







ORIGINAL RESEARCH

Automated detection of Hainan gibbon calls for passive acoustic monitoring

Emmanuel Dufourq^{1,2} , Ian Durbach^{3,4,1} , James P. Hansford^{5,6} , Amanda Hoepfner⁷, Heidi Ma⁵ , Jessica V. Bryant⁸ , Christina S. Stender⁹, Wenying Li¹⁰, Zhiwei Liu¹⁰, Qing Chen¹⁰, Zhaoli Zhou¹⁰ & Samuel T. Turvey⁵ 

¹African Institute for Mathematical Sciences, Muizenberg, South Africa

²Stellenbosch University, Stellenbosch, South Africa

³Centre for Research into Ecological and Environmental Modelling, University of St Andrews, St Andrews, UK

⁴Centre for Statistics in Ecology, the Environment, and Conservation, University of Cape Town, Rondebosch, South Africa

⁵Institute of Zoology, Zoological Society of London, Regent's Park, London NW1 4RY, UK

⁶Department of Biological Sciences, Northern Illinois University, DeKalb Illinois, 60115,

⁷School of Biological Sciences, University of Utah, Salt Lake City Utah, 84112,

⁸Department of Life Sciences, University of Roehampton, London SW15 4JD, UK

⁹Living Collections, Zoological Society of London, Regent's Park, London NW1 4RY, UK

¹⁰Bawangling National Nature Reserve, Changjiang Li Autonomous County, Hainan, China

Keywords

Bioacoustics, convolutional neural networks, deep learning, Hainan gibbons, passive acoustic monitoring, species identification

Correspondence

Emmanuel Dufourq, African Institute for Mathematical Sciences, Muizenberg, South Africa. Tel: +27 21 787 9320; Fax: +27 21 7879321; E-mail: edufourq@gmail.com

Editor: Nathalie Pettorelli

Associate Editor: Christos Astaras

Received: 29 October 2020; Revised: 22 January 2021; Accepted: 4 March 2021

doi: 10.1002/rse.2.201

Abstract

Extracting species calls from passive acoustic recordings is a common preliminary step to ecological analysis. For many species, particularly those occupying noisy, acoustically variable habitats, the call extraction process continues to be largely manual, a time-consuming and increasingly unsustainable process. Deep neural networks have been shown to offer excellent performance across a range of acoustic classification applications, but are relatively underused in ecology. We describe the steps involved in developing an automated classifier for a passive acoustic monitoring project, using the identification of calls of the Hainan gibbon *Nomascus hainanus*, one of the world's rarest mammal species, as a case study. This includes preprocessing—selecting a temporal resolution, windowing and annotation; data augmentation; processing—choosing and fitting appropriate neural network models; and post-processing—linking model predictions to replace, or more likely facilitate, manual labelling. Our best model converted acoustic recordings into spectrogram images on the mel frequency scale, using these to train a convolutional neural network. Model predictions were highly accurate, with per-second false positive and false negative rates of 1.5% and 22.3%. Nearly all false negatives were at the fringes of calls, adjacent to segments where the call was correctly identified, so that very few calls were missed altogether. A post-processing step identifying intervals of repeated calling reduced an 8-h recording to, on average, 22 min for manual processing, and did not miss any calling bouts over 72 h of test recordings. Gibbon calling bouts were detected regularly in multi-month recordings from all selected survey points within Bawangling National Nature Reserve, Hainan. We demonstrate that passive acoustic monitoring incorporating an automated classifier represents an effective tool for remote detection of one of the world's rarest and most threatened species. Our study highlights the viability of using neural networks to automate or greatly assist the manual labelling of data collected by passive acoustic monitoring projects. We emphasize that model development and implementation be informed and guided by ecological objectives, and increase accessibility of these tools with a series of notebooks that allow users to build and deploy their own acoustic classifiers.

Introduction

Deep learning holds enormous promise for automating the labelling of bioacoustic data. The number of applications is growing (Christin et al., 2019), but the majority of datasets are still labelled manually (Fairbrass et al., 2019; Kiskin et al., 2020; Pamula et al., 2019), even as the rate of data collection makes this approach increasingly unsustainable. The mismatch between the potential of deep learning approaches and their actual uptake among practitioners occurs because getting models to perform as well as an experienced human is difficult. Human-like performance usually requires substantial amounts of training data or relatively stable background environments, conditions that are often absent in ecological applications. Model tuning and data manipulation is often required, and while guidelines are emerging (Patterson & Gibson, 2017; Stowell et al., 2019b), these can, with some justification, appear subjective and case specific. A lack of computing resources and user-friendly software can also be a barrier to entry. Case studies reporting successful applications play an important role in developing and disseminating best practices, and in discriminating between those tasks that current deep learning methods are able to automate and those they cannot. Previous applications have used convolutional neural networks (CNNs; LeCun et al. (2015)) to identify various bird (Grill & Schlüter, 2017; Kahl et al., 2017; Stowell et al., 2019b) and whale species (Bergler et al., 2019; Bermant et al., 2019; Jiang et al., 2019; Shiu et al., 2020), bees (Kulyukin et al., 2018; Nolasco et al., 2019), as well as anomalous acoustic events in soundscapes (Sethi et al., 2020). These have shown, for example, that a generally good approach is to represent data as spectrograms and treat the problem as an image classification one, as well as providing specialized approaches for data augmentation on spectrogram inputs, such as pitch and time shifting and introducing background noise (Bergler et al., 2019; Sprengel et al., 2016).

Despite this, no studies report the process of applying deep learning within the scope of a typical acoustic monitoring project designed to answer a well-defined research question. Most applications are either smaller – using data collected for the purpose of testing a deep learning approach, and often written for a machine learning rather than ecological audience (e.g. Kiskin et al., 2020; Kulyukin et al., 2018); or larger – aggregating datasets across several independent studies to investigate if models generalize (Bergler et al., 2019; Shiu et al., 2020; Stowell et al., 2019b) – than most monitoring projects. In this paper we address this gap, describing the development of a classifier for identifying Hainan gibbon *Nomascus hainanus* calls in

passive acoustic recordings collected as part of a long-term monitoring project, with the aim of providing practitioners with a realistic and relatable idea of the process, and modelling choices, involved, as well as guidelines for these choices.

The Hainan gibbon is the world's rarest primate and one of the world's rarest mammals, with only a single population of about 30 individuals surviving in Bawangling National Nature Reserve (BNNR), Hainan, China (Chan et al., 2005; Liu et al., 2020; Turvey et al., 2015). Improved monitoring of this population using novel methods, to understand factors affecting successful dispersal, breeding group formation and colonization of new habitat, has been identified as an urgent short-term conservation goal for the species (Turvey et al., 2015; Zhang et al., 2020). Gibbons call regularly to advertise territory and maintain group cohesiveness against rivals, using a complex structure consisting of short individual vocal syllables or 'notes' of ca. 0.2–2.75 s assembled together into longer 'phrases' consisting of one to six notes, which are themselves organized into 'songs' of several minutes (Deng et al., 2014). Gibbon population surveys are usually conducted by detecting this daily song using a fixed-point count survey method, whereby researchers listen opportunistically for calls at elevated listening posts (Brockelman & Srikosamatara, 1993; Kidney et al., 2016). However, this traditional monitoring approach is labour intensive and is only conducted for discrete survey periods. Gibbons are therefore prime candidates for passive acoustic monitoring and recent studies have used data collected in this way to model occupancy (Vu & Tran, 2019) and to discriminate between individuals using spectral features (Clink et al., 2019; Zhou et al., 2019). All of these studies, however, have relied on an initial manual extraction of calls.

