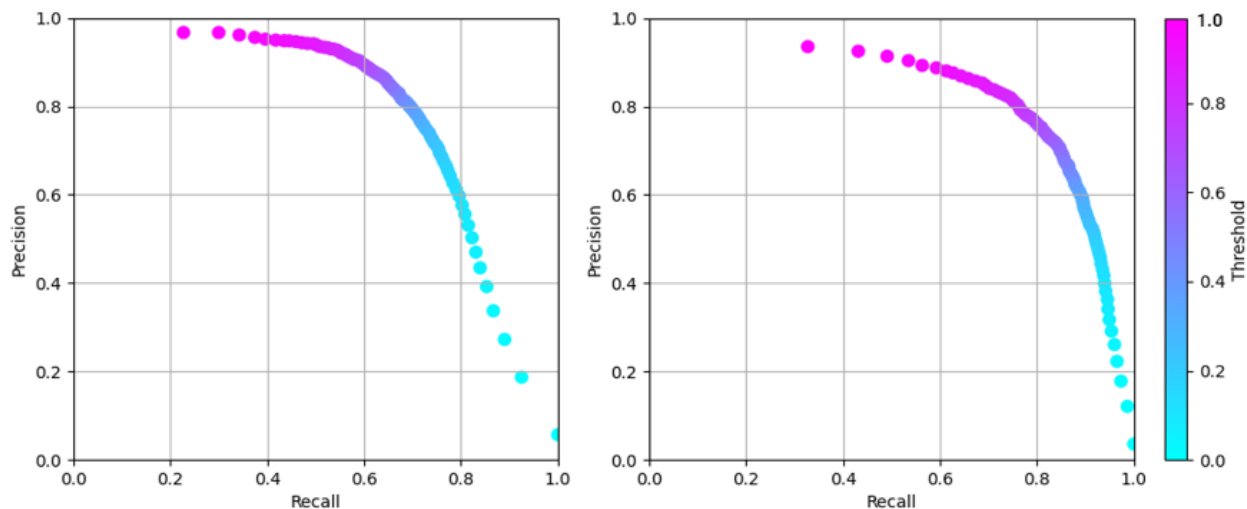
### 269 III. RESULTS

#### 270 A. Preliminary assessment

271 The preliminary assessment of the best performing model, trained on only "clear" calls,  
272 showed precision ranging from 0.19 to 0.97 and recall ranging from 0.93 to 0.23 at thresholds  
273 from 0.01 to 0.99 (Fig. 3). This gave a maximum F score of 0.75 at a threshold of 0.36 (precision:  
274 0.80, recall: 0.70).

275 As the aim was to optimise the model for "clear" and "faint" calls, the testing was re-run  
276 excluding "very faint" calls. In this case, precision ranged from 0.12 to 0.94 and recall ranged  
277 from 0.99 to 0.33 at thresholds from 0.01 to 0.99 (Fig. 3). The maximum F score was improved  
278 to 0.78 at a threshold of 0.78 (precision: 0.80, recall: 0.76). To maximise performance for

279 “clear” and “faint” calls while minimising false positives, a threshold of 0.78 was therefore  
280 chosen.



281  
282 FIG. 3. Precision-recall curves of the best performing model tested against all calls (left) and  
283 both “clear” and “faint” calls (right).

