

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/347310404>

# Bioacoustics and machine learning for automated avian species monitoring in global biodiversity hotspots

Article in *The Journal of the Acoustical Society of America* · October 2020

DOI: 10.1121/1.5146736

CITATIONS

0

READS

141

7 authors, including:



Ming Zhong

Microsoft

14 PUBLICATIONS 233 CITATIONS

SEE PROFILE



Naomi Bates

Future Generations University

5 PUBLICATIONS 20 CITATIONS

SEE PROFILE



Rahul Dodhia

Microsoft

48 PUBLICATIONS 706 CITATIONS

SEE PROFILE



Juan Miguel Lavista Ferres

Microsoft

115 PUBLICATIONS 819 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



EDC Oceans of Data Institute and IBM [View project](#)



Risks of Using Non-verified Open Data: A case study on using Machine Learning techniques for predicting Pregnancy Outcomes in India [View project](#)

# 1 **Acoustic Detection of Regionally Rare Bird Species Through** 2 **Deep Convolutional Neural Networks**

3 Ming Zhong<sup>1</sup>, Ruth Taylor<sup>2</sup>, Naomi Bates<sup>2</sup>, Damian Christey<sup>2</sup>, Hari Basnet<sup>2</sup>, Jennifer Flippin<sup>2</sup>,  
4 Shane Palkovitz<sup>2</sup>, Rahul Dodhia<sup>1</sup>, Juan Lavista Ferres<sup>1</sup>

5 <sup>1</sup> AI for Good Research Lab, Microsoft

6 <sup>2</sup> Songs of Adaptation, Future Generations University

7

8

9 Corresponding Author:

10 Juan Lavista Ferres<sup>1</sup>

11 AI for Good Research Lab, Microsoft, Redmond, WA 98052, USA

12 Email address: [jlavista@microsoft.com](mailto:jlavista@microsoft.com)

13

14 **Abstract:** Bioacoustic monitoring with machine learning (ML) models can provide valuable  
15 insights for informed decision-making in conservation efforts. In this study, the team built deep  
16 convolutional neural networks to analyze field recordings and classify calls of Yellow-vented  
17 warbler (*Phylloscopus cantator*) and Rufous-throated wren-babbler (*Spelaeornis caudatus*), both  
18 of which are regionally rare in Nepal. Data augmentation techniques for calls of the two bird  
19 species were utilized to effectively increase the size of the training set and thus boost model  
20 performance. Nepali ornithologists were engaged in iterative data labeling from field recordings,  
21 leveraging ML technology in conjunction with expert manual labeling and verification. The  
22 model output provides insights of species activity and abundance throughout 2018-2019 in  
23 multiple ecosystems along an elevational transect in the Barun River Valley, Nepal. The results  
24 of this study may help conservationists better understand species distribution, behavior, diversity,  
25 and habitat preference. Additionally, the results provide baseline data to quantify future changes

26 due to habitat disruption or climate change. This modeling methodology and its framework can  
27 be easily adopted by other acoustic classification problems.

28 **Keywords:** Deep learning, Convolutional Neural Networks (CNN), Bioacoustic classification,  
29 Transfer learning, Species population, Presence survey

## 30 **I. Introduction**

31 In recent decades, the populations of various animals, including birds, amphibians, insects, and  
32 mammals, have exhibited steep declines worldwide. While many decreases are due to habitat  
33 loss and overutilization, other unidentified processes threaten 48% of rapidly declining species  
34 and are driving species most quickly to extinction [1]. As biodiversity plays a critical role in  
35 many aspects, well-designed monitoring programs provide a basis for identifying the species,  
36 sites and threats of most significant concern. Such monitoring programs also provide reliable  
37 tools when evaluating the integrity of ecosystems and their responses to disturbances, assessing  
38 progress in efforts to conserve biodiversity, and measuring the success of actions taken to  
39 preserve or recover biodiversity. However, manual observation remains limited and challenging  
40 in many scenarios, especially in the areas that are difficult to access physically or when the focus  
41 is to study animals' night-time behavior. In such scenarios, passive acoustic monitoring is highly  
42 appropriate, as many birds, including rare species, are most readily detectable by their sounds,  
43 often more so than by vision. With modern remote monitoring stations, it can continuously  
44 monitor large remote areas for avian community composition and tracking migratory and  
45 seasonal changes in populations ([2] – [7]).

