



# Acoustic detection of regionally rare bird species through deep convolutional neural networks

Ming Zhong<sup>a</sup>, Ruth Taylor<sup>b</sup>, Naomi Bates<sup>b</sup>, Damian Christey<sup>b</sup>, Hari Basnet<sup>b</sup>, Jennifer Flippin<sup>b</sup>, Shane Palkovitz<sup>b</sup>, Rahul Dodhia<sup>a</sup>, Juan Lavista Ferres<sup>a,\*</sup>

<sup>a</sup> AI for Good Research Lab, Microsoft, USA

<sup>b</sup> Songs of Adaptation, Future Generations University, USA

## ARTICLE INFO

### Keywords:

Deep learning  
Convolutional neural networks (CNN)  
Bioacoustic classification  
Transfer learning  
Species population  
Presence survey

## ABSTRACT

Bioacoustic monitoring with machine learning (ML) models can provide valuable insights for informed decision-making in conservation efforts. In this study, the team built deep convolutional neural networks to analyze field recordings and classify calls of Yellow-vented warbler (*Phylloscopus cantator*) and Rufous-throated wren-babbler (*Spelaeornis caudatus*), both of which are regionally rare in Nepal. Data augmentation techniques for calls of the two bird species were utilized to effectively increase the size of the training set and thus boost model performance. Nepali ornithologists were engaged in iterative data labeling from field recordings, leveraging ML technology in conjunction with expert manual labeling and verification. The model output provides insights of species activity and abundance throughout 2018–2019 in multiple ecosystems along an elevational transect in the Barun River Valley, Nepal. The results of this study may help conservationists better understand species distribution, behavior, diversity, and habitat preference. Additionally, the results provide baseline data to quantify future changes due to habitat disruption or climate change. This modeling methodology and its framework can be easily adopted by other acoustic classification problems.

## 1. I. Introduction

In recent decades, the populations of various animals, including birds, amphibians, insects, and mammals, have exhibited steep declines worldwide. While many decreases are due to habitat loss and over-utilization, other unidentified processes threaten 48% of rapidly declining species and are driving species most quickly to extinction (Stuart et al., 2004). As biodiversity plays a critical role in many aspects, well-designed monitoring programs provide a basis for identifying the species, sites and threats of most significant concern. Such monitoring programs also provide reliable tools when evaluating the integrity of ecosystems and their responses to disturbances, assessing progress in efforts to conserve biodiversity, and measuring the success of actions taken to preserve or recover biodiversity. However, manual observation remains limited and challenging in many scenarios, especially in the areas that are difficult to access physically or when the focus is to study animals' night-time behavior. In such scenarios, passive acoustic monitoring is highly appropriate, as many birds, including rare species, are most readily detectable by their sounds, often more so than by

vision. With modern remote monitoring stations, it can continuously monitor large remote areas for avian community composition and tracking migratory and seasonal changes in populations (Aide et al., 2013; Frommolt, 2017a; Furnas and Callas, 2015a; Hill et al., 2017; Knight et al., 2017; Matsubayashi et al., 2017).

Earlier applications that have employed such technology either performed automatic recording but relied on manual analysis of sound recordings (Frommolt, 2017b; Furnas and Callas, 2015b) or were based on low-complexity signal processing such as template matching (Colonna et al., 2015; Towsey et al., 2012), feature extraction (Mesaros et al., 2018), or traditional machine learning methods (Keen et al., 2014; Zhao et al., 2017).

With the constant increase in computing power and the development of more efficient codes, high-performance computing helps the extremely fast growth of deep learning in recent years, which has been shown to outperform previous state-of-the-art techniques in several tasks. Deep learning has fueled great strides in a variety of computer vision problems, and in particular, Convolutional Neural Networks (CNN) have demonstrated great potential and success in image

\* Corresponding author at: AI for Good Research Lab, Microsoft, Redmond, WA 98052, USA.

E-mail address: [jlavista@microsoft.com](mailto:jlavista@microsoft.com) (J. Lavista Ferres).

classification tasks and thus drawn much attention in constructing the automatic bird sound classification systems. Some popular CNN architectures applied to bioacoustics classification include AlexNet (Krizhevsky et al., 2012), LeNet-5 (LeCun, 2015), VGG16 (Simonyan and Zisserman, 2015), ResNet50 (He et al., 2016), among others.