## 284 B. Comparison with manual annotations

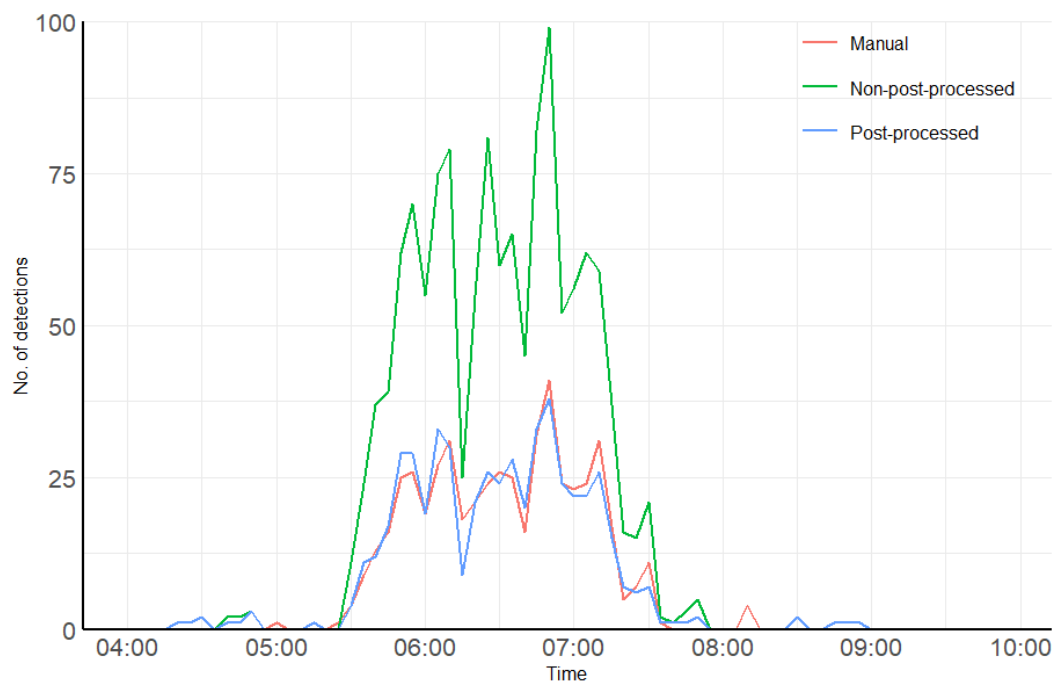
285 After re-running the chosen model on the test dataset at the desired threshold and  
286 processing the output, it produced 535 detections. These were then visually analysed in Raven  
287 Pro and compared to the manually-annotated dataset to discern the true number of TPs, FPs,  
288 and FNs. In this case, a TP instance was defined as any model detection that overlapped with a  
289 great call. The model was found to have produced 511 TPs and 24 FPs, missing a further 133  
290 calls (FNs). This gives a precision of 0.96, a recall of 0.79, and an F score of 0.87.

291 Upon closer inspection of the post-processed model output, 86 of the TP detections  
292 were the result of repeated detections for singular great call events. Additionally, 22 FNs were

293 classed as TPs before the post-processing stage. These were created as a result of grouping  
294 together detections where two great calls overlapped in time or were adjacent to one another,  
295 despite them being within the bounds of the model's detections prior to processing the output.

296 Overall, the non-post-processed model identified 409 (78%) of the 522 manually-  
297 annotated calls. This included 98% of all "clear" calls, 73% of "faint" calls, and 44% of "very  
298 faint" calls. Out of the FN instances before post-processing, 77% were "very faint" calls with  
299 only 0.05% representing 6 missed "clear" annotations. A further 38 great calls were identified  
300 which had been missed in the manual annotation stage, including 6 "clear", 8 "faint" and 24  
301 "very faint" calls.

302 The frequency of great call detections over time for both the non-post-processed and  
303 post-processed model outputs were compared against the manually-annotated dataset. A  
304 Kolmogorov–Smirnov test indicated no significant difference ( $D = 0.036$ ,  $p > 0.05$  and  $D = 0.045$ ,  
305  $p > 0.05$  respectively). Also, we found no significant difference between pre-processed and  
306 post-processed data ( $D = 0.020$ ,  $p > 0.05$ ).



307

308 FIG. 4. Histogram showing the number of great calls detected every 5 minutes between 04:00  
309 and 10:00 for the manual annotations, the non-post-processed model output, and the post-  
310 processed model output.

#### 311 IV. DISCUSSION

312 The results show how our best performing model was effective at detecting high-quality  
313 *H. albibarbis* great calls with a low rate of false positive instances. The best performing model (F  
314 score: 0.87) exceeded previously reported SVM models for detecting gibbon vocalisations e.g.,  
315 *Hylobates funereus* detector, F score: 0.78 (Clink et al. 2023). It also performed comparatively  
316 to other CNN models, e.g., *Nomascus hainanus* detector, F score: 0.91 (Dufourq et al. 2021),  
317 *Nomascus concolor* detector, F score: 0.92 (Zhou et al. 2023).