In order to develop a continuous monitoring protocol for Hainan gibbons we conducted long-term passive acoustic monitoring and developed an automated classifier able to identify whether gibbons were calling in the vicinity of a particular recorder, with the aim of establishing whether the area proximal to the recorder was occupied that day. It was therefore important to be able to detect individual gibbon calling bouts, but not necessarily to be able to discriminate every phrase made during the bout. We address issues that are important to the overall usefulness of a classifier, including deciding how much data to manually label, data augmentation, operationally meaningful definitions of classifier success and the development of user-friendly software. Our study provides an effective new monitoring method for the world's rarest primate, and also has wider applicability for applying deep learning to develop passive acoustic monitoring frameworks for other conservation-priority loud-call

species such as cetaceans, elephants or other primates (Crunchant et al., 2020).

Materials and Methods

Data collection

Eight Song Meter SM3 recorders (Wildlife Acoustics, Maynard, Massachusetts) were used to collect acoustic data from 1 March to 20 August 2016 within BNNR. Recorders were attached to trees at a height of approximately 1.5 m in tropical evergreen forest. Four recorders were situated within the known home ranges of the four Hainan gibbon social groups existing during the study period (Groups A–D; see Bryant et al. (2017)), three were situated at locations intermediate between known home ranges, and a further recorder was placed in an area where a solitary male gibbon was thought to occur (Bryant et al., 2016). They were placed at locations that were used as regular listening posts for monitoring gibbons by reserve staff (Fig. 1). The peak Hainan gibbon calling period is 06:00–07:00, with calling continuing at decreasing regularity for several hours (Chan et al., 2005). Recorders were therefore set to record for 8 h each day from the time of sunrise, which varied between approximately 05:00 and 06:00 during the study period. Memory cards and batteries were changed every 40 days. Devices did not record continuously throughout the entire survey period due to logistical and technical issues; in total, survey days per recorder varied between 79 and 129 days, and roughly 6000 h of recordings were collected. The majority of recordings were made with a sampling rate of 9600 Hz and bit depth of 16, with isolated recordings at 28 800 Hz.

Data analysis

We manually labelled 32 8-h recordings by inspecting spectrograms and listening to audio using Sonic Visualiser (Cannam et al., 2010), and end times, and the number of notes, of each observed gibbon phrase. Four files containing no gibbon calls were discarded, as periods without gibbon calls are readily available from the remaining 28 files. This process yielded 1246 gibbon phrases.

To construct the fixed-length inputs required by CNNs, we divided each 8-h recording into segments with window length 10 s and hop length 1 s (starting times of consecutive 10 s segments differ by 1 s, Fig. 2). This window length was chosen so that even the longest phrase (8 s, Supplementary Material A) fits within a single segment; using a slightly longer segment length allows for potentially longer unseen phrases, and results in more

positive segments after windowing. All audio was converted into mono, as done in various applications (e.g. Bergler et al., 2019; Qazi et al., 2018; Stowell et al., 2019a). By cross-referencing the time intervals of each segment with the logged start and end times of known gibbon phrases, each segment was labelled as (a) a ‘presence’, if its time interval completely contained the interval of at least one labelled phrase, (b) an ‘absence’, if its time interval contained no part of any phrase or (c) a ‘partial presence’, if its time interval intersected but did not completely contain the interval of at least one labelled phrase (Fig. 2). Partial presences were excluded from further analysis.

Each recording was downsampled to 4800 Hz, so that the Nyquist rate was higher than the maximum frequency of Hainan gibbon calls (2000 Hz). No anti-aliasing was performed although, because we downsampled entire recordings, we would expect any artefacts to be unrelated to the presence of gibbon calls. This was confirmed by a post hoc comparison of aliased and non-aliased versions of a 5% sample of segments. The downsampled inputs – each segment a time series of 48 000 sample points – used as inputs to the 1-D CNNs described in the next section. In addition, we converted each audio segment into a mel-scale spectrogram (Bergler et al., 2019; Huang et al., 2001), to be used as an input image to a 2-D CNN, using a Hann analysis window size of 1024 samples (213 ms), a hop size of 256 samples (53.3 ms, 75% overlap) and 128 mel frequency bins with centres uniformly spaced between 1 and 2k Hz, a conservative interval following Deng et al. (2014) and our own exploratory analyses. This results in 188 time steps by 128 frequency bands. These were computed using the Librosa library. These values were chosen on the basis of preliminary investigations, although the results were not particularly sensitive to these choices. The spectrogram images had a size of 128 × 188 pixels; larger image sizes can capture greater detail but typically require more network parameters and computation time to do so.

After processing, our dataset consisted of 5285 segments containing at least one complete phrase. While the vast majority of segments do not contain any gibbon calls, we restricted the number of absence segments to the same number as presences, to avoid a large class imbalance. Absence segments were initially collected by randomly sampling, but we found that better results were obtained by specifically including absence segments that contained typical ambient noise, such as bird calls, rain events and other background noises that could potentially confuse the classifier (Stowell et al., 2019a). Extracting these required additional manual processing of the audio data.

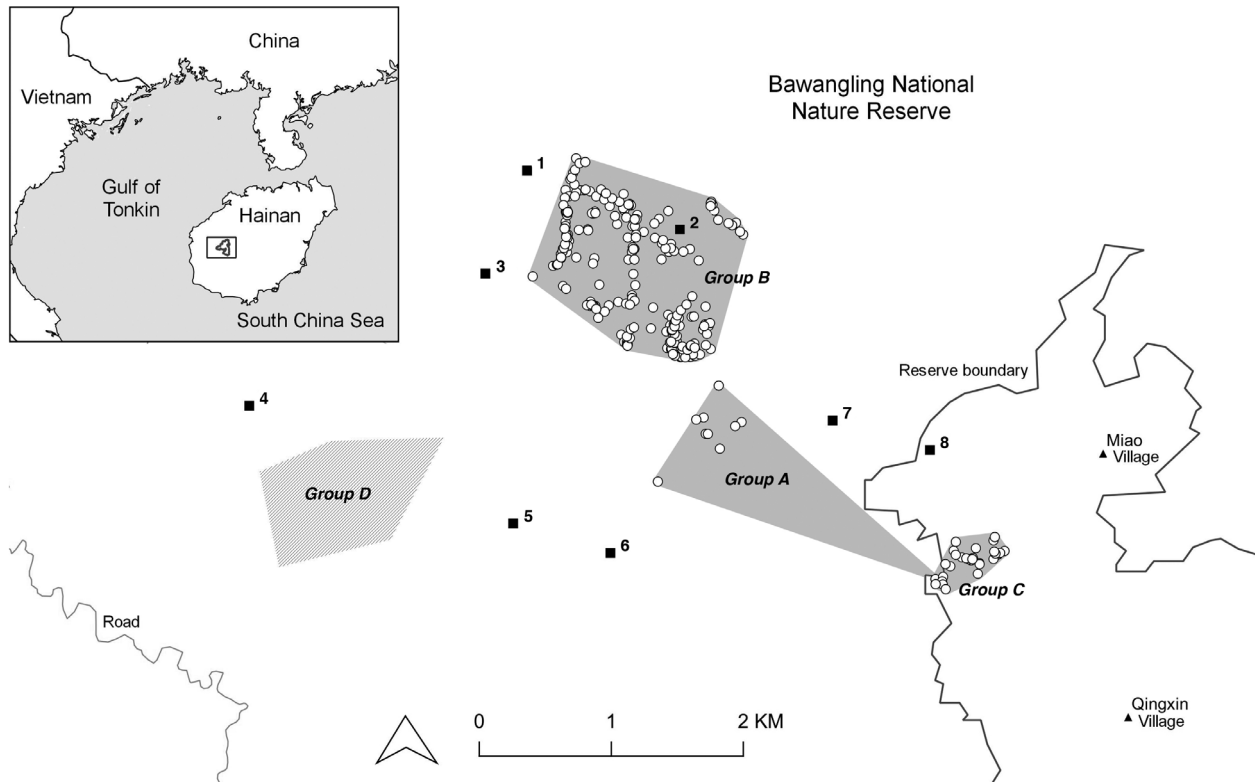


Figure 1. Locations of eight Song Meter SM3 recorders (labelled 1–8) used to detect gibbons in 2016 within Bawangling National Nature Reserve, Hainan, China, in relation to approximate distributions of four Hainan gibbon social groups (A–D). Mapped distributions of groups A–C are based on field data collected in 2010–2011 (see Bryant et al., (2017)); the groups all changed their location slightly between 2011 and 2016, but data on exact group locations in 2016 are unavailable. Approximate location of Group D indicated with hatching based on Bryant et al., (2016).