46 Earlier applications that have employed such technology either performed automatic recording  
47 but relied on manual analysis of sound recordings ([8], [9]) or were based on low-complexity

48 signal processing such as template matching ([10], [11]), feature extraction ([12]), or traditional  
49 machine learning methods ([13], [14]).

50 With the constant increase in computing power and the development of more efficient codes,  
51 high-performance computing helps the extremely fast growth of deep learning in recent years,  
52 which has been shown to outperform previous state-of-the-art techniques in several tasks. Deep  
53 learning has fueled great strides in a variety of computer vision problems, and in particular,  
54 Convolutional Neural Networks (CNN) have demonstrated great potential and success in image  
55 classification tasks and thus drawn much attention in constructing the automatic bird sound  
56 classification systems. Some popular CNN architectures applied to bioacoustics classification  
57 include AlexNet [15], LeNet-5 [16], VGG16 [17], ResNet50 [18], among others.

58 In this study, two regionally rare species were chosen: Yellow-vented warbler (*Phylloscopus*  
59 *cantator*) and Rufous-throated wren-babbler (*Spelaeornis caudatus*). The Rufous-throated Wren  
60 Babbler is a very rare bird that has an extremely limited range in Nepal. The species is Near  
61 Threatened globally; it is listed within Nepal as a Critically Endangered species on a national  
62 level ([19], [20]). The nationally endangered Yellow-vented Warbler can be found in the East of  
63 Nepal. It is recorded between 75m and 1525m in a few locations, including Makalu Barun  
64 National Park [20]. These species provided a proof of concept demonstrating that with limited  
65 training samples, deep learning models can classify rare species calls.

66 As a research project of Future Generations University, this project brings expertise in  
67 community development with decades-long global partnerships that ensure long-term data  
68 collection and research permissions, data labeling, and collaboration for sustainable, just, and  
69 lasting climate action. Protecting 100,000,000 acres of land, Future Generations leadership in

70 community-based conservation established multiple national parks across Asia. Over the past 27  
71 years, the University has employed key indicators (quick, easy-to-use measurements), shaped to  
72 fit specific communities' contexts, to empower community members to measure change over  
73 time for themselves. This project is unique in its commitment to community engagement.

