

# BirdNET-Annotator: AI-assisted Strong Labelling of Bird Sound Datasets

Anonymous submission

## Abstract

Monitoring biodiversity in biosphere reserves is challenging due to the vast regions to be monitored. Thus, conservationists have resorted to employing passive acoustic monitoring (PAM), which automates the audio recording process. Passive acoustic monitoring can create large, unlabeled datasets, but deriving knowledge from such recordings is usually still done manually. In the past, data science methods have been employed to generate low-dimensional summaries, so-called acoustic indices, which are however difficult to interpret even for human experts.

Today, machine learning enables us to detect vocalizations of species automatically, which means we can summarize the biodiversity in an area in terms of species richness. While there are powerful neural networks available for bird vocalization detection, they have not been adopted by taxonomists widely, who instead oftentimes still perform manual labeling of audio files from passive acoustic recordings. This is due to even the most complex audio analysis tools not employing deep learning models, because running them is expensive on servers and time-consuming on GPU-less client PCs.

In this paper, we present BirdNET-Annotator, a tool for AI-assisted labeling of audio datasets. BirdNET-Annotator was developed in co-design with taxonomists, who are coauthors of this publication. BirdNET-Annotator runs in the cloud free of charge, enabling taxonomists to scale beyond the limitations of their local hardware. We evaluated the performance of our solution in the context of its intended workflow and found a slight uplift in annotation times. While our results show that our application now meets the user requirements, there are still opportunities to seize for additional performance and usability improvement.

Our application demonstrates that large, pretrained neural models can be used by taxonomists when packaged in a user-friendly manner. We observe that although our solution adds a step to the preexisting workflow, the overall annotation speed is significantly increased. This hints at further improvement to be realized by integrating more steps of the workflow in fewer tools.

## Introduction

Passive acoustic monitoring (PAM) has become widely used in ecological research, as it is the most scalable data acquisition scheme (Sugai et al. 2018). This method involves surveying and monitoring wildlife and environments using sound recorders (acoustic sensors) deployed in the field

for extended periods to record acoustic data on a specified schedule. In contrast to active recording, PAM uses omnidirectional microphones, which capture sounds in any direction with equal sensitivity. Use cases for passive acoustic monitoring data include monitoring wildlife reserves, detecting poaching and observing not previously described species.

Machine learning (ML) has proven to be a powerful tool that can achieve great performance in virtually every data analysis task, thus enabling scalable automation for many knowledge work tasks. In contrast to classical programming, where the system behavior is derived from a bottom-up data model, machine learning derives the underlying structure from examples of the data during training, before artificial neural models can be used to make predictions. While machine learning models can be trained on unlabeled data to learn the data structure, associating meaning to data points requires labelled datasets, in which each example is labelled with its meaning in the language of the domain in question (Serrano 2021).

While machine learning systems for automated audio data analysis exist, their applicability is limited to datasets with strong labels (Cobos et al. 2022). While passive acoustic monitoring enables collecting large amounts of data as required by machine learning, labelling the data can quickly become the bottleneck in building datasets. Common approaches to the challenge labelling large datasets poses include crowdsourced, where anonymous users of a service contribute labels to examples. Such libraries include the Macaulay library<sup>1</sup>, Xeno-Canto<sup>2</sup>, and iNaturalist<sup>3</sup>. While the approach of spreading the labelling workload to many people works, their time is still valuable and assistance hence welcome. To address this need, we propose an automated tool which can assist in labelling audio files for the creation of labelled datasets from passive acoustic monitoring data.

## Related Work

Birds play a crucial role in soundscape ecology due to their prevalence in the recorded biophony. As birds feed on insects, nectar, and fruits, they are also part of many food

---

<sup>1</sup><https://www.macaulaylibrary.org/>

<sup>2</sup><https://xeno-canto.org/>

<sup>3</sup><https://inaturalist.org>

chains, making their species richness and abundance good proxies for observing ecosystems (Xie et al. 2023). Solutions for detecting birds from their sounds have advanced from classical pattern recognition to deep learning, significantly improving automatic recognition, e.g. (Potamitis 2016), (Eichinski et al. 2022). Current deep learning models for bird vocalization detection can discern more than 6,000 species, a significant portion of both the 10,000 described and the estimated 20,000 total bird species (Kahl et al. 2021), which demonstrates both the power of this approach and the need for more data to increase the species coverage.