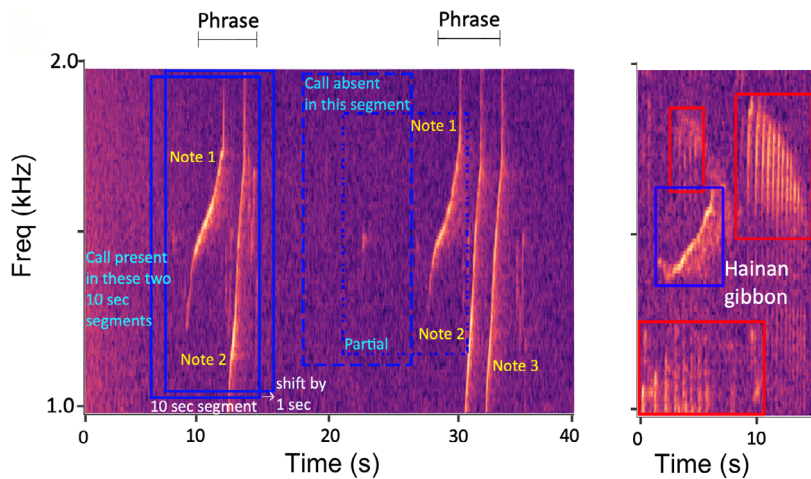


Figure 2. Hainan gibbon calls consist of a sequence of ‘phrases’, each phrase consisting of variable (typically, 1–6) ‘notes’ and often with relatively large intervals between phrases. Left: a two-note phrase followed by a three-note phrase. A single calling bout may last anywhere from a few to dozens of minutes. Our model divides the recording interval into sliding 10 s windows or ‘segments’ (blue boxes), with 80% overlap between adjacent segments. Segments are classified as contained at least one full gibbon phrase (Present; solid line), a partial phrase (Partial; dotted line), or no part of a phrase (Absent; dashed line). Partial presences were excluded from further analysis, creating a two-class audio classification problem. Right: a gibbon phrase partially obscured by noisy background conditions, in this case other species calling (red boxes).

Data augmentation

Data augmentation – boosting sample sizes by adding new samples artificially created by manipulating existing ones, for example using geometric operations like translations and rotation – is commonly used to improve classifier performance, particularly when the training dataset is relatively small (Hestness et al., 2017; Sun et al., 2017). We used data augmentation to create either one or two copies of each 10 s segment in both presence and absence classes. For each presence segment $\mathbf{x}^{(pre)}$, we randomly selected two absence segments, $\mathbf{x}_i^{(abs)}$, $i = 1, 2$. We randomly shifted the starting time of each absence segment forward by $0 < t_i < 9$ s, with the absence segment wrapping back on itself so that it remained 10 s long (Fig. 3 C), to obtain the shifted segment $\mathbf{x}_i^{(shift)}$. Presence segments were not shifted, as this already occurred during the windowing process used to create the original segments. Segments contain amplitude values and thus allow for arithmetic operations to be performed on them. We blended the presence segment with each shifted segment to create augmented presence segments $\mathbf{x}_i^{(aug)} = \alpha \mathbf{x}^{(pre)} + (1 - \alpha) \mathbf{x}_i^{(shift)}$, where α is a mixing parameter, here chosen to be 0.9 (Fig. 3D). We created augmented absence segments using the same approach, that is, combining pairs of absence segments to create a mixture of background scenes.

After augmenting the original segments, we obtained 18 992 segments (9496 presence, 9496 absence) from 19 recordings to train the neural networks. We randomly selected 60% of the data for training (5697 presence, 5697 absence) and used the remaining 40% for validation (3799 presence, 3799 absence). Non-augmented segments from nine separate recordings (2231 presence, 23 689 absence) were kept aside for testing. The files which were used for training and testing were randomly

selected and each file contained at least one presence of a gibbon call.

Neural networks

We considered two kinds of CNN architectures: a 1-D CNN using preprocessed amplitudes of 10 s segments as inputs, and a 2-D CNN that had inputs consisting of spectrogram images constructed from the preprocessed amplitudes. A CNN with a large number of network parameters (e.g. MobileNetV2 (Sandler et al., 2018) which has over 3 million parameters) can result in overfitting – due to degree of freedom given the large number of parameters – if the network is trained on a relatively small number of examples. This observation is often reported in the literature and has also been reported in applications of CNNs in ecology (Chilson et al., 2019). As we had relatively little training data by deep learning standards, we chose these networks as they use simple architectures requiring relatively few parameters. Both 1-D and 2-D CNNs use up to three convolutional layers, each followed by a max pooling layer that reduces the size of the intermediate input passed to the next layer of the network. We used 16×1 and 16×16 convolutional kernels for 1-D and 2-D CNNs respectively. The stack of convolutional layers was followed by one or two dense layers (Fig. 4). The resulting model outputs a detection probability that the input segment (1-D or 2-D) contains at least one complete gibbon phrase.

We chose model hyperparameters using a grid search over the number of convolutional (1, 2, 3) and dense (1, 2, 3) layers, nodes in each of the dense layers (8, 16, 32), filters in each convolutional layer (8, 16, 32), kernel size in each convolutional and max pooling layer (4, 8, 16), and dropout rate (0, 0.2, 0.4, 0.6). Each model was trained for 50 epochs using the Adam optimizer (Kingma

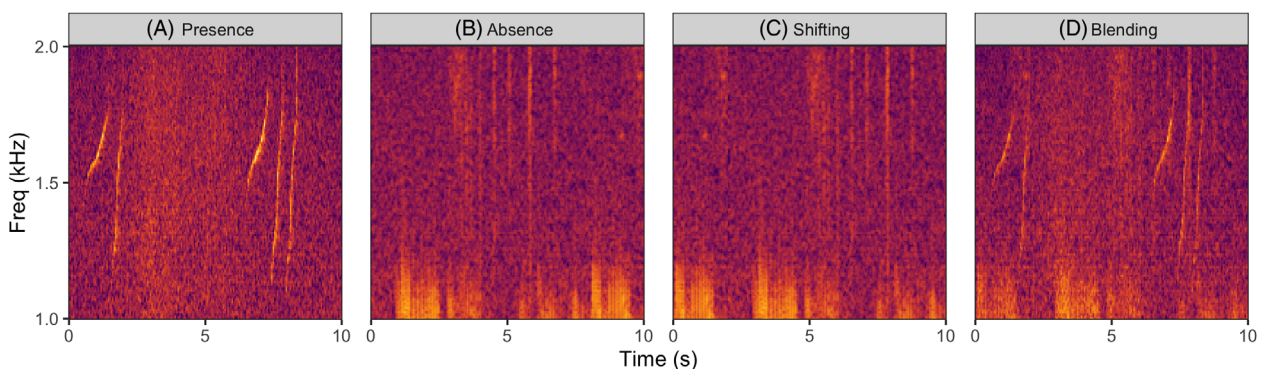


Figure 3. Data augmentation steps involve (a) selecting a presence segment containing a Hainan gibbon phrase, (b) randomly selecting a segment containing only background noise, (c) shifting the starting time of the absence segment forward by a random amount, here 2 s and (d) blending together the presence and shifted absence segments.