In this study, two regionally rare species were chosen: Yellow-vented warbler (*Phylloscopus cantator*) and Rufous-throated wren-babbler (*Spelaornis caudatus*). The Rufous-throated Wren Babbler is a very rare bird that has an extremely limited range in Nepal. The species is Near Threatened globally; it is listed within Nepal as a Critically Endangered species on a national level (BirdLife International, 2020; Inskipp et al., 2016). The nationally endangered Yellow-vented Warbler can be found in the East of Nepal. It is recorded between 75 m and 1525 m in a few locations, including Makalu Barun National Park (Inskipp et al., 2016). These species provided a proof of concept demonstrating that with limited training samples, deep learning models can classify rare species calls.

As a research project of Future Generations University, this project brings expertise in community development with decades-long global partnerships that ensure long-term data collection and research permissions, data labeling, and collaboration for sustainable, just, and lasting climate action. Protecting 100,000,000 acres of land, Future Generations leadership in community-based conservation established multiple national parks across Asia. Over the past 27 years, the University has employed key indicators (quick, easy-to-use measurements), shaped to fit specific communities' contexts, to empower community members to measure change over time for themselves. This project is

unique in its commitment to community engagement.

## 2. Data collection and pre-processing

Audio data were collected from 8 different sites along an elevational transect in the mountains of Makalu Barun National Park, Nepal, between 2018 and 2019. Audios were recorded into wav format using Song Meter SM4 Acoustic Recorders (Wildlife Acoustics) at a sampling rate of 48 kHz and 24-bit rate, and recorders were programmed to record 5-min audio every 15 min 24 h per day.

The training and test datasets were initially generated using pattern recognition clustering software (Kaleidoscope Pro Analysis Software, Wildlife Acoustics (Wildlife Acoustics, 2019)), and local avian experts subsequently analyzed clusters in Nepal to identify calls from the two species of interest, Yellow-vented warbler (*Phylloscopus cantator*) and Rufous-throated wren-babbler (*Spelaornis caudatus*). A target of 100+ positive detections for each species and 300+ negative detections (ex. rain wind river, insects, other bird species, etc.) were set as the minimum amount of data necessary for model training and development. The positive and negative samples were used as input to train CNN models. Spectrogram images containing a target positive or negative sample were standardized using a 4-s audio clip beginning at each detection's start-time.

Spectrograms were extracted from audio files (with NFFT = 256, Hanning window) using Python 3.6 and then resized to 224 by 224 pixels with RGB channels and stored as color PNG images (see Fig. 1 for example). The color spectrograms were the input for the machine