318 We found that the best performing model detected 38 great calls that had been missed  
319 during the manual annotation stage, amounting to 0.07% of the total number of target events.  
320 Human listening is subject to error (Brauer et al. 2016; Knight et al. 2017), and this study relied  
321 on multiple human observers with differing levels of training to construct the manually-  
322 annotated dataset. While it has been shown that human observers with less experience may  
323 perform worse than some automated detectors (Jennings et al. 2008), the consensus from  
324 multiple observers may have actually reduced the level of human error (Drake et al. 2016).  
325 Furthermore, signals with low signal to noise ratio (SNR) are difficult to detect for both humans  
326 and automated detectors (Knight et al. 2017). Precision and recall are often measured relative  
327 to manually-annotated datasets, but these are not always perfect. With this in mind, it has  
328 been recommended to view the manually-annotated dataset as the output from an alternative  
329 detector rather than a ground truth set (Knight et al. 2017).

330 Due to the apparent inevitability for great calls to be missed by manual annotation, it is  
331 unrealistic to presume that all target calls were identified in the 210-hour training dataset. This  
332 could pose a problem when training the model if any of the randomly selected unannotated  
333 time periods for the negative class contained target events, potentially increasing the rate of  
334 false negatives. Dufourq et al. (2021) found that better results were obtained by specifically  
335 including negative class segments with typical ambient noise, such as other species'  
336 vocalisations, which could potentially confuse the classifier. This method of 'hand picking' the  
337 negative class could reduce the false negative rate while also reducing the false positive rate by  
338 negatively labelling potentially confusing information. The low false positive rate reported in  
339 this study, as well as the low false negative rate for "clear" calls, suggests that our method of

340 randomly selecting the negative class was sufficient for our aims. Where reducing the false  
341 negative rate has greater importance, such as identifying infrequent vocalisations, a more  
342 thorough approach could be preferable. Either way, it is important to recognise the trade-off  
343 between ensuring that the negative class contains as little erroneous information as possible  
344 and the time required to construct an adequate training dataset.

345           It is important to acknowledge that the distinction of great calls into “clear”, “faint” and  
346 “very faint” categories was not based on exact measurements. These categories cannot be  
347 translated directly to distance; however, in most cases, it is likely that “clear” calls were  
348 recorded from gibbons singing closer to the ARUs. Overall, our best performing model  
349 identified 98% of all “clear” calls, missing only 6 from the manually-annotated dataset. This was  
350 comparable to the human observers, who also missed 6 “clear” calls. The model performed  
351 worse than the human detector at detecting “faint” and “very faint” calls, however, picking up  
352 73% and 44% of the manually-annotated instances respectively. With this in mind, the results  
353 suggest that detection likelihood is affected by the caller-ARU distance (Spillmann et al. 2015).  
354 This supports the suggestion by Jahn et al. (2017) that the difference in recall between an  
355 automated detector and a human listener is caused by the former having a smaller detection  
356 radius. Future studies should aim to apply relationships between signal strength and the  
357 distance of the source from the receiver to estimate call detection probability over distance.  
358 This will help to determine the effective area being sampled by PAM studies.

359           Although the model did not have perfect recall, the importance of detecting all great  
360 calls within a recording will depend on the research question. With regards to calling activity  
361 over time, the distribution of calling frequency detected by the model was not significantly

362 different to that of the manual annotations. Therefore, our model can be used to estimate  
363 spatiotemporal variations in *H. albibarbis* calling activity. Through analysing recordings from  
364 multiple habitats, our model could be applied to understand the relative importance of forest  
365 subtypes for the species, as singing behaviour is density dependent (Cheyne et al. 2008), with  
366 less singing activity at lower group densities. This could operate on a continuous long-term time  
367 frame which would be hard to achieve when relying on human observers in the field.

368 For an in-depth understanding of gibbon group abundance and density, further  
369 information is necessary. One method is to use individuals as the sampling unit (Buckland 2006)  
370 by analysing their call structure (Clink et al. 2023), or by localising vocalisations using estimates  
371 of direction to the source from multiple ARUs (Stevenson et al. 2015). This may prove highly  
372 effective at estimating gibbon abundance and density over the short term yet could prove too  
373 complex over the course of hundreds, or thousands of hours of audio. An alternative method is  
374 to estimate vocalisation density per unit time, apply an estimation of vocalisation rate, and  
375 then convert vocalisation density into group density (Marques et al. 2013). This does require a  
376 knowledge of the area covered by each ARU, yet, for this task, Marques et al. (2013) notes that  
377 automated detectors need not perform extraordinarily well so long as true positive and false  
378 positive rates are characterised accurately. This method does not require for the effective area  
379 of ARUs to overlap, so in theory could monitor a larger area with the same resources.