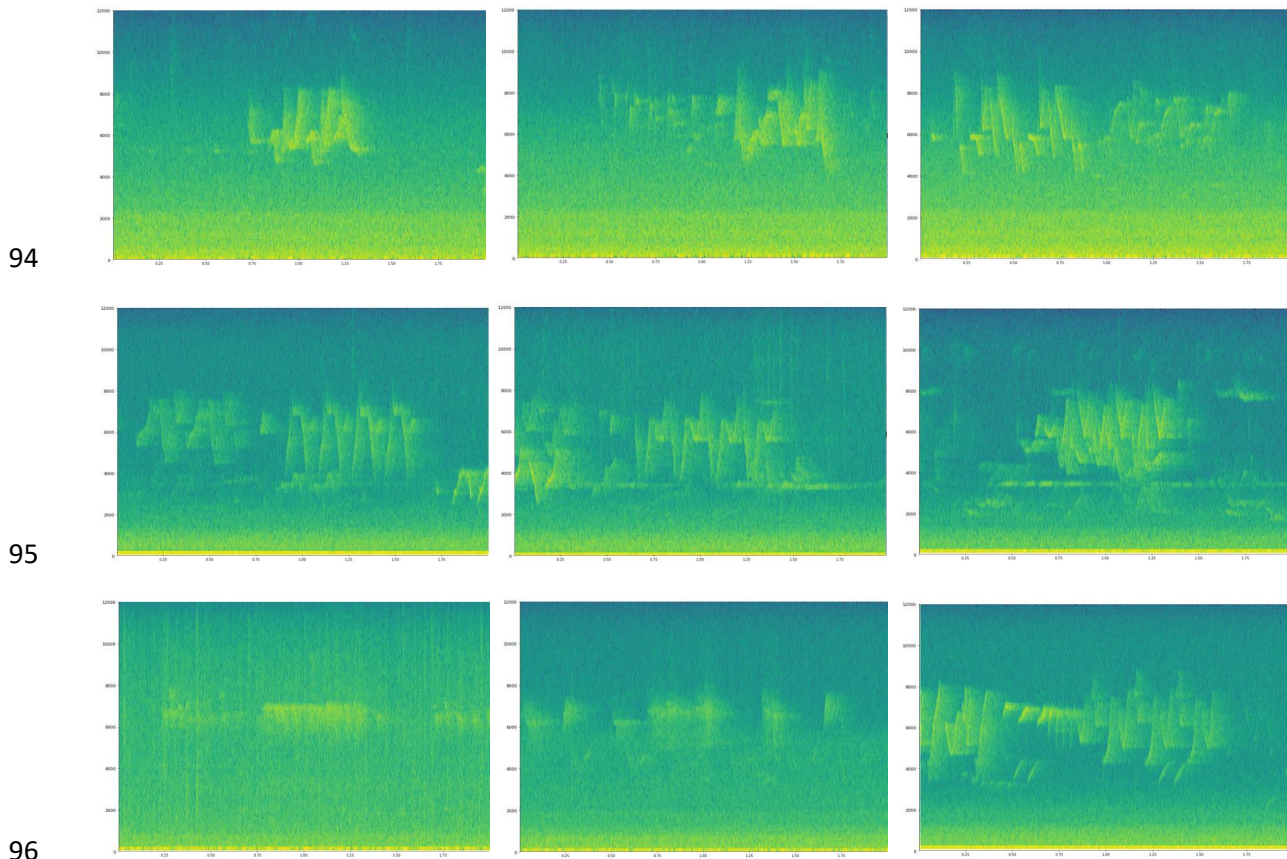
## 74 **II. Data collection and pre-processing**

75 Audio data were collected from 8 different sites along an elevational transect in the mountains of  
76 Makalu Barun National Park, Nepal, between 2018 and 2019. Audios were recorded into wav  
77 format using Song Meter SM4 Acoustic Recorders (Wildlife Acoustics) at a sampling rate of  
78 48kHz and 24-bit rate, and recorders were programmed to record 5-minute audio every 15 minutes  
79 24 hours per day.

80 The training and test datasets were initially generated using pattern recognition clustering software  
81 (Kaleidoscope Pro Analysis Software, Wildlife Acoustics [21]), and local avian experts  
82 subsequently analyzed clusters in Nepal to identify calls from the two species of interest, Yellow-  
83 vented warbler (*Phylloscopus cantator*) and Rufous-throated wren-babbler (*Spelaeornis*  
84 *caudatus*). A target of 100+ positive detections for each species and 300+ negative detections (ex.  
85 rain wind river, insects, other bird species, etc.) were set as the minimum amount of data necessary  
86 for model training and development. The positive and negative samples were used as input to train  
87 CNN models. Spectrogram images containing a target positive or negative sample were  
88 standardized using a 4-second audio clip beginning at each detection's start-time.

89 Spectrograms were extracted from audio files (with NFFT = 256, Hanning window) using Python  
90 3.6 and then resized to 224 by 224 pixels with RGB channels and stored as color PNG images (see  
91 Fig. 1 for example). The color spectrograms were the input for the machine learning model, and

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



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

### 100 III. Approaches

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

102 Here the neural network model ResNet50 was applied to classify the calls of the 2 bird species.  
103 This ResNet50 CNN architecture is a variant of ResNet model which has 48 Convolution layers  
104 along with 1 Max Pooling and 1 Average Pooling layer. It begins with the RGB images (size 224

105  $\times 224 \times 3$ ) as input and performs the initial convolution and max-pooling using  $7 \times 7$  and  $3 \times 3$   
106 kernel sizes, respectively. Afterward, it stacks a series of residual blocks. With the skip  
107 connection of residual blocks, it allows the model to propagate larger gradients to initial layers.  
108 These layers are able to learn as fast as the final layers, in order to train deeper networks. Finally,  
109 the network has an average pooling layer, followed by a fully connected layer. When training the  
110 ResNet50 model, the Adam optimizer algorithm was applied, and an initial learning rate of  $1e-4$   
111 with a decay factor of  $1e-7$ .