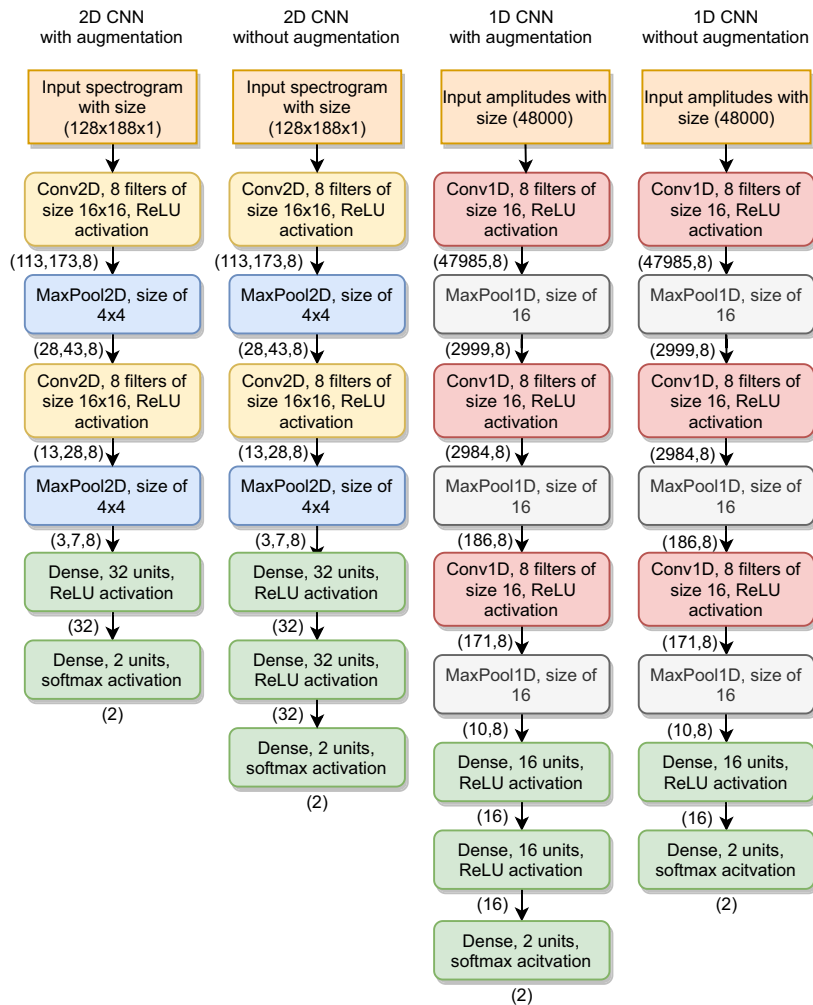


Figure 4. Best architectures for 1-D and 2-D CNNs, for both augmented and non-augmented training datasets. Selected architectures were those with intermediate numbers of free parameters, particularly for 2-D CNNs. The dimensions of the data after each operation is provided in parentheses.

& Ba, 2014) a batch size of eight segments, and a learning rate of 0.001. Models were evaluated based on test set accuracy (proportion of all predictions that were correct), sensitivity (recall) (proportion of true positives divided by positive examples), specificity (proportion of true negatives divided by negative examples), precision (portion of true positives divided by true positives and false positives) and F1-score (harmonic mean between precision and F1-score). Optimal thresholds for converting detection probabilities into binary classifications were those that minimized the ratio of sensitivity and false discovery rate in the validation dataset.

Models were implemented in Python 3 using the TensorFlow (Abadi et al., 2015) library with Keras (Chollet, 2015) for the neural network component, and the Librosa library for audio processing and spectrogram construction

(McFee et al., 2020). Model training and testing was done on a machine running Ubuntu 16.04 LTS with an Intel i7-6700K CPU, 16GB of RAM, and an Nvidia GTX 1070 8GB Graphics Processing Unit. Code and analysis scripts are available online at <https://github.com/emmanueldufourq/GibbonClassifier>.

Post-processing

For an audio recording of arbitrary duration, our approach was to break that recording into overlapping 10 s segments, and to use a trained CNN to output, for each segment starting at second $s = 0, 1, 2, \dots$, a detection probability indicating the likelihood that at least one complete gibbon phrase is contained in the next 10 s. These probabilities are based only on the acoustic content

of their associated segments, and can give rise to biologically unrealistic call patterns. We used a post-processing step to remove isolated detected presence segments which are highly likely to be false positives rather than actual calls, and to obtain start and end times for each detected calling bout, to facilitate manual verification and support the main research objective of detecting and monitoring gibbon activity.

To do this, we formed connected components of presence segments that occur close together in time and in sufficient numbers that, given known gibbon call characteristics (i.e. song duration, inter-phrase duration), they are likely to be part of a single calling bout (Supplementary Material A). With presence segments arranged in temporal order, presence segment i is included in the same component as segment $i-1$ if they are separated by less than 200 s; otherwise segment i begins a new component. This process allocates each presence segment to exactly one component. Any component consisting of fewer than 20 segments (equivalent to roughly four phrases of length 5 s) are automatically removed. This was done given our analysis of the characteristics of the calls which revealed that the calls are typically repetitive over a period of time and the total duration was never less than 20 s. Additionally, any component where the average time between consecutive presence segments in the component was greater than 10 s, was removed (suggesting a 'chain' of isolated presence predictions, since calls usually persist over multiple consecutive segments).

The first and last presence segment in each remaining component give the start and end times of each predicted gibbon calling bout. To evaluate the potential usefulness of the post-processing step, we recalculated accuracy measures under the assumption that all detected bouts were subsequently passed to an observer for manual processing, and that this observer correctly identified all presence segments within the bout. This mimics the intended application of our approach, but means that post-processing

accuracy measures are conditional on the use of additional, error-free manual verification.

Results

Hainan gibbon calls could be detected with a high degree of accuracy. Without post-processing, nearly 80% of segments containing gibbon calls were correctly identified, with very few false positives (Table 1). Even with false negative rates of 20% very few gibbon phrases were missed altogether, because phrases occur across multiple overlapping segments and nearly all segments incorrectly identified as absences occurred at the beginning and end of a phrase, abutted by several segments where the phrase was correctly detected (Fig. 5). After post-processing, fewer than 2% of all presence segments occurred outside of detected call bouts (Table 1), and all 20 call bouts across nine test set recordings were detected, with two predicted call bouts being false positives (Supplementary Material B). In the training set, 34 of 35 call bouts were correctly recognised with two false positive call bouts.

The best performing approach was a 2-D CNN with both data augmentation and post-processing. Data augmentation improved specificity by 5.6%, a relative reduction in false positives of 79% but without associated relative reduction in sensitivity; post-processing further improved both sensitivity (20.6%) and specificity (0.9%, Table 1). Accuracy was substantially higher when treating the task as an image (spectrogram) classification problem than if the preprocessed acoustic data were directly used as input to a 1-D CNN. Using the 2-D CNN with both data augmentation and post-processing, an 8 h test file took on average 6 min to process of which 3 min 10 s were used for reading in the audio file and 2 min 42 s to convert to spectrograms; the remaining time was used to compute the CNN predictions.

We applied the 2-D CNN with both data augmentation and post-processing on the entire monitoring project and

Table 1. Average performance and parameter settings for the best 2-D and 1-D CNN models across 72 h of test recordings (2231 segments containing gibbon phrases, 23 689 without). Gibbon calls can be identified with very high accuracy, and performance is improved by data augmentation and a post-processing heuristic.