380 While post-processing greatly reduced the number of repeated detections for single  
381 great call events, these would still account for many false positives if only one true positive per  
382 great call is allowed. The post-processing protocol did not seek to judicate between detections  
383 if there was a tie in the highest scoring instances within a group, as in some cases this

384 represented two great call events close in time. In fact, 22 false negatives derived from  
385 instances when the grouping of detections failed to take this into account, and limiting the  
386 number of detections per group to one would have increased the false negative rate further. An  
387 alternative approach would be an improved post-processing stage to better interpret the  
388 output of the model. In audio classification, CNNs inspect audio recordings as image-like  
389 segments and so are unable to use broader-scale contextual information, such as whether the  
390 current point in a recording is preceded by a target call (Wang et al. 2022). It has therefore  
391 been proposed to combine CNNs with other machine learning techniques, such as Hidden  
392 Markov Models (HMMs) or deep learning techniques, such as Recurrent Neural Networks  
393 (RNNs). Postprocessing the output of a CNN with a HMM or an RNN has been shown to improve  
394 F score (Madhusudhana et al. 2021; Wang et al. 2022). For instance, Wang et al. (2022) showed  
395 that the application of a combined Convolutional Recurrent Neural Network (CRNN) reduced  
396 error rates arising from overestimation of gibbon calls (49-54%) to 0.5%. It was noted, however,  
397 that while post-processing the model output with a HMM performed second best, it required  
398 much less computational power. Future work should therefore consider both of these options  
399 when attempting to improve on our post-processing method.

400 A key aim of this study was to help address the analysis bottleneck evident in PAM data  
401 by improving on the time taken to manually analyse recordings. During the batch processing  
402 stage, it took 19 seconds for our trained model to process each hour of test recordings. This  
403 greatly improved on a human processing speed of minimum 1 hour per hour of audio for this  
404 study, varying depending on the level of observer experience. One caveat is the significant  
405 amount of time required to construct a manually-annotated dataset to train and test a CNN

406 when compared to other machine learning approaches (Stowell 2022). Despite our study  
407 showing how data augmentation can be effective where training data is limited, careful  
408 consideration should be taken when under time-pressure if no training datasets are already  
409 available. In some cases, it may be better to adopt an approach that requires less data to  
410 develop, such as an SVM or a GMM (Clink et al. 2023).

411 Stowell (2022) noted that it was increasingly common to evaluate deep learning models  
412 on test sets specifically designed to differ in some respects from the training data, such as  
413 location, signal to noise ratio, or by season. In this case, the model was designed with  
414 application to the MBEF bioacoustic dataset in mind, and so it is appropriate that the test data  
415 came from the same location as the training data. Test inputs spanned across all times of the  
416 year and were recorded from three different habitats. This ensured a variety of potential sound  
417 environments were included in the testing stage. However, for applications in other locations,  
418 especially outside the Rungan Forest Landscape, it would be advantageous to first test the  
419 model on recordings captured elsewhere within *H. albibarbis*' range.

## 420 V. CONCLUSION

421 Our study demonstrates how an open-source deep-learning framework can be adapted  
422 to produce a CNN capable of detecting *H. albibarbis* great calls, performing at a comparative  
423 level to similar CNN approaches (see Table 1). Our model performed best on the highest-quality  
424 calls and yielded a low false-positive rate, meeting the objectives of this study. There was a  
425 much lower likelihood of successful detection for the lowest-quality calls, however, and future

426 studies should aim to estimate call detection probability over distance to determine the  
427 effective area being sampled.

428 Further development of the post-processing stage could help to reduce the number of  
429 repeat detections for each call. However, the current output can be used to estimate calling  
430 rate over time. Further work should seek to apply this model to long-term acoustic datasets  
431 over a variety of habitats to study spatial and temporal variation in gibbon calling activity.  
432 Furthermore, in combination with future studies on sound propagation of gibbon vocalisations,  
433 this represents an opportunity to monitor *H. albibarbis*' populations on an ever-greater  
434 spatiotemporal scale. Our work presents some key considerations to inform decision-making  
435 for such projects, and a full workflow script to visualise how these can be implemented in  
436 developing an automated detector.