112 In the context of deep learning, most models include millions of parameters. ResNet50, for  
113 example, has 23 million parameters. To train such complex models, it typically requires an  
114 extensive dataset to achieve an optimal parameter configuration. However, in practice it may be  
115 very difficult to collect large amounts of labeled data, especially if a species rarely calls or if the  
116 species is endangered and there are few individuals. Besides, using experts to obtain a large  
117 number of labeled samples in acoustics is an expensive and time-consuming endeavor. Given  
118 this scenario, transfer learning with fine-tuning [22] is a useful technique when there is only a  
119 small number of labeled data available.

120 Transfer learning is a machine learning technique where a model trained on one task (or domain)  
121 is re-purposed on a second related task (or domain). Pre-trained models are usually shared in the  
122 form of the millions of parameters/weights the model achieved while being trained to an optimal  
123 state. In this study, the model weights were initially trained on ImageNet [23] dataset with 1000  
124 classes of objects, but their pre-trained weights can be leveraged by a different task or domain  
125 [24]. This approach is effective because the source model was trained on a large number of  
126 images and made predictions on a relatively large number of classes. In turn, it required the

127 model to extract distinct features from images in order to perform well. With fine-tuning, some  
128 layers are frozen from the pre-trained model, and it is sufficient to train the last several layers  
129 only, instead of having to train the whole model with random initialization of all parameters.

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

### 132 *B. K-Fold Cross-Validation*

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

142 Step 1: Split the entire data randomly into K folds (here, we use  $K=5$ ).

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

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

### 147 *C. Data Augmentation*



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

155 Among various data augmentation methods for image processing, some basic ones include flips,  
156 rotations, shifts, noise injections, color space transformations, sharpening or blurring, and  
157 random erasing or cropping. Specifically, for audio recordings, there are methods such as time-  
158 stretching, pitch shifting, and mixing multiple audios [25]. Beyond them, there are more  
159 advanced techniques, for example, generative adversarial network (GAN)-based methods [26],  
160 which can be used to generate synthetic images.

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

164 Another effective method we adopted to boost the training data size is to use spectrograms with  
165 smaller time-windows. While the detected calls for the two regionally rare species, *Phylloscopus*  
166 *cantator* and *Spelaeornis caudatus*, usually last for 2 seconds or longer (see Fig. 1. as an  
167 example), our baseline model was fit based on spectrograms generated from a 4-second time  
168 window. In order to boost the size of training data, we break down each 4-second detection into  
169 3 shorter detections, where each detection lasts for 2 seconds (that is, to create 3 spectrograms

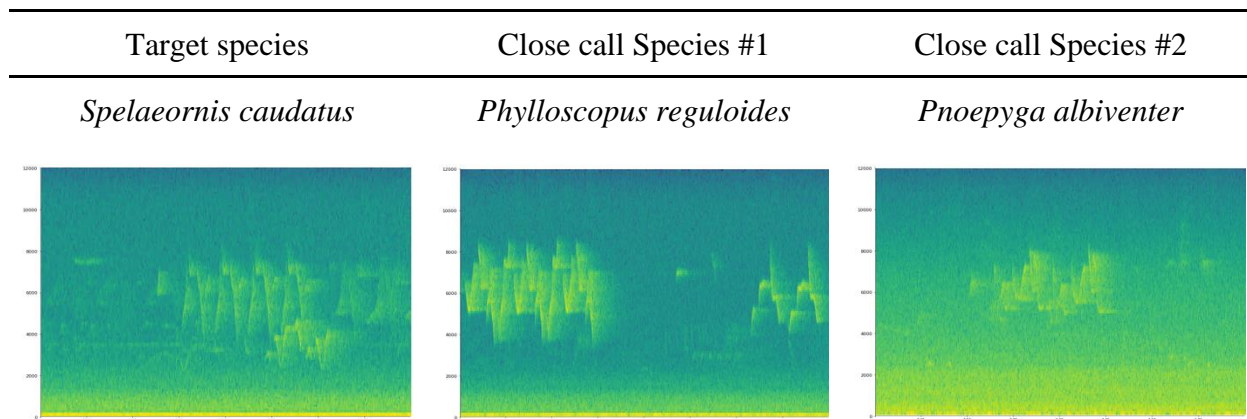
170 starting at second 0, 1, and 2, respectively, from each original 4-second detection). Even though  
 171 breaking down spectrograms into 2-second windows may not include one complete call within  
 172 each spectrogram and may bring some noisy labels during model training, but with this  
 173 implementation, the size of data available for training tripled.