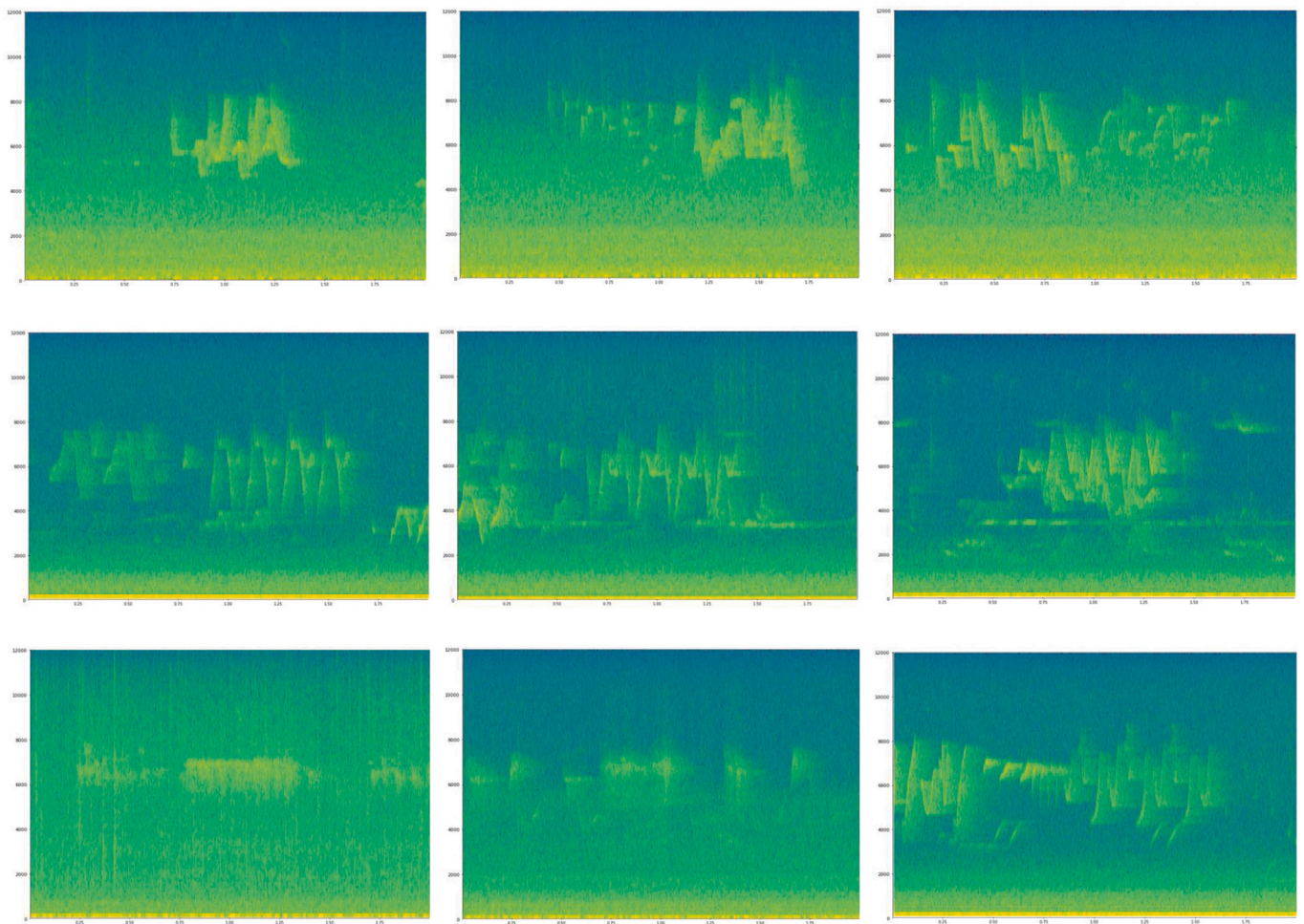


Fig. 1. Sample spectrograms for 4-s audio recordings. First row: calls from the species *Phylloscopus cantator*; Second row: calls from the species *Spelaornis caudatus*; Third row from left to right: rain background, river background, unidentified species.

learning model, and the corresponding single-species labels for each image (i.e. species present (positive) or absent (negative)) were used as the ground truth data for training and evaluating the classification model.

### 3. Approaches

#### 3.1. Transfer learning and fine-tuning with a pre-trained CNN model

Here the neural network model ResNet50 was applied to classify the calls of the 2 bird species. This ResNet50 CNN architecture is a variant of ResNet model which has 48 Convolution layers along with 1 Max Pooling and 1 Average Pooling layer. It begins with the RGB images (size  $224 \times 224 \times 3$ ) as input and performs the initial convolution and max-pooling using  $7 \times 7$  and  $3 \times 3$  kernel sizes, respectively. Afterward, it stacks a series of residual blocks. With the skip connection of residual blocks, it allows the model to propagate larger gradients to initial layers. These layers are able to learn as fast as the final layers, in order to train deeper networks. Finally, the network has an average pooling layer, followed by a fully connected layer. When training the ResNet50 model, the Adam optimizer algorithm was applied, and an initial learning rate of  $1e-4$  with a decay factor of  $1e-7$ .

In the context of deep learning, most models include millions of parameters. ResNet50, for example, has 23 million parameters. To train such complex models, it typically requires an extensive dataset to achieve an optimal parameter configuration. However, in practice it may be very difficult to collect large amounts of labeled data, especially if a species rarely calls or if the species is endangered and there are few individuals. Besides, using experts to obtain a large number of labeled samples in acoustics is an expensive and time-consuming endeavor. Given this scenario, transfer learning with fine-tuning (Shin et al., 2016) is a useful technique when there is only a small number of labeled data available.

Transfer learning is a machine learning technique where a model trained on one task (or domain) is re-purposed on a second related task (or domain). Pre-trained models are usually shared in the form of the millions of parameters/weights the model achieved while being trained to an optimal state. In this study, the model weights were initially trained on ImageNet (Deng et al., 2009) dataset with 1000 classes of objects, but their pre-trained weights can be leveraged by a different task or domain (Huh et al., 2016). This approach is effective because the source model was trained on a large number of images and made predictions on a relatively large number of classes. In turn, it required the model to extract distinct features from images in order to perform well. With fine-tuning, some layers are frozen from the pre-trained model, and it is sufficient to train the last several layers only, instead of having to train the whole model with random initialization of all parameters.

In this study, the model design included pre-trained weights of ResNet50 and fine-tuned parameters, adding a fully connected layer, a dropout layer and an output layer.