437 **SUPPLEMENTARY MATERIAL**

438 For the full Koogu workflow, see Supplementary Material A (SupPub1.txt). For the post-  
439 processing script, see Supplementary Material B (SupPub2.txt).

440 **ACKNOWLEDGEMENTS**

441 We thank Universitas Muhammadiyah Palangkaraya and the Borneo Nature Foundation for  
442 access and support at the Mungku Baru Education Forest. We thank the Indonesian  
443 government for permission to carry out this research (RISTEK foreign research permits:  
444 189/SIP/FRP/E5/Dit.KI/VI/2018 and 43/E5/E5.4/SIP.EXT/2019 [Wendy M. Erb], BRIN foreign  
445 research permits: 223/SIP/IV/FR/5/2023 [F. J. F. van Veen] and 225/SIP/IV/FR/5/2023 [A. F.  
446 Owens]). We thank Erik Estrada and Rido for their invaluable support in deploying and  
447 maintaining the ARUs and data in the field. We also thank Georgia Allen, Amy Barron, Sophie  
448 Carpenter, and Elena Gough for their contributions in creating the manually-annotated dataset.  
449 This project has been a collaboration between international and Indonesian researchers from  
450 the start and this is recognised in joint authorship here. In addition, the Indigenous Dayak Ngaju  
451 community of Mungku Baru have been closely involved in the field research, sharing their  
452 knowledge, and discussing the implications with researchers and with local conservation  
453 practitioners from Yayasan Borneo Nature.

454 AFO, FJFvV, WME, MAI and KH developed the project and concept; MAN, TMS and SMS  
455 provided permissions and guidance for field study design and execution; AFO, SM and MS  
456 carried out building and analysis of the model; AFO drafted the manuscript to which all other  
457 authors then contributed to produce the final version.

458 **AUTHOUR DECLARATIONS**

459 The authors declare no competing interests.

460 **DATA AVAILABILITY**

461 [Recordings of great calls for both the training and testing datasets will be uploaded to a public  
462 repository that issues datasets with DOIs before publication.]

463 **ETHICS APPROVAL**

464 Ethical approval was provided by the University of Exeter (application ID: 1845574), BRIN  
465 (application number: 22022023000026) and Rutgers University (IACUC ID: PROTO201800073).

466 **APPENDIX**

467 i. Song Meter SM4 default settings: sensitivity of  $-35 \pm 4$  dB (0 dB = 1V/pa@1kHz), dynamic  
468 range of 14 to 100 dB SPL at 0 dB gain, microphone gain of 16 dB, and inbuilt preamplifier gain  
469 of 26 dB.

470 ii. Spectrograms were loaded into Raven Pro 1.6 for manual annotation. These were generated  
471 using a 3462-sample Hann window with a 90% overlap and a 4096-sample Discrete Fourier  
472 Transformation.

473 iii. Spectrograms were computed to provide inputs to the CNN model. These had an analysis  
474 window of 0.192s and a 75% overlap. The bandwidth was also restricted to between 200 and  
475 2200 Hz to exclude noise outside of the target frequency range.