174 *D. Model Training*

175 Our manually validated dataset consists of 195 positive detections for *Phylloscopus cantator*,  
 176 320 positive detections for *Spelaeornis caudatus*, and 1060 negative detections composed of  
 177 various types of noises (rain, wind, river, bugs, other bird species, etc.), where each detection  
 178 lasts for 4 seconds.

179 Finding sufficient clips of exemplar training data from the field recordings was challenging,  
 180 because of call volume variations, overlapping calls with other species, and background noises.

181 In addition, to distinguish a species with multiple and varying calls, it was also essential and  
 182 challenging to determine other species that had similar calls to the target species and label these  
 183 close calls as negative training data. Particularly, for the target species *Spelaeornis caudatus*,  
 184 there are two other species that have acoustically similar calls (Fig. 2).



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

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

## 191 **IV. Results**

### 192 *A. Model Performance*

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

200 **TABLE I:** Classification results (sensitivity, specificity, and AUC) for both target species by  
201 each CNN model. The results are based on the average score of conducting 5-fold cross-  
202 validation, with a neutral threshold score 0.5.

203

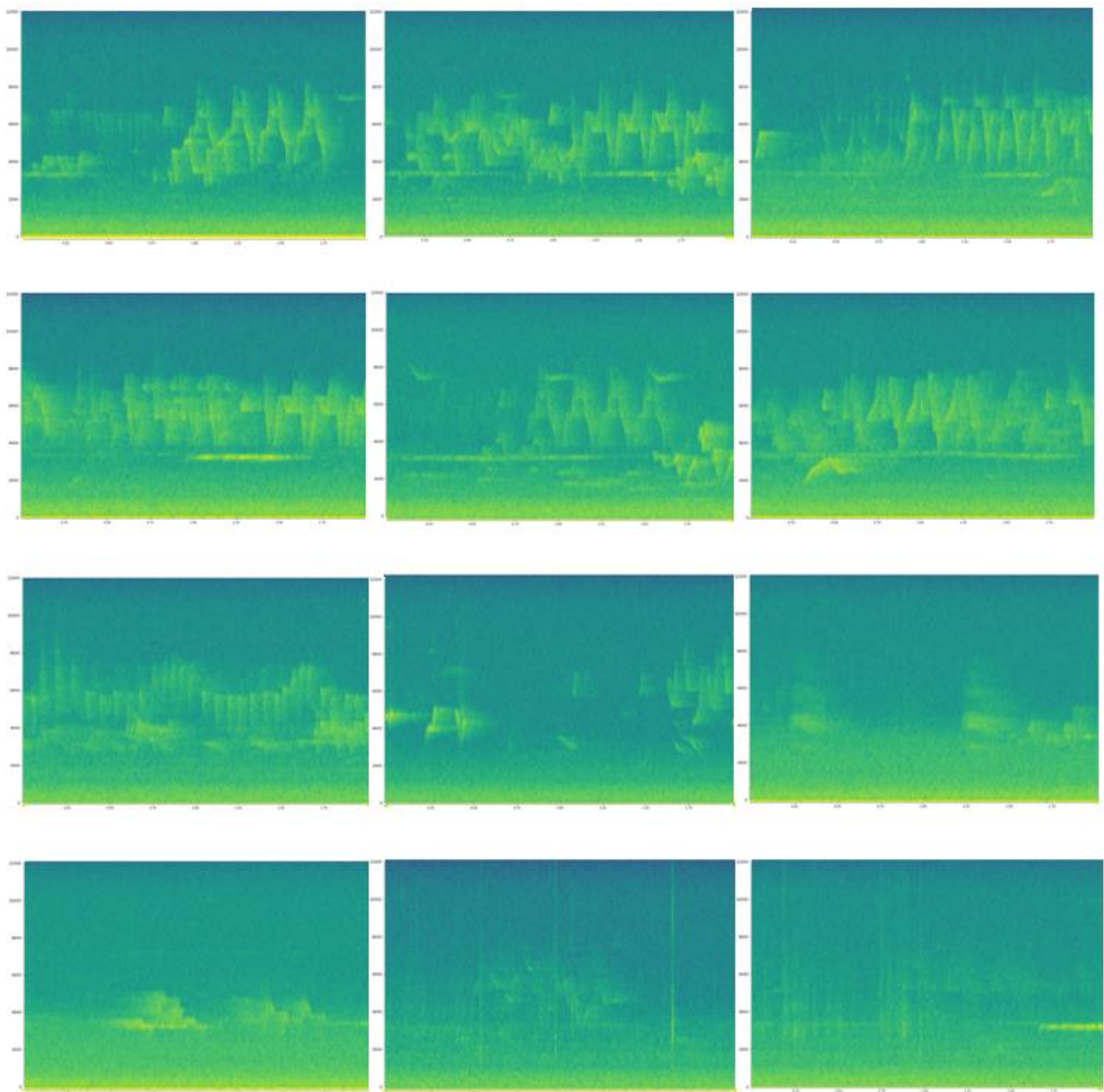
204

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

206

207 For both species *Phylloscopus cantator* and *Spelaeornis caudatus*, the model based on 2-second  
208 spectrograms performed significantly better, especially sensitivity, compared to the model based  
209 on 4-second spectrograms. Using data augmentation made further improvement for the model  
210 based on 2-second spectrograms (Table I). The sensitivity for classifying the species *Spelaeornis*

211 *caudatus* was not as good as that of the model classifying the species *Phylloscopus cantator*, and  
212 resulted in about 10% of detections that were misclassified as negative. A closer investigation of  
213 the data revealed that the labeled calls for *Spelaeornis caudatus* included detections with various  
214 levels of clarity, different call stereotypes, and maybe some incorrectly labeled detections. It  