#### 3.2. K-fold cross-validation

In this dataset, there are only a few hundreds of detected calls for the two target species, *Phylloscopus cantator* and *Spelaornis caudatus*, that include different stereotypes of calls from each species. By partitioning the available data into three sets (training, validation and testing), we drastically reduce the number of samples which can be used for learning the model, and the results that depend on a particular random choice for the three sets are not stable. A solution to this problem is a procedure called K-fold cross-validation, which generally results in a less biased model compared to other methods. With this procedure, it ensures every observation from the original dataset has the chance of appearing in the training and test set. This is one of the best approaches if we have limited input data. This method follows the below steps:

Step 1: Split the entire data randomly into K folds (here, we use  $K =$

5).

Step 2: Fit the model (training and validation) using the  $K - 1$  ( $K$  minus 1) folds and test the model using the remaining Kth fold. Note down the scores/errors.

Step 3: Repeat this process until every K-fold serves as the test set. Then take the average of all recorded scores. That will be the performance metric for the model.

#### 3.3. Data augmentation

While many deep neural network models have parameters in the order of millions, they are heavily reliant on big data to avoid overfitting. Unfortunately, in many real-world applications, the amount of data that can be used for training is rather limited, either due to the huge manual efforts required to collect data, or due to the fact that it is almost impossible to acquire large amounts of data in some cases. As an effective data-space solution to the problem of limited data, data augmentation encompasses a suite of techniques that enhance the size and quality of training datasets such that better deep learning models can be built using them.

Among various data augmentation methods for image processing, some basic ones include flips, rotations, shifts, noise injections, color space transformations, sharpening or blurring, and random erasing or cropping. Specifically, for audio recordings, there are methods such as time-stretching, pitch shifting, and mixing multiple audios (Salamon and Bello, 2017). Beyond them, there are more advanced techniques, for example, generative adversarial network (GAN)-based methods (Shorten and Khoshgofaar, 2019), which can be used to generate synthetic images.

For this model implementation, basic techniques were applied to increase the size of data that can be used for model training: rotation (up to 5 degrees), shifting (width and height shifting up to 10% of the original spectrogram), and cropping.

Another effective method we adopted to boost the training data size is to use spectrograms with smaller time-windows. While the detected calls for the two regionally rare species, *Phylloscopus cantator* and *Spelaornis caudatus*, usually last for 2 s or longer (see Fig. 1. as an example), our baseline model was fit based on spectrograms generated from a 4-s time window. In order to boost the size of training data, we break down each 4-s detection into 3 shorter detections, where each detection lasts for 2 s (that is, to create 3 spectrograms starting at second 0, 1, and 2, respectively, from each original 4-s detection). Even though breaking down spectrograms into 2-s windows may not include one complete call within each spectrogram and may bring some noisy labels during model training, but with this implementation, the size of data available for training tripled.

#### 3.4. Model training

Our manually validated dataset consists of 195 positive detections for *Phylloscopus cantator*, 320 positive detections for *Spelaornis caudatus*, and 1060 negative detections composed of various types of noises (rain, wind, river, bugs, other bird species, etc.), where each detection lasts for 4 s.

Finding sufficient clips of exemplar training data from the field recordings was challenging, because of call volume variations, overlapping calls with other species, and background noises. In addition, to distinguish a species with multiple and varying calls, it was also essential and challenging to determine other species that had similar calls to the target species and label these close calls as negative training data. Particularly, for the target species *Spelaornis caudatus*, there are two other species that have acoustically similar calls (Fig. 2).

After scoring, all the false positives in the training data were verified by experts and correctly labeled to retrain the model after the first round of model training. External training data (exemplar calls manually verified from [xeno-canto.org](http://xeno-canto.org)) were added to supplement the project's