476 **REFERENCES**

- 477 Acevedo, M. and Villanueva-rivera, L. 2010. From the Field: Using Automated Digital  
478 Recording Systems as Effective Tools for the Monitoring of Birds and Amphibians. *Wildlife*  
479 *Society Bulletin* 34(1), pp. 211–214.
- 480 Mac Aodha, O. et al. 2018. Bat detective—Deep learning tools for bat acoustic signal  
481 detection. *PLOS Computational Biology* 14(3), p. e1005995. doi: 10.1371/journal.pcbi.1005995.
- 482 Bender, D. J., Contreras, T. A. and Fahrig, L. 1998. Habitat loss and population decline:  
483 a meta-analysis of the patch size effect. *Ecology* 79(2), pp. 517–533.
- 484 Bengio, Y., Goodfellow, I. and Courville, A. 2016. *Deep Learning*. Cambridge,  
485 Massachusetts: MIT Press.
- 486 Bjorck, J., Rappazzo, B.H., Chen, D., Bernstein, R., Wrege, P.H. and Gomes, C.P. 2019.  
487 Automatic Detection and Compression for Passive Acoustic Monitoring of the African Forest  
488 Elephant. *Proceedings of the AAAI Conference on Artificial Intelligence* 33(01), pp. 476–484.  
489 doi: 10.1609/aaai.v33i01.3301476.
- 490 Borowiec, M.L., Dikow, R.B., Frandsen, P.B., McKeeken, A., Valentini, G. and White,  
491 A.E. 2022. Deep learning as a tool for ecology and evolution. *Methods in Ecology and Evolution*  
492 13(8), pp. 1640–1660. doi: 10.1111/2041-210X.13901.
- 493 Brauer, C.L., Donovan, T.M., Mickey, R.M., Katz, J. and Mitchell, B.R. 2016. A  
494 comparison of acoustic monitoring methods for common anurans of the northeastern United  
495 States. *Wildlife Society Bulletin* 40(1), pp. 140–149. doi: 10.1002/wsb.619.
- 496 Buckland, S.T. 2006. Point-Transect Surveys for Songbirds: Robust Methodologies. *The*  
497 *Auk* 123(2), pp. 345–357. doi: 10.1093/auk/123.2.345.

498 Buckley, B. et al. 2018. Biodiversity, Forest Structure & Conservation Importance of the  
499 Mungku Baru Education Forest, Rungan, Central Kalimantan. *Borneo Nature Foundation*.

500 Cheyne, S. et al. 2016. Population mapping of gibbons in Kalimantan, Indonesia:  
501 correlates of gibbon density and vegetation across the species' range. *Endangered Species*  
502 *Research* 30, pp. 133–143. doi: 10.3354/esr00734.

503 Cheyne, S.M., Thompson, C.J.H., Phillips, A.C., Hill, R.M.C. and Limin, S.H. 2008.  
504 Density and population estimate of gibbons (*Hylobates albibarbis*) in the Sabangau catchment,  
505 Central Kalimantan, Indonesia. *Primates* 49(1), pp. 50–56. doi: 10.1007/s10329-007-0063-0.

506 Clink, D.J., Kier, I., Ahmad, A.H. and Klinck, H. 2023. A workflow for the automated  
507 detection and classification of female gibbon calls from long-term acoustic recordings. *Frontiers*  
508 *in Ecology and Evolution* 11. doi: 10.3389/fevo.2023.1071640.

509 Colonna, J., Peet, T., Ferreira, C.A., Jorge, A.M., Gomes, E.F. and Gama, J. 2016.  
510 Automatic Classification of Anuran Sounds Using Convolutional Neural Networks. In:  
511 *Proceedings of the Ninth International C\* Conference on Computer Science & Software*  
512 *Engineering - C3S2E '16*. New York, New York, USA: ACM Press, pp. 73–78. doi:  
513 10.1145/2948992.2949016.

514 Drake, K.L., Frey, M., Hogan, D. and Hedley, R. 2016. Using digital recordings and  
515 sonogram analysis to obtain counts of yellow rails. *Wildlife Society Bulletin* 40(2), pp. 346–354.  
516 doi: 10.1002/wsb.658.

517 Dufourq, E. et al. 2021. Automated detection of Hainan gibbon calls for passive acoustic  
518 monitoring. *Remote Sensing in Ecology and Conservation* 7(3), pp. 475–487. doi:  
519 10.1002/rse2.201.