215 appears that the neural network model did not find enough commonalities among these detected  
216 calls to make correct classification. Some examples of spectrograms are shown in Fig. 3.



217

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

### 221 *B. Scoring on Unlabeled Data*

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

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

234 Table II. Number of model results returned for three selected probability ranges.

Species	Predicted Probability Range	# of 2-sec 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

---

235

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



240

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

## 244 **V. Discussion**

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

250 With limited labeled data, especially for rare species, the CNN model performs reasonably well.  
251 While transfer learning leverages the learning from one task which is generally trained on a large  
252 size dataset, it does not require learning from scratch for the new task, which is motivated by the  
253 observation that the earlier features of a CNN model contain more generic features (e.g. edge  
254 detectors or color blob detectors) that should be useful for many tasks. In this study, we used a  
255 pre-trained ResNet50 model to implement transfer learning with fine-tunings, and there are other  
256 options of pre-trained CNN models, such as VGG16 or DenseNet ([27]), that can be used to  
257 achieve comparable results. Except for these pre-trained models, which are based on ImageNet,  
258 transferring learned knowledge from networks trained on audio data (for example, SoundNet  
259 ([28]) or SincNet ([29])) is another reasonable choice.

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



262 more complicated data augmentation methods such as adding or removing noises, image  
263 sharpening or masking, changing audio loudness, and audio mixing,  
264 Finally, the methodology and implementation framework presented in this study can be easily  
265 adopted by other similar bioacoustics applications, where target signals require manual  
266 validation. This study sets initial steps for placing deep learning CNN analysis as the natural  
267 evolution of analysis methods for passive acoustic monitoring data.