CNN	2-D	2-D	2-D	1-D	1-D	1-D
+ Augmentation	Yes	Yes	No	Yes	Yes	No
+ Post-processing	Yes	No	No	Yes	No	No
Accuracy	99.37%	97.60%	92.32%	94.30%	94.76%	94.76%
Sensitivity	98.30%	77.68%	79.65%	54.21%	40.98%	25.56%
Specificity	99.42%	98.51%	92.92%	95.96%	96.91%	97.60%
Precision	85.30%	70.28%	45.78%	49.14%	41.35%	44.58%
F1-score	90.55%	72.18%	53.14%	46.36%	38.42%	27.09%
Model Parameters	23 922	23 922	24 978	2650	2650	2378
Train Duration (sec)	644	643	265	628	627	117

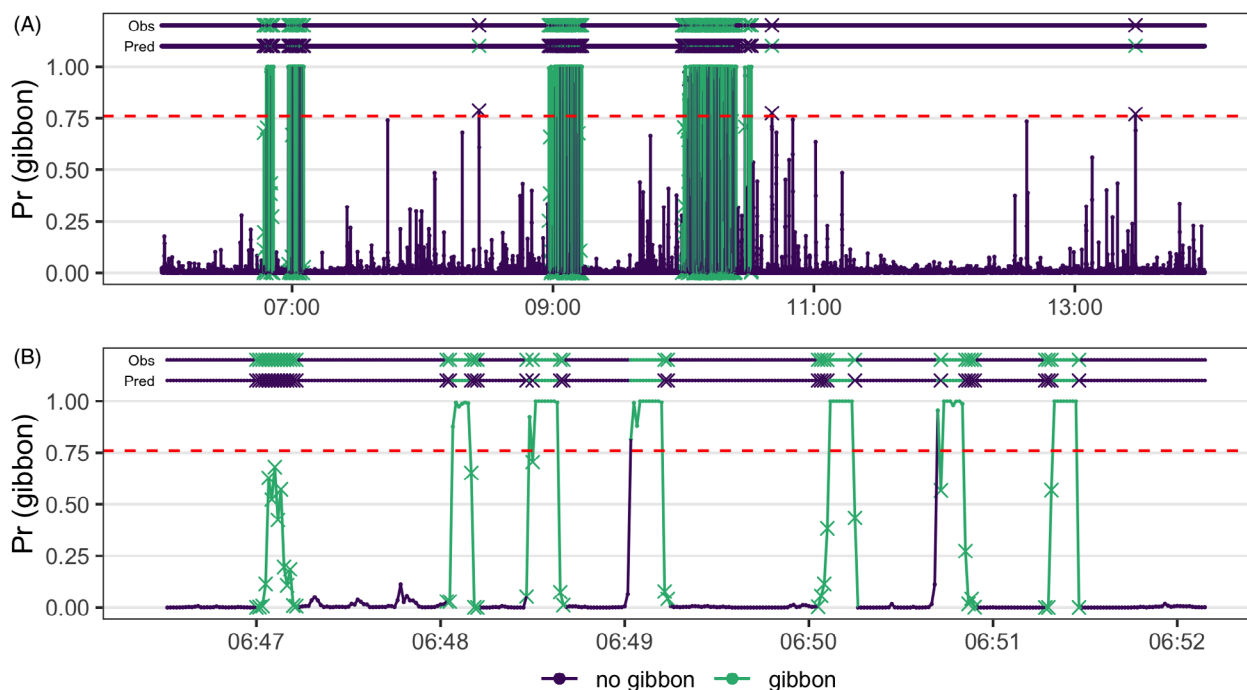


Figure 5. Per-second detected probabilities that a gibbon phrase is contained within the next 10 s of audio, over (A) an 8-h file, (B) a 5-min window. Segments with detected probabilities above an optimized threshold of 0.76 (red line) are classified as containing a gibbon phrase, with misclassifications denoted by crosses. Observed and detected classes are plotted above the probabilities, using the same notation. Colour is used to denote the observed class. Most incorrect false negative classifications are at the beginning and end of phrases, separated by segments that correctly identify the call. In this way, nearly all phrases are clearly identified, and a practitioner can be pointed to those regions that contain calls.

gibbon calls were detected on 71% of recording days across all locations. Gibbons were detected regularly at all locations, with recorders situated within known group or solitary home ranges detecting calls on 33–86% of recording days, and those situated between home ranges detecting calls on 46–89% of recording days. Mean durations of calling bouts per recorder varied between 24.2 and 40.8 min (overall mean = 29.7 min), with mean starting times of 06:16–07:56 AM and mean finishing times of 09:12–10:15 AM (Fig. 6; Table 2). Calls were detected less frequently during the wet season (March–April) than the dry season (May–August), with inter-season differences varying substantially between locations (Supplementary Table C).

Discussion

Long-term monitoring will generate thousands of hours of recordings across multiple survey sites, and manually labelling these recordings is typically infeasible given logistical constraints. Our results demonstrate that passive acoustic monitoring incorporating an automated classifier can be an effective tool for remote detection of calling species, potentially enabling systematic monitoring while

saving time, funds and manpower. Our approach, applied to Hainan gibbons, is general and easily extended to other calling species.

Our models allow new recordings to be classified on a per-second basis, to a high degree of accuracy. Although perhaps false negative rates of 1.7% may not be sufficiently low for full automation of Hainan gibbon call monitoring, they greatly facilitate the process of manually annotating these datasets by ruling out large portions of recordings that have a relatively low probability of containing gibbon song. In our test datasets, this reduced the amount of audio to be manually processed by 95%. Our model clearly detected all calling bouts in the test data, at the cost of two false positives. Where false negatives are particularly costly, this is easily incorporated by lowering the threshold required for manual verification. We expect that with more, and more diverse, training data, error rates would decline further.

Where environmental conditions were similar to those used to train the model, predictions were almost perfect and could be used to identify start and end times of call phrases and bouts, returning almost identical values to a human observer. It is impossible to know in advance whether environmental conditions are similar enough to