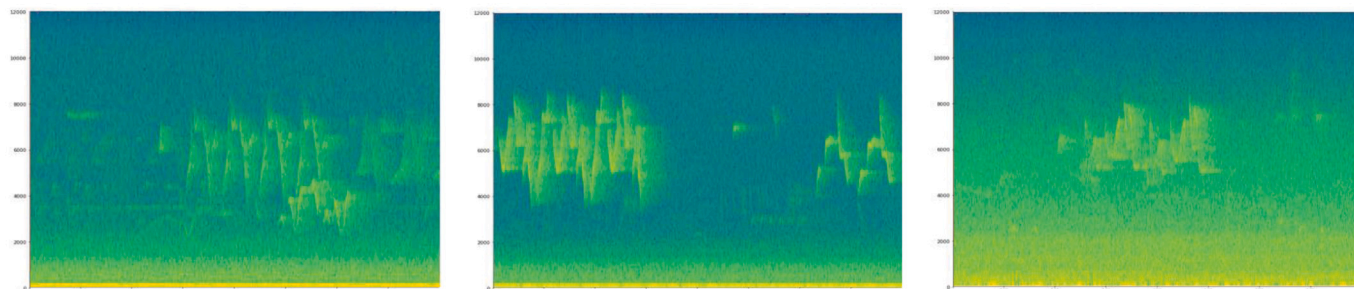


Fig. 2. Spectrogram of *Spelaeornis caudatus* and two other species (*Phylloscopus reguloides* and *Pnoepyga albiventer*) with acoustically similar calls.

data.

## 4. Results

### 4.1. Model performance

Three key metrics are reported to evaluate and compare the performance of the model on the testing data set: 1) sensitivity (true positive rate, recall); 2) specificity (true negative rate), and 3) area under a curve (AUC). Sensitivity measures the proportion of true positives that were identified correctly; and specificity measures the proportion of true negatives that were identified correctly. While sensitivity and specificity are dependent on the choice of threshold score, the area under a curve (AUC) provides an aggregate measure of performance across all possible classification thresholds. It is not affected by the class imbalance.

For both species *Phylloscopus cantator* and *Spelaeornis caudatus*, the model based on 2-s spectrograms performed significantly better, especially sensitivity, compared to the model based on 4-s spectrograms. Using data augmentation made further improvement for the model based on 2-s spectrograms (Table 1). The sensitivity for classifying the species *Spelaeornis caudatus* was not as good as that of the model classifying the species *Phylloscopus cantator*, and resulted in about 10% of detections that were misclassified as negative. A closer investigation of the data revealed that the labeled calls for *Spelaeornis caudatus* included detections with various levels of clarity, different call stereotypes, and maybe some incorrectly labeled detections. It appears that the neural network model did not find enough commonalities among these detected calls to make correct classification. Some examples of spectrograms

Table 1

Classification results (sensitivity, specificity, and AUC) for both target species by each CNN model. The results are based on the average score of conducting 5-fold cross-validation, with a neutral threshold score 0.5. The highest performance for each measure and acoustic population is in boldface type.

Species	CNN Model Description	Sensitivity (%)	Specificity (%)	AUC (%)
<i>Phylloscopus cantator</i>	based on 4-s spectrograms, no data augmentation	86.15	98.91	98.62
	based on 2-s spectrograms, no data augmentation	94.92	99.92	99.02
	based on 2-s spectrograms, with data augmentation	<b>95.94</b>	<b>99.92</b>	<b>99.58</b>
<i>Spelaeornis caudatus</i>	based on 4-s spectrograms, no data augmentation	53.12	94.82	91.50
	based on 2-s spectrograms, no data augmentation	78.46	<b>95.96</b>	97.05
	based on 2-s spectrograms, with data augmentation	<b>90.15</b>	93.92	<b>97.85</b>

are shown in Fig. 3.

### 4.2. Scoring on unlabeled data

The model was run using over one year of data with dates ranging from 3/2018 to 7/2019 using data from three stations in Makalu Barun National Park around the elevations where the target species were expected - Hinju Camp (elevation 1820 m), Deurali danda (elevation 2100 m), and Tutin Camp (elevation 2300 m).

In order for results to be analyzed, a threshold needs to be chosen for what probability will be counted as the presence of the species. Table 2 shows the number of detected calls for three sample threshold probability ranges (clip numbers rounded to the nearest ten). While the model predicts the probability of target species calls for each extracted spectrogram from the corresponding audio clip, the probability itself does not give a definite answer of presence/absence of species calls. As our next step, we will send those results to the local ecologists and conduct output validation by sampling spectrograms with different predicted probability ranges and then choosing the optimal threshold.