- 520 Guyot, P., Alix, F., Guerin, T., Lambeaux, E. and Rotureau, A. 2021. Fish migration  
521 monitoring from audio detection with CNNs. In: *Audio Mostly 2021*. New York, NY, USA: ACM,  
522 pp. 244–247. doi: 10.1145/3478384.3478393.
- 523 Heinicke, S., Kalan, A.K., Wagner, O.J.J., Mundry, R., Lukashevich, H. and Kühl, H.S.  
524 2015. Assessing the performance of a semi-automated acoustic monitoring system for primates.  
525 *Methods in Ecology and Evolution* 6(7), pp. 753–763. doi: 10.1111/2041-210X.12384.
- 526 Hibino, S., Suzuki, C. and Nishino, T. 2021. Classification of singing insect sounds with  
527 convolutional neural network. *Acoustical Science and Technology* 42(6), p. E2152. doi:  
528 10.1250/ast.42.354.
- 529 Huang, G., Liu, Z., van der Maaten, L. and Weinberger, K.Q. 2016. Densely Connected  
530 Convolutional Networks.
- 531 Jahn, O., Ganchev, T.D., Marques, M.I. and Schuchmann, K.-L. 2017. Automated Sound  
532 Recognition Provides Insights into the Behavioral Ecology of a Tropical Bird. *PLOS ONE* 12(1),  
533 p. e0169041. doi: 10.1371/journal.pone.0169041.
- 534 Jennings, N., Parsons, S. and Pocock, M.J.O. 2008. Human vs. machine: identification  
535 of bat species from their echolocation calls by humans and by artificial neural networks.  
536 *Canadian Journal of Zoology* 86(5), pp. 371–377. doi: 10.1139/Z08-009.
- 537 Jetz, W. et al. 2019. Essential biodiversity variables for mapping and monitoring species  
538 populations. *Nature Ecology & Evolution* 3(4), pp. 539–551. doi: 10.1038/s41559-019-0826-1.
- 539 Josh Patterson and Adam Gibson. 2017. *Deep Learning*. Loukides, M. and McGovern,  
540 T. eds. O'Reilly Media.
- 541 Kingma, D.P. and Ba, J. 2014. Adam: A Method for Stochastic Optimization.

- 542 Knight, E.C., Hannah, K.C., Foley, G.J., Scott, C.D., Brigham, R.M. and Bayne, E. 2017.  
543 Recommendations for acoustic recognizer performance assessment with application to five  
544 common automated signal recognition programs. *Avian Conservation and Ecology* 12(2), p.  
545 art14. doi: 10.5751/ACE-01114-120214.
- 546 Kong, Q., Xu, Y. and Plumbley, M.D. 2017. Joint detection and classification  
547 convolutional neural network on weakly labelled bird audio detection. In: *2017 25th European*  
548 *Signal Processing Conference (EUSIPCO)*. IEEE, pp. 1749–1753. doi:  
549 10.23919/EUSIPCO.2017.8081509.
- 550 Madhusudhana. 2023. shyamblast/Koogu: v0.7.2.
- 551 Madhusudhana, S. et al. 2021. Improve automatic detection of animal call sequences  
552 with temporal context. *Journal of The Royal Society Interface* 18(180), p. 20210297. doi:  
553 10.1098/rsif.2021.0297.
- 554 Marques, T.A. et al. 2013. Estimating animal population density using passive acoustics.  
555 *Biological Reviews* 88(2), pp. 287–309. doi: 10.1111/brv.12001.
- 556 Mielke, A. and Zuberbühler, K. 2013. A method for automated individual, species and  
557 call type recognition in free-ranging animals. *Animal Behaviour* 86(2), pp. 475–482. doi:  
558 10.1016/j.anbehav.2013.04.017.
- 559 Miller, B.S., Madhusudhana, S., Aulich, M.G. and Kelly, N. 2023. Deep learning  
560 algorithm outperforms experienced human observer at detection of blue whale D-calls: a  
561 double-observer analysis. *Remote Sensing in Ecology and Conservation* 9(1), pp. 104–116. doi:  
562 10.1002/rse2.297.

563 Morgan, M.M. and Braasch, J. 2021. Long-term deep learning-facilitated environmental  
564 acoustic monitoring in the Capital Region of New York State. *Ecological Informatics* 61, p.  
565 101242. doi: 10.1016/j.ecoinf.2021.101242.

566 Narasimhan, R., Fern, X.Z. and Raich, R. 2017. Simultaneous segmentation and  
567 classification of bird song using CNN. In: *2017 IEEE International Conference on Acoustics,  
568 Speech and Signal Processing (ICASSP)*. IEEE, pp. 146–150. doi:  
569 10.1109/ICASSP.2017.7952135.