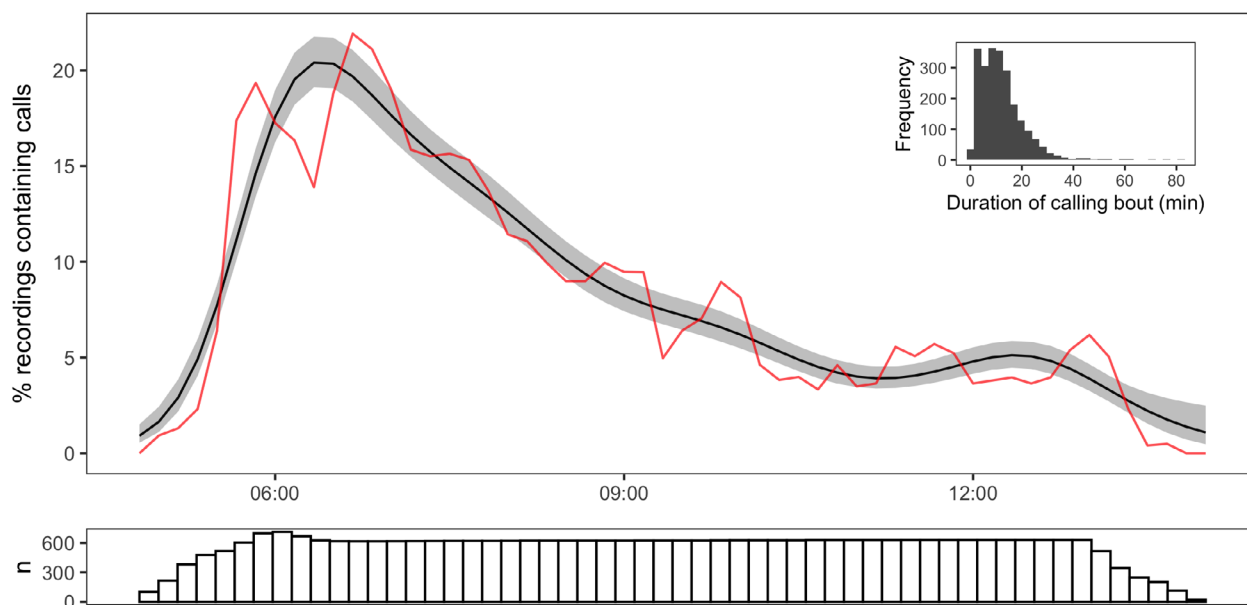


Figure 6. Daily patterns in gibbon calling activity. The red line denotes, per 10 min, the proportion of recordings across all locations in which a call was detected (e.g. 05:00–05:10, 05:10–05:20,...). The black line smooths the observed proportions using a GAM (see Supplementary Material D for details). The bottom plot shows the number of recordings per 10-min segment, showing the survey effort from 05:00 to 14:00. Peak activity occurs shortly after dawn, dropping rapidly but with some calling activity recorded throughout the morning. Plot inset shows the duration of independent call bouts detected by the classifier. Call bouts are intervals of regular calling, with no detected call 200 s either side of the bout. Daily calling typically consists of a number of calling bouts.

Table 2. Calling behaviour across eight survey locations for the 161 day survey period March–August 2016.

Location	Survey days	% days calls detected	Mean calling time per day (min)	Mean start time of first bout	Mean end time of last bout
1	87	70	24.2	07:34	09:41
2	90	46	29.9	06:58	09:12
3	103	82	31.3	07:30	10:15
4	105	86	26.5	07:44	09:52
5	79	33	29.9	07:31	09:23
6	103	79	24.4	07:56	10:15
7	129	89	30.9	06:53	09:54
8	105	65	40.8	06:16	10:01

Recorders were situated within the known home ranges of the four Hainan gibbon social groups existing during the study period, at locations intermediate between known home ranges, and in an area where a solitary male gibbon was thought to occur. Locations of home ranges are indicated by numbers 1, 2, 3 and 4. 6 = solitary.

warrant confidence in the associated predictions, but these results suggest that, as more training data covering a range of environmental conditions are added, model applications may go beyond gibbon detection, by automatically extracting inputs for more detailed behavioural analyses, for example of gibbon call syntax (Clarke et al., 2006).

Practically, developing an acoustic classifier such as ours requires a number of steps: deciding on an appropriate unit of analysis; manually labelling data; augmenting data and allocating it between training, validation and test sets; choosing and fitting appropriate neural network

models; and selecting a preferred model and using it to process the unlabelled portion of the data. Our study illustrates how model development and implementation are informed and guided by ecological objectives, here primarily detecting gibbon vocalizations over time scales of minutes or hours, and domain knowledge of Hainan gibbon call behaviour.

We based our classifier on phrases, rather than shorter notes or longer calling bouts, to balance ease of identification with data availability and computational requirements. Individual notes are easily confused with other sources (see Fig. 2B). While calling bouts are highly

distinctive, there are relatively few of them and, being longer in duration, they require more parameters to capture the same degree of detail. Phrases are far more numerous, less variable and require fewer parameters.

Given this choice, segment duration was chosen to be longer than the longest phrase across all training data (8 s). The slightly longer segment length provides more presence segments – for example, an 8 s phrase results in three 10 s presence segments, but would only result in a single segment if the segment length was restricted to 8 s. Preliminary runs based on shorter segments of 0.5–2 s and *partial* phrases did not yield good performance, with many false positives, probably because a small segment is not enough to distinguish gibbons from other species calling within the same frequency range.

Even using phrases, we have relatively few positive examples and these occur within a highly variable background environment, which is likely to be a common situation for ecological applications. The amount of data available to train neural networks is important, and CNNs tend to require relatively large amounts of data to generalize well. While preempting the exact amount of data required to train CNNs is challenging, one approach is simply to attempt to train a network and evaluate its performance on a test set and iteratively add data if need be. It may often be possible, as in our case, to collect or label additional data, but data augmentation is a valuable low-cost strategy for increasing sample sizes in conjunction with these other more effort-intensive approaches (Bergler et al., 2019; Hestness et al., 2017; Kahl et al., 2017; Sun et al., 2017). In practice the process can be an iterative one guided by subjective judgement. We initially annotated only 40 h across five recordings, but models based on these were poor, even with augmentation. Model performance (on the same test set) improved as we add more training data; we were also able to create more complex neural networks. Gains in accuracy decreased with additional annotations, and we stopped when these became marginal, but presumably further increases are possible as novel environments are included.

Training, validation and test datasets should be constructed by allocating longer contiguous sequences of audio to each of these, and then preprocessing each of these, rather than randomly allocating the segments themselves, which are highly autocorrelated and will thus overstate test accuracy. Wherever possible, we recommend using entirely independent recordings in the test dataset.

We found that 2-D CNNs based on spectrograms performed substantially better than 1-D CNNs that use amplitude time series following some initial preprocessing, mirroring Stowell et al. (2019b). Deep neural networks are often motivated by an argument that they learn salient features, rather than having to have these provided

to them, but where intermediate features (here, spectral densities) can be provided, these speed up the learning process and provide measurable benefits. Beyond the 2-D/1-D distinction, we found that there was little impact on the model performance when different configurations (i.e. changing the number of filters or units) to small networks were explored. Large networks with a much larger number of layers did not improve the performance and we achieved good performance using relatively small, simple network architectures. We used few dense layers, each with only a small number of nodes, as these are particularly parameter hungry. Our basic approach was to start with simple architectures, evaluate them, and then add complexity in an iterative manner.

Traditional performance metrics such as precision and sensitivity (recall), while important, are not the only relevant measures of classifier success. Practically, classifiers such as ours can be used to point to audio segments that possibly contain gibbon calls, and that require manual verification. Where classification accuracy lags behind that of human experts, or where errors are costly – that is, in many ecological applications – attention shifts from replacing manual annotation to facilitating it. Probability cutoffs can be calibrated to balance the costs of false positives and negatives, and, even if the model is wrong by a few seconds, the amount of time spent in manual verification, compared to that required to processing the entire file manually, is minimal. Our classifier reduces an 8-h recording to on average 22 min with false positive and negative rates under 2%. This time can be further reduced by playing back only those 10 s segments that are predicted to contain phrases, although in our case the reduction in overall time was offset by the difficulty of manually verifying segments that are often not contiguous in time.

Analysis of our multi-month dataset demonstrated that gibbons could be detected regularly across all selected survey points, with call detection consistent with known patterns of gibbon behaviour and ecology. Calls were detected at expected times (Chan et al., 2005), and our dataset provides a more precise baseline on Hainan gibbon call timing and duration. Hainan gibbon calling bouts were also generally detected less frequently during the wet season, a period when other gibbon species are also known to sing less frequently (Cheyne, 2008; Clink et al., 2020). Interestingly, call bouts recorded within the area occupied by a solitary male gibbon were among the shortest recorded bouts, and started and finished later than bouts from known social groups. While we cannot exclude the possibility of detecting group calls at this location, this finding suggests important new information on the behavioural ecology of solitary Hainan gibbons that may assist future monitoring and conservation planning.