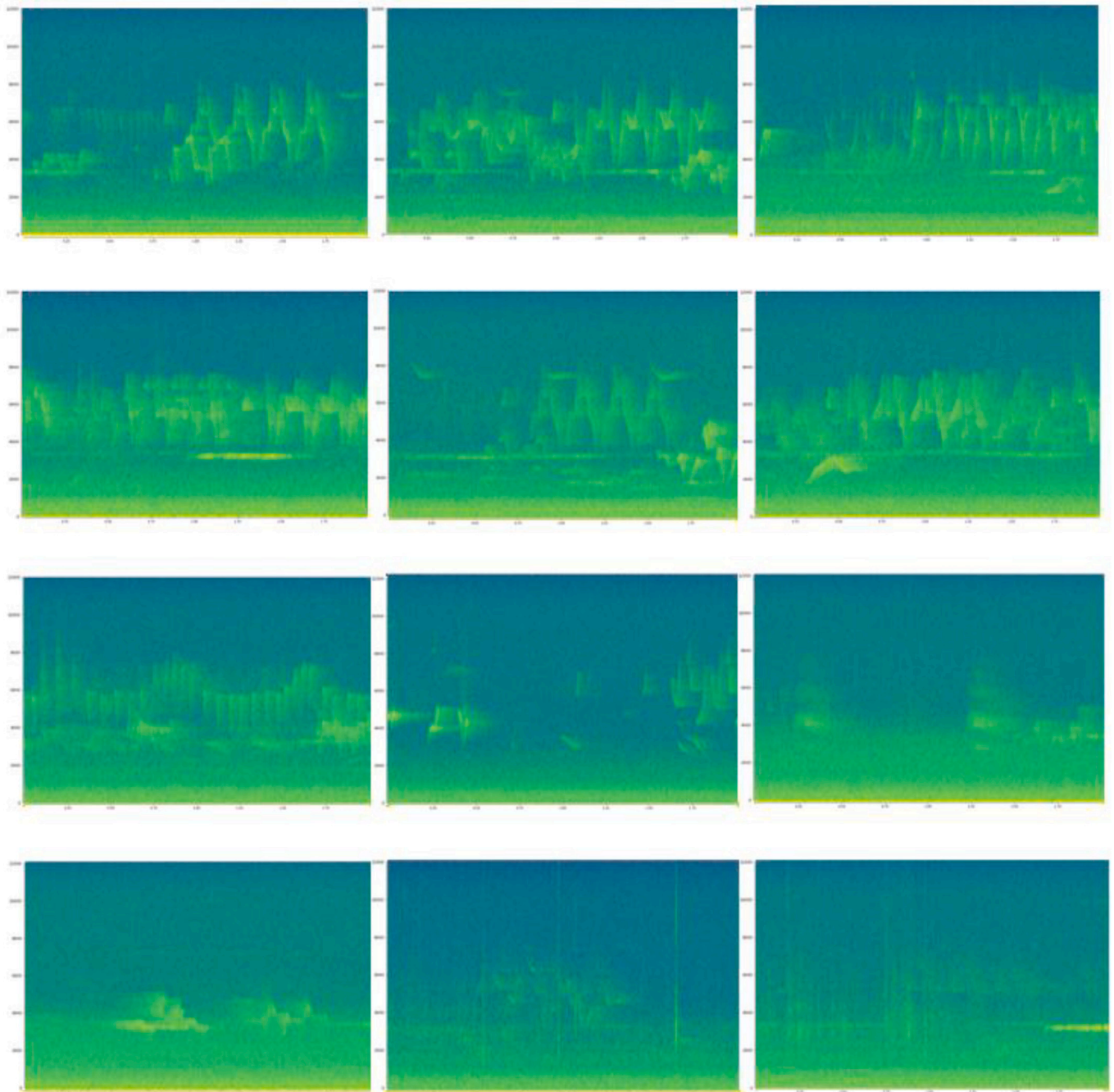
Visualization of these big data results is a helpful tool for data analysis, as well as further verification and spot checking of results. Utilizing Plotly Dash (<https://plotly.com/dash/>), a web interface was created to visualize daily and hourly count (Fig. 4), with interactive options to filter results by species, predicted probability range (threshold), model iteration, and station.

## 5. Discussion

In this study, we demonstrate how deep convolutional neural networks (CNN) and transfer learning can achieve higher accuracy for the classification of calls from the targeted rare species with limited training data. We provide both methodological and practical contributions by testing the performance of a machine learning approach to augment the manual validation process, which is time-consuming and labor-intensive.

With limited labeled data, especially for rare species, the CNN model performs reasonably well. While transfer learning leverages the learning from one task which is generally trained on a large size dataset, it does not require learning from scratch for the new task, which is motivated by the observation that the earlier features of a CNN model contain more generic features (e.g. edge detectors or color blob detectors) that should be useful for many tasks. In this study, we used a pre-trained ResNet50 model to implement transfer learning with fine-tunings, and there are other options of pre-trained CNN models, such as VGG16 or DenseNet (Huang et al., 2017), that can be used to achieve comparable results. Except for these pre-trained models, which are based on ImageNet, transferring learned knowledge from networks trained on audio data (for example, SoundNet (Aytar et al., 2016) or SincNet (Ravanelli and Bengio, 2018)) is another reasonable choice.