570 Noda, J.J., Travieso, C.M., Sanchez-Rodriguez, D., Dutta, M.K. and Singh, A. 2016.  
571 Using bioacoustic signals and Support Vector Machine for automatic classification of insects. In:  
572 *2016 3rd International Conference on Signal Processing and Integrated Networks (SPIN)*. IEEE,  
573 pp. 656–659. doi: 10.1109/SPIN.2016.7566778.

574 Norouzzadeh, M.S., Nguyen, A., Kosmala, M., Swanson, A., Palmer, M.S., Packer, C.  
575 and Clune, J. 2018. Automatically identifying, counting, and describing wild animals in camera-  
576 trap images with deep learning. *Proceedings of the National Academy of Sciences* 115(25). doi:  
577 10.1073/pnas.1719367115.

578 Piel, A.K. and Wich, S.A. 2021. *Conservation Technology*. Wich, S. A. and Piel, A. K.  
579 eds. Oxford University PressOxford. doi: 10.1093/oso/9780198850243.001.0001.

580 Purnama, A. and Afitah, I. 2021. Motivasi Masyarakat Terhadap Pengelolaan Khdtk  
581 Mungku Baru, Palangka Raya. *Anterior Jurnal* 20(2), pp. 43–49. doi:  
582 10.33084/anterior.v20i2.2162.

583 Rammer, W. and Seidl, R. 2019. Harnessing Deep Learning in Ecology: An Example  
584 Predicting Bark Beetle Outbreaks. *Frontiers in Plant Science* 10. doi: 10.3389/fpls.2019.01327.

585 Schlüter, J. and Grill, T. 2015. Exploring Data Augmentation for Improved Singing Voice  
586 Detection with Neural Networks. In: , *16th International Society for Music Information Retrieval*  
587 *Conference*.

588 Simon S. Haykin. 2009. *Neural Networks: A Comprehensive Foundation*. 3rd ed.  
589 McMaster University, Canada: Pearson Education.

590 Spillmann, B., van Noordwijk, M.A., Willems, E.P., Mitra Setia, T., Wipfli, U. and van  
591 Schaik, C.P. 2015. Validation of an acoustic location system to monitor Bornean orangutan (  
592 *Pongo pygmaeus wurmbii*) long calls. *American Journal of Primatology* 77(7), pp. 767–776. doi:  
593 10.1002/ajp.22398.

594 Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I. and Salakhutdinov, R. 2014.  
595 Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine*  
596 *Learning Research* 15, pp. 1929–1958.

597 Stevenson, B.C., Borchers, D.L., Altwegg, R., Swift, R.J., Gillespie, D.M. and Measey,  
598 G.J. 2015. A general framework for animal density estimation from acoustic detections across a  
599 fixed microphone array. *Methods in Ecology and Evolution* 6(1), pp. 38–48. doi: 10.1111/2041-  
600 210X.12291.

601 Stowell, D. 2022. Computational bioacoustics with deep learning: a review and roadmap.  
602 *PeerJ* 10, p. e13152. doi: 10.7717/peerj.13152.

603 Sugai, L.S.M., Silva, T.S.F., Ribeiro, J.W. and Llusia, D. 2019. Terrestrial Passive  
604 Acoustic Monitoring: Review and Perspectives. *BioScience* 69(1), pp. 15–25. doi:  
605 10.1093/biosci/biy147.

606 Verma, A., van der Wal, R. and Fischer, A. 2016. Imagining wildlife: New technologies  
607 and animal censuses, maps and museums. *Geoforum* 75, pp. 75–86. doi:  
608 10.1016/j.geoforum.2016.07.002.

609 Vu, T.T. and Tran, L.M. 2019. An Application of Autonomous Recorders for Gibbon  
610 Monitoring. *International Journal of Primatology* 40(2), pp. 169–186. doi: 10.1007/s10764-018-  
611 0073-3.

612 Wang, Y., Ye, J. and Borchers, D.L. 2022. Automated call detection for acoustic surveys  
613 with structured calls of varying length. *Methods in Ecology and Evolution* 13(7), pp. 1552–1567.  
614 doi: 10.1111/2041-210X.13873.

615 Zhou, X. et al. 2023. Methods for processing and analyzing passive acoustic monitoring  
616 data: An example of song recognition in western black-crested gibbons. *Ecological Indicators*  
617 155, p. 110908. doi: 10.1016/j.ecolind.2023.110908.