## 268 **VI. Further Research**

269 In order for the results to be accurately used for species presence survey data, more iterations of  
270 label verification and model retraining are needed. The next step for this research is to establish a  
271 pipeline for verifying the ML results, determining when to re-run the model with additional  
272 verified training data, and ultimately choosing a threshold per species that represents accurate  
273 species presence survey data.

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

## 280 **Acknowledgements**

281 The authors would like to thank everybody who participated in the experiment for their support.  
282 This work was supported by AI for Earth grants at Microsoft. Our appreciation to Dan Morris for

283 connecting different parties for fruitful discussions and useful online materials. Our gratitude to  
284 the Nepal Government Department of National Parks and Wildlife Conservation, Makalu Barun  
285 National Park, The East Foundation (TEF), and the Barun Bachaon (“Save the Barun”)  
286 Taskforce for their partnership.

## 287 **References**

288 [1] Stuart, S.N., J.S. Chanson, N.A. Cox, B.E. Young, A.S.L. Rodrigues, D.L. Fischman, and  
289 R.W. Waller. Status and trends of amphibian declines and extinctions worldwide. *Science*, 2004,  
290 306, 1783–1786.

291 [2] Aide, T. M., Corrada-Bravo, C., Campos-Cerqueira, M., Milan, C., Vega, G., & Alvarez, R.  
292 Real-time bioacoustics monitoring and automated species identification. *PeerJ*, 2013, 1, e103.

293 <https://doi.org/10.7717/peerj.103>

294 [3] Furnas, B. J., & Callas, R. L. Using automated recorders and occupancy models to monitor  
295 common forest birds across a large geographic region. *The Journal of Wildlife Management*,  
296 2015, 79, 325–337. <https://doi.org/10.1002/jwmg.821>

297 [4] Frommolt, K.-H. Information obtained from long-term acoustic recordings: Applying  
298 bioacoustic techniques for monitoring wetland birds during breeding season. *Journal of*  
299 *Ornithology*, 2017, 158, 1–10. <https://doi.org/10.1007/s10336-016-1426-3>

300 [5] Hill, A. P., Prince, P., Piña Covarrubias, E., Patrick Doncaster, C., Snaddon, J. L., & Rogers,  
301 A. AudioMoth: Evaluation of a smart open acoustic device for monitoring biodiversity and the  
302 environment. *Methods in Ecology and Evolution*, 2017, 9, 1199–1211.

303 <https://doi.org/10.1111%2F2041-210x.12955>

- 304 [6] Knight, E., Hannah, K., Foley, G., Scott, C., Brigham, R., & Bayne, E. Recommendations for  
305 acoustic recognizer performance assessment with application to five common automated signal  
306 recognition programs. *Avian Conservation and Ecology*, 2017, 12(2), 14.  
307 <https://doi.org/10.5751/ACE-01114-120214>
- 308 [7] Matsubayashi, S., Suzuki, R., Saito, F., Murate, T., Masuda, T., Yamamoto, K., & Okuno, H.  
309 G. Acoustic monitoring of the great reed warbler using multiple microphone arrays and robot  
310 audition. *Journal of Robotics and Mechatronics*, 2017, 29, 224–235.
- 311 [8] Furnas, B. J., & Callas, R. L. Using automated recorders and occupancy models to monitor  
312 common forest birds across a large geographic region. *The Journal of Wildlife Management*,  
313 2015, 79, 325–337. <https://doi.org/10.1002/jwmg.821>
- 314 [9] Frommolt, K.-H. Information obtained from long-term acoustic recordings: Applying  
315 bioacoustic techniques for monitoring wetland birds during breeding season. *Journal of*  
316 *Ornithology*, 2017, 158, 1–10. <https://doi.org/10.1007/s10336-016-1426-3>
- 317 [10] Towsey, M., Planitz, B., Nantes, A., Wimmer, J., & Roe, P. A toolbox for animal call  
318 recognition. *Bioacoustics*, 2012, 21, 107–125. <https://doi.org/10.1080/09524622.2011.648753>
- 319 [11] Colonna, J. G., Cristo, M., Júnior, M. S., & Nakamura, E. F. An incremental technique for  
320 real-time bioacoustic signal segmentation. *Expert Systems with Applications*, 2015, 42, 7367–  
321 7374. <https://doi.org/10.1016/j.eswa.2015.05.030>
- 322 [12] Mesaros A, Heittola T, Benetos E, Foster P, Lagrange M, Virtanen T, Plumbley MD.