It is uncertain whether within-recorder and between-recorder variation in calling bout detections represents variation in calling frequency between groups, and/or variation in detection effectiveness by recorders, with the latter possibility likely associated with specific recorder placement, local terrain, specific gibbon movement patterns across landscapes, and group home range size (cf. Bryant et al., (2017)). Future work could investigate detection likelihood in relation to specific environmental parameters and local weather conditions (e.g. rainfall, wind, temperature), data on which were not available for our survey period but are known to affect calling behaviour in other gibbons (Coudrat et al., 2015; Yin et al., 2016).

Where calls can be detected across multiple recording locations, acoustic spatial capture–recapture methods provide a means of estimating animal abundance (Stevenson et al., 2015). While our locations are too far apart for this to be feasible, this represents an important next step in monitoring a critically endangered population. Classifiers capable of discriminating between groups or individuals can be valuable inputs to this process (Augustine et al., 2020), as well as providing insight into the behavioural ecology of groups or individuals. We also recommend that call detection ranges should be determined for the specific field conditions at BNNR (e.g. slope, vegetation density), to calibrate monitoring effectiveness of specific recorders, and determine effective recorder placement (grid area/density) to ensure saturation of monitoring coverage. However, passive acoustic monitoring can now be introduced as an important component of the Hainan gibbon conservation toolkit, both for future use at BNNR and also to potentially detect unknown remnant gibbon populations elsewhere across Hainan (Turvey et al., 2017). Our classifier permits rapid and potentially real-time monitoring of Hainan gibbons, and we hope that the approach we describe in developing this classifier can serve as a roadmap for practitioners to implement their own classifier for other passive acoustic monitoring projects, and contribute to the effective conservation of calling species.

Acknowledgments

We thank the Management Office of Bawangling National Nature Reserve for logistical assistance in the field. Fieldwork was funded by an Arcus Foundation grant to STT and a Wildlife Acoustics grant to JVB. ID is supported in part by funding from the National Research Foundation of South Africa (Grant ID 90782, 105782). ED is supported by a postdoctoral fellowship from the African Institute for Mathematical Sciences South Africa, Stellenbosch University and the Next Einstein Initiative. This work was carried out with the aid of a grant from the

International Development Research Centre, Ottawa, Canada (www.idrc.ca), and with financial support from the Government of Canada, provided through Global Affairs Canada (GAC; www.international.gc.ca). We also thank the following rangers who contributed to data collection: Guang Wei, Zhong Zhao, Qing Lin, Jinbing Zhang, Zhicheng Zhang, Qianjin Li, Xiaoliang Fu, Zhengchong Zhou, Lubiao Huang, Zhengkun Ye, Zhenghai Zou, Jinqiang Wang, Wentao Han and Zengnan Xie.

Conflict of Interest

The authors declare no competing interests.

Authors' Contributions

ST, JB and HM conceived the passive monitoring project and developed the study designs and protocols. ED, ID and JH conceived the development of an automated classifier and designed the methodology. WL, ZL, QC, ZZ, HM and JB were responsible for fieldwork and data collection. ED, AH and JB annotated the data. ED constructed the classifier and performed the analysis. ED, ID and ST wrote the paper. All authors contributed critically to the drafts and gave final approval for publication.

Data Availability Statement

All code for training and testing the neural networks and conducting additional analyses is available at <https://github.com/emmanueldufourq/GibbonClassifier> (Supplementary Material E). A subset of acoustic recordings, including training and testing labels, has been stored on Zenodo: <https://doi.org/10.5281/zenodo.3991714>.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C.. et al. (2015) *Large-scale machine learning on heterogeneous systems*. Retrieved from <https://www.tensorflow.org>
- Augustine, B.C., Royle, J.A., Linden, D.W. & Fuller, A.K. (2020) Spatial proximity moderates genotype uncertainty in genetic tagging studies. *Proceedings of the National Academy of Sciences*, **117**(30), 17903–17912.
- Bergler, C., Schröter, H., Cheng, R.X., Barth, V., Weber, M., Nöth, E. et al. (2019) Orca-spot: an automatic killer whale sound detection toolkit using deep learning. *Scientific Reports*, **9**(1), 1–17.
- Bermant, P.C., Bronstein, M.M., Wood, R.J., Gero, S. & Gruber, D.F. (2019) Deep machine learning techniques for the detection and classification of sperm whale bioacoustics. *Scientific Reports*, **9**(1), 1–10.
- Brockelman, W.Y. & Srikosamatara, S. (1993) Estimation of density of gibbon groups by use of loud songs. *American*

- Journal of Primatology*, **29**(2), 93–108. <https://doi.org/10.1002/ajp.1350290203>
- Bryant, J.V., Brulé, A., Wong, M.H., Hong, X., Zhou, Z., Han, W. et al. (2016) Detection of a new Hainan gibbon (*Nomascus hainanus*) group using acoustic call playback. *International Journal of Primatology*, **37**(4–5), 534–547.
- Bryant, J.V., Zeng, X., Hong, X., Chatterjee, H.J. & Turvey, S.T. (2017) Spatiotemporal requirements of the Hainan gibbon: does home range constrain recovery of the world's rarest ape? *American Journal of Primatology*, **79**(3), e22617.
- Cannam, C., Landone, C. & Sandler, M. (2010) Sonic visualiser: an open source application for viewing, analysing, and annotating music audio files. In *Proceedings of the ACM Multimedia 2010 International Conference, Firenze, Italy*. New York, NY: Association for Computing Machinery, pp. 1467–1468.
- Chan, B.P.L., Fellowes, J., Geissmann, T. & Zhang, J. (2005) *Hainan gibbon status survey and conservation action plan*. technical report 3
- Cheyne, S.M. (2008) Effects of meteorology, astronomical variables, location and human disturbance on the singing apes: *Hylobates albibarbis*. *American Journal of Primatology*, **70**(4), 386–392.
- Chilson, C., Avery, K., McGovern, A., Bridge, E., Sheldon, D. & Kelly, J. (2019) Automated detection of bird roosts using Nexrad Radar data and convolutional neural networks. *Remote Sensing in Ecology and Conservation*, **5**(1), 20–32.
- Chollet, F. (2015) *Keras*. <https://keras.io>
- Christin, S., Hervet, E. & Lecomte, N. (2019) Applications for deep learning in ecology. *Methods in Ecology and Evolution*, **10**(10), 1632–1644.
- Clarke, E., Reichard, U.H. & Zuberbühler, K. (2006) The syntax and meaning of wild gibbon songs. *PLoS One*, **1**(1), e73.
- Clink, D.J., Ahmad, A.H. & Klinck, H. (2020) Gibbons aren't singing in the rain: presence and amount of rainfall influences ape calling behavior in Sabah, Malaysia. *Scientific Reports*, **10**(1), 1–13.
- Clink, D.J., Crofoot, M.C. & Marshall, A.J. (2019) Application of a semi-automated vocal fingerprinting approach to monitor Bornean gibbon females in an experimentally fragmented landscape in Sabah, Malaysia. *Bioacoustics*, **28**(3), 193–209. <https://doi.org/10.1080/09524622.2018.1426042>
- Coudrat, C., Nanthavong, C., Ngoprasert, D., Suwanwaree, P. & Savini, T. (2015) Singing patterns of white-cheeked gibbons (*Nomascus* sp.) in the annamite mountains of Laos. *International Journal of Primatology*, **36**(4), 691–706.
- Crunchant, A.-S., Borchers, D., Kuhl, H. & Piel, A. (2020) Listening and watching: do camera traps or acoustic sensors more efficiently detect wild chimpanzees in an open habitat? *Methods in Ecology and Evolution*, **11**(4), 542–552. <https://doi.org/10.1111/2041-210X.13362>
- Deng, H., Zhou, J. & Yang, Y. (2014) Sound spectrum characteristics of songs of Hainan Gibbon (*Nomascus hainanus*). *International Journal of Primatology*, **35**(2), 547–556.
- Fairbrass, A.J., Firman, M., Williams, C., Brostow, G.J., Titheridge, H. & Jones, K.E. (2019) CityNet-deep learning tools for urban ecoacoustic assessment. *Methods in Ecology and Evolution*, **10**(2), 186–197.
- Grill, T. & Schlüter, J. (2017) *Two convolutional neural networks for bird detection in audio signals*. 1764–1768.
- Hestness, J., Narang, S., Ardalani, N., Diamos, G., Jun, H., Kianinejad, H. et al. (2017) *Deep learning scaling is predictable, empirically*. Preprint <https://arxiv.org/abs/1712.00409>
- Huang, X., Acero, A. & Hon, H.-W. (2001) *Spoken language processing: a guide to theory, algorithm, and system development*. Upper Saddle River, NJ: Prentice Hall PTR.
- Jiang, J.-J., Bu, L.-R., Duan, F.-J., Wang, X.-Q., Liu, W., Sun, Z.-B. et al. (2019) Whistle detection and classification for whales based on convolutional neural networks. *Applied Acoustics*, **150**, 169–178.
- Kahl, S., Wilhelm-Stein, T., Hussein, H., Klinck, H., Kowerko, D., Ritter, M. et al. (2017) Large-scale bird sound classification using convolutional neural networks. In *Working notes of CLEF*.
- Kidney, D., Rawson, B.M., Borchers, D.L., Stevenson, B.C., Marques, T.A. & Thomas, L. (2016) An efficient acoustic density estimation method with human detectors applied to gibbons in Cambodia. *PLoS One*, **11**(5), 1–16. <https://doi.org/10.1371/journal.pone.0155066>
- Kingma, D.P. & Ba, J. (2014) *Adam: a method for stochastic optimization*. Preprint <https://arxiv.org/abs/1412.6980>
- Kiskin, I., Zilli, D., Li, Y., Sinka, M., Willis, K. & Roberts, S. (2020) Bioacoustic detection with wavelet-conditioned convolutional neural networks. *Neural Computing and Applications*, **32**(4), 915–927.
- Kulyukin, V., Mukherjee, S. & Amlathe, P. (2018) Toward audio beehive monitoring: deep learning vs. standard machine learning in classifying beehive audio samples. *Applied Sciences*, **8**(9), 1573.
- LeCun, Y., Bengio, Y. & Hinton, G. (2015) Deep learning. *Nature*, **521**(7553), 436–444.
- Liu, H., Ma, H., Cheyne, S.M. & Turvey, S.T. (2020) Recovery hopes for the world's rarest primate. *Science*, **368**(6495), 1074.
- McFee, B., Lostanlen, V., McVicar, M., Metsai, A., Balke, S., Thome, C. et al. (2020) *Librosa*. <https://doi.org/10.5281/zenodo.3606573>
- Nolasco, I., Terenzi, A., Cecchi, S., Orcioni, S., Bear, H.L. & Benetos, E. (2019) Audio-based identification of beehive states. In *2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, UK, pp. 8256–8260.
- Pamula, H., Pocha, A. & Klaczynski, M. (2019) Towards the acoustic monitoring of birds migrating at night. *Biodiversity Information Science and Standards*, **3**, e36589.