Data augmentation is another effective way to increase the training sample size in order to achieve better classification performance. Beyond the ones that we used in our model, there are more complicated



**Fig. 3.** Examples of spectrograms for 4-s audio recordings with detected calls from the species *Spelaeornis caudatus*. Row 1–2: examples of detections that the model can correctly classify; Row 3–4: examples of detections that the model wrongly classified as “no call”.

**Table 2**  
Number of model results returned for three selected probability ranges.

Species	Predicted Probability Range	# of 2-s clips ML results show species presence
<i>Spelaeornis caudatus</i>	0.99–1	240,300 clips
	0.7–1	982,230 clips
	0.5–1	1,247,170 clips
<i>Phylloscopus cantator</i>	0.99–1	51,550 clips
	0.7–1	189,380 clips
	0.5–1	237,260 clips

data augmentation methods such as adding or removing noises, image sharpening or masking, changing audio loudness, and audio mixing,

Finally, the methodology and implementation framework presented in this study can be easily adopted by other similar bioacoustics applications, where target signals require manual validation. This study sets initial steps for placing deep learning CNN analysis as the natural evolution of analysis methods for passive acoustic monitoring data.

### 5.1. Further research

In order for the results to be accurately used for species presence survey data, more iterations of label verification and model retraining are needed. The next step for this research is to establish a pipeline for



Fig. 4. Web interface that can visualize the number of detected calls in multiple monitoring stations over time for certain targeted species. The interface allows the users to choose different probability ranges from model predictions.

verifying the ML results, determining when to re-run the model with additional verified training data, and ultimately choosing a threshold per species that represents accurate species presence survey data.

One tool that will aid this verification is being added to the interface and will be tested with further research. 10% stratified sample of the results will be returned for experts to spot check and compare with model analysis in order to determine if the model needs to be retrained or if the results are accurate for species presence research. A framework for this verification is essential because each species call will require different amounts of training data and/or a different threshold that returns accurate results.

#### Declaration of Competing Interest

None.

#### Acknowledgements

The authors would like to thank everybody who participated in the experiment for their support. This work was supported by AI for Earth grants at Microsoft. Our appreciation to Dan Morris for connecting different parties for fruitful discussions and useful online materials. Our gratitude to the Nepal Government Department of National Parks and Wildlife Conservation, Makalu Barun National Park, The East Foundation (TEF), and the Barun Bachaon (“Save the Barun”) Taskforce for their partnership.

#### References

- Aide, T.M., Corrada-Bravo, C., Campos-Cerqueira, M., Milan, C., Vega, G., Alvarez, R., 2013. Real-time bioacoustics monitoring and automated species identification. *PeerJ* 1, e103. <https://doi.org/10.7717/peerj.103>.
- Aytar, Y., Vondrick, C., Torralba, A., 2016. *SoundNet: Learning Sound Representations from Unlabeled Video*. NIPS.
- BirdLife International. IUCN Red List for Birds. Downloaded from <http://www.birdlife.org> on 4/23/2020. 2020.
- Colonna, J.G., Cristo, M., Júnior, M.S., Nakamura, E.F., 2015. An incremental technique for real-time bioacoustic signal segmentation. *Expert Syst. Appl.* 42, 7367–7374. <https://doi.org/10.1016/j.eswa.2015.05.030>.

- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., Fei-Fei, L., 2009. *ImageNet: A Large-Scale Hierarchical Image Database*. CVPR.
- Frommolt, K.-H., 2017a. Information obtained from long-term acoustic recordings: applying bioacoustic techniques for monitoring wetland birds during breeding season. *J. Ornithol.* 158, 1–10. <https://doi.org/10.1007/s10336-016-1426-3>.
- Frommolt, K.-H., 2017b. Information obtained from long-term acoustic recordings: applying bioacoustic techniques for monitoring wetland birds during breeding season. *J. Ornithol.* 158, 1–10. <https://doi.org/10.1007/s10336-016-1426-3>.
- Furnas, B.J., Callas, R.L., 2015a. Using automated recorders and occupancy models to monitor common forest birds across a large geographic region. *J. Wildl. Manag.* 79, 325–337. <https://doi.org/10.1002/jwmg.821>.
- Furnas, B.J., Callas, R.L., 2015b. Using automated recorders and occupancy models to monitor common forest birds across a large geographic region. *J. Wildl. Manag.* 79, 325–337. <https://doi.org/10.1002/jwmg.821>.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. *Deep Residual Learning for Image Recognition*. CVPR, pp. 770–778.
- Hill, A.P., Prince, P., Piña Covarrubias, E., Patrick Doncaster, C., Snaddon, J.L., Rogers, A., 2017. *AudioMoth: evaluation of a smart open acoustic device for monitoring biodiversity and the environment*. *Methods Ecol. Evol.* 9, 1199–1211. <https://doi.org/10.1111/2041-210X.12955>.
- Huang, G., Liu, Z., van der Maaten, L., Weinberger, K.Q., Densely, K.Q., 2017. *Connected Convolutional Networks*. CVPR.
- Huh, M., Agrawal, P., Efos, A.A., 2016. What makes imagenet good for transfer learning? *arXiv (1608.08614)*.
- Inskipp, C., Baral, H.S., Phuyal, S., Bhatt, T.R., Khatiwada, M., Inskipp, T., Khatiwada, A., Gurung, S., Singh, P.B., Murray, L., Poudyal, L., Amin, R., 2016. The status of Nepal's Birds: the national red list series. *Zool. Soc. Lond., UK*.
- Keen, S., Ross, J.C., Griffiths, E.T., Lanzone, M., Farnsworth, A., 2014. A comparison of similarity-based approaches in the classification of flight calls of four species of north American wood-warblers (Parulidae). *Ecol. Inform.* 21, 25–33.
- Knight, E., Hannah, K., Foley, G., Scott, C., Brigham, R., Bayne, E., 2017. Recommendations for acoustic recognizer performance assessment with application to five common automated signal recognition programs. *Avian Conserv. Ecol.* 12 (2), 14. <https://doi.org/10.5751/ACE-01114-120214>.
- Krizhevsky, A., Sutskever, I., Hinton, G., 2012. *Imagenet Classification with Deep Convolutional Neural Networks*. NIPS.
- LeCun, Y., 2015. *Lenet-5, convolutional neural networks*. URL: <http://yann.lecun.com/exdb/lenet>.
- Matsubayashi, S., Suzuki, R., Saito, F., Murate, T., Masuda, T., Yamamoto, K., Okuno, H. G., 2017. Acoustic monitoring of the great reed warbler using multiple microphone arrays and robot audition. *J. Robot. Mechatron.* 29, 224–235.
- Mesaros, A., Heittola, T., Benetos, E., Foster, P., Lagrange, M., Virtanen, T., Plumbley, M. D., 2018. Detection and classification of acoustic scenes and events: outcome of the DCASE. 2016 challenge. *IEEE/ACM Trans. Audio Speech Lang. Process.* 26, 379–393. <https://doi.org/10.1109/TASLP.2017.2778423>.
- Ravanelli, M., Bengio, Y., 2018. *Speaker recognition from raw waveform with sinchnet*. *arXiv (1808.00158)*.

- Salamon, J., Bello, J.P., 2017. Deep convolutional neural networks and data augmentation for environmental sound classification. *Sig. Proces. Lett. (SPL)* 24 (3), 279–283.
- Shin, H.-C., Roth, H.R., Gao, M., Lu, L., Xu, Z., Nogues, I., Yao, J., Mollura, D., Summers, R.M., 2016. Deep convolutional neural networks for computer-aided detection: CNN architectures, dataset characteristics and transfer learning. *IEEE Trans. Med. Imaging* 35, 1285–1298.
- Shorten, C., Khoshgoftaar, T.M., 2019. A survey on image data augmentation for deep learning. *J. Big Data* 6, 60.
- Simonyan, K., Zisserman, A., 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. *ICLR*.
- Stuart, S.N., Chanson, J.S., Cox, N.A., Young, B.E., Rodrigues, A.S.L., Fischman, D.L., Waller, R.W., 2004. Status and trends of amphibian declines and extinctions worldwide. *Science* 306, 1783–1786.
- Towsey, M., Planitz, B., Nantes, A., Wimmer, J., Roe, P., 2012. A toolbox for animal call recognition. *Bioacoustics* 21, 107–125. <https://doi.org/10.1080/09524622.2011.648753>.
- Wildlife Acoustics, 2019. Kaleidoscope Pro Analysis Software. Boston, MA. Available at: <https://www.wildlifeacoustics.com/products/kaleidoscope-pro>.
- Zhao, Z., Zhang, S.-h., Xu, Z.-y., Bellisario, K., Dai, N.-h., Omrani, H., Pijanowski, B.C., 2017. Automated bird acoustic event detection and robust species classification. *Ecol. Inform.* 39, 99–108.