323 Detection and classification of acoustic scenes and events: outcome of the DCASE. 2016  
324 challenge. *IEEE/ACM Trans. Audio Speech Lang. Process*, 2018, 26, 379–393.  
325 (doi:10.1109/TASLP.2017.2778423)

326 [13] Keen, S., Ross, J. C., Griffiths, E. T., Lanzone, M., and Farnsworth, A. A comparison of  
327 similarity-based approaches in the classification of flight calls of four species of North American  
328 wood-warblers (Parulidae). *Ecological Informatics*, 2014, 21, 25–33.

329 [14] Zhao, Z., S.-h. Zhang, Z.-y. Xu, K. Bellisario, N.-h. Dai, H. Omrani, and B. C. Pijanowski.  
330 Automated bird acoustic event detection and robust species classification. *Ecological*  
331 *Informatics*, 2017, 39, 99–108.

332 [15] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional  
333 neural networks. In *NIPS*, 2012.

334 [16] Y. LeCun. Lenet-5, convolutional neural networks. URL: [http://yann. lecun.](http://yann.lecun.com/exdb/lenet)  
335 [com/exdb/lenet](http://yann.lecun.com/exdb/lenet), 2015.

336 [17] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image  
337 recognition. In *ICLR*, 2015.

338 [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In  
339 *CVPR*, 2016, pp. 770-778.

340 [19] BirdLife International. IUCN Red List for birds. Downloaded from <http://www.birdlife.org>  
341 on 4/23/2020. 2020.

- 342 [20] C. Inskipp, H.S. Baral, S. Phuyal, T.R. Bhatt, M. Khatiwada, T. Inskipp, A. Khatiwada, S.  
343 Gurung, P.B. Singh, L. Murray, L. Poudyal. & R. Amin. The status of Nepal's Birds: The  
344 national red list series. *Zoological Society of London*, UK, 2016.
- 345 [21] Wildlife Acoustics, Kaleidoscope Pro Analysis Software. Boston, MA. Available at:  
346 <https://www.wildlifeacoustics.com/products/kaleidoscope-pro>, 2019.
- 347 [22] H.-C. Shin, H. R. Roth, M. Gao, L. Lu, Z. Xu, I. Noguees, J. Yao, D. Mollura, and R. M.  
348 Summers. Deep convolutional neural networks for computer-aided detection: CNN architectures,  
349 dataset characteristics and transfer learning. *IEEE Transactions on Medical Imaging*, 2016, 35,  
350 pp. 1285-1298.
- 351 [23] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale  
352 Hierarchical Image Database. In *CVPR*, 2009.
- 353 [24] M. Huh, P. Agrawal, and A. A. Efros. What makes imagenet good for transfer learning?  
354 *arXiv:1608.08614*, 2016.
- 355 [25] J. Salamon, J. P. Bello, Deep convolutional neural networks and data augmentation for  
356 environmental sound classification, *Signal Processing Letters (SPL)* 24 (3) (2017) 279–283.
- 357 [26] Shorten, C., Khoshgoftaar, T.M. A survey on image data augmentation for deep learning.  
358 *Journal of Big Data*, 2019, 6, 60.
- 359 [27] Huang, G., Liu, Z., van der Maaten, L. and K. Q. Weinberger, K. Q. Densely Connected  
360 Convolutional Networks. In *CVPR*, 2017.
- 361 [28] Y. Aytar, C. Vondrick, and A. Torralba. SoundNet: Learning sound representations from  
362 unlabeled video. In *NIPS*, 2016.

363 [29] M. Ravanelli and Y. Bengio, Speaker recognition from raw waveform with sincnet.  
364 arXiv:1808.00158, 2018.