- Patterson, J. & Gibson, A. (2017) *Deep learning: a practitioner's approach* (M. Loukides & T. McGovern (Eds.)). Sebastopol, CA: O'Reilly Media Inc.
- Qazi, K.A., Tabassam Nawaz, Z.M., Rashid, M. & Habib, H.A. (2018) A hybrid technique for speech segregation and classification using a sophisticated deep neural network. *PLoS One*, **13**(3), e0194151.
- Sandler, M., Howard, A., Zhu, M., Zhmoginov, A. & Chen, L. (2018) Mobilenetv 2: inverted residuals and linear bottlenecks. In *IEEE conference on computer vision and pattern recognition (CVPR)*, pp. 4510–4520.
- Sethi, S.S., Jones, N.S., Fulcher, B.D., Picinali, L., Clink, D.J., Klinck, H. et al. (2020) Characterizing soundscapes across diverse ecosystems using a universal acoustic feature set. *Proceedings of the National Academy of Sciences*, **117**(29), 17049–17055.
- Shiu, Y., Palmer, K., Roch, M.A., Fleishman, E., Liu, X., Nosal, E.-M. et al. (2020) Deep neural networks for automated detection of marine mammal species. *Scientific Reports*, **10**(1), 1–12.
- Sprengel, E., Jaggi, M., Kilcher, Y. & Hofmann, T. (2016) *Audio based bird species identification using deep learning techniques*. 2016 Conference and Labs of the Evaluation Forum.
- Stevenson, B.C., Borchers, D.L., Altwegg, R., Swift, R.J., Gillespie, D.M. & Measey, G.J. (2015) A general framework for animal density estimation from acoustic detections across a fixed microphone array. *Methods in Ecology and Evolution*, **6**(1), 38–48.
- Stowell, D., Petrusková, T., Šálek, M. & Linhart, P. (2019) Automatic acoustic identification of individuals in multiple species: improving identification across recording conditions. *Journal of the Royal Society Interface*, **16**(153), 20180940.
- Stowell, D., Wood, M.D., Pamuła, H., Stylianou, Y. & Glotin, H. (2019) Automatic acoustic detection of birds through deep learning: the first bird audio detection challenge. *Methods in Ecology and Evolution*, **10**(3), 368–380.
- Sun, C., Shrivastava, A., Singh, S. & Gupta, A. (2017) Revisiting unreasonable effectiveness of data in deep learning era. In: *2017 IEEE international conference on computer vision (ICCV)*, pp. 843–852. <https://doi.org/10.1109/ICCV.2017.97>
- Turvey, S.T., Bryant, J.V., Duncan, C., Wong, M.H., Guan, Z., Fei, H. et al. (2017) How many remnant gibbon populations are left on Hainan? Testing the use of local ecological knowledge to detect cryptic threatened primates. *American Journal of Primatology*, **79**(2), e22593.
- Turvey, S., Traylor-Holzer, K., Wong, M., Bryant, J., Zeng, X., Hong, X. & et al. (2015) *International conservation planning workshop for the hainan gibbon: final report*. Zoological Society of London, London, UK IUCN SSC Conservation Breeding Specialist Group, Apple Valley, MN, USA
- Vu, T.T. & Tran, L.M. (2019) An application of autonomous recorders for gibbon monitoring. *International Journal of Primatology*, **40**(2), 169–186. <https://doi.org/10.1007/s10764-018-0073-3>
- Yin, L.-Y., Fei, H.-L., Chen, G.-S., Li, J.-H., Cui, L.-W. & Fan, P.-F. (2016) Effects of group density, hunting, and temperature on the singing patterns of eastern Hoolock gibbons (*Hoolock leuconedys*) in gaoligongshan, southwest china. *American Journal of Primatology*, **78**(8), 861–871.
- Zhang, H., Wang, C., Turvey, S.T., Sun, Z., Tan, Z., Yang, Q. et al. (2020) Thermal infrared imaging from drones can detect individuals and nocturnal behavior of the world's rarest primate. *Global Ecology and Conservation*, **23**, e01101. <https://doi.org/10.1016/j.gecco.2020.e01101>
- Zhou, X., Guan, Z., Zhong, E., Dong, Y., Li, H. & Hu, K. (2019) *Automated monitoring of western black crested gibbon population based on voice characteristics*. 5th International Conference on Computer and Communications, pp. 1383–1387.

Supporting Information

Additional supporting information may be found online in the Supporting Information section at the end of the article.

Data S1. Supplementary information and results which provides details of the characteristics of the call bouts in the training data. Specific model predictions for the testing data is provided. An analysis on the seasonal differences in the detections is presented. Details regarding the software pipeline is provided.