Audio data labeling tools play a crucial role in streamlining the labeling process and ensuring the accuracy of labeled data, thereby impacting the overall efficiency and accuracy of a machine learning pipeline. A common workflow for labelling audio files is to open the audio file in a desktop application to play it back and annotate regions of its spectrogram representation with the relevant labels. The common tool for this workflow is still the open-source application Audacity, which provides a label track for annotations in the time-domain<sup>4</sup>. Extensions to the functionality existing in Audacity like caching the boundaries of selections, automatically inferring binary labels from one another have been proposed as a means of increasing usability. (Li, Burgoyne, and Fujinaga 2006). Also, replacing the audio editor Audacity with a wholly new tool that is specialized for the labelling has been suggested as a more thorough solution (Gibbons et al. 2023). Label Studio is a data labeling platform that has a similar feature set for audio, but also built-in automatic classification tools<sup>5</sup>.

Both in the interest of usability and in the interest of model performance it is advisable to identify a narrow context of use and domain of audio to consider. In this work, we propose therefore to add a web tool to the labelling workflow, which employs state-of-the-art machine learning models for automating the labelling step of building passive-acoustic bird sound datasets.

## Methods

We connected as computer science experts with taxonomy experts to form a consortium which could both identify and build the tool needed for leveraging machine learning for improving the labelling workflow. To guide the process of developing the tool iteratively, we followed Co-design, which refers to a process by which practitioners and researchers come together to either adapt or create and test new materials, exhibits, programs, or technology tools. (Mitchell et al. 2015) In our case, the taxonomy experts can be considered actual end users of our proposed application, as they are labelling files from datasets with gold-standard quality on a regular basis. The computer science researchers had been working with passive acoustic monitoring before and most recently built an application for automatic classification of bird song sound snippets, so they were also already familiar with the relevant technologies.

To bring knowledge from both taxonomy and computer science experts together, we conducted a series of semi-structured interviews, which were aimed at understanding the current workflow, its context of use, and the persistent pain points. It turned out that the taxonomists are labelling audio files in sessions of about one or two hours once to twice a week, as part of their jobs, in their offices at university or home. Next, we built a first version of a tool addressing the identified issues within the constraints of the identified context. As we identified installing the complex machine learning applications as a challenge to the taxonomists, we designed the tool as a web service that the users can access using their preexisting web browser. The taxonomy experts used the tool for one of their weekly annotation sessions and took notes on their observations, which were discussed with the computer science experts in a following meeting. The taxonomists' feedback was considered in a revision of the software artifact. This led to the addition of an integration with Xeno-Canto, which can be used as a ground-truth source for reference sounds of each detected species. Finally, we conducted a simple user study where the taxonomy experts would label some files both with and without our tool, yielding some promising quantitative results.

## Implementation

Packaging a complex machine learning application in a way that can easily be deployed to a client computer is still a challenge from a software engineering perspective. This is rooted in the machine learning applications usually being written in Python, which does not provide standard tools for packaging applications. We explored third-party solutions like PyInstaller<sup>6</sup>, which has been used with success for machine learning applications before, e.g. in the BirdNET-Analyzer<sup>7</sup>. However, this solution did not yield a deployable package for us. Instead, we choose to package our application as an OCI container image<sup>8</sup> using Google Kaniko<sup>9</sup> and publish it on Docker Hub<sup>10</sup>. This allowed us to deliver the application as a Docker image to a Hugging Face Space for easy and free perpetual hosting.<sup>11</sup> Our application is designed to be used interactively, so responses from the machine learning model usually take about 1 minute on 1-minute input files. Since our Hugging Face Space can easily be forked using the methods of the web interface of Hugging Face, anyone can create a copy for themselves and attach compute resources from Hugging Face's catalog<sup>12</sup> to it, should that be necessary.

Our Docker image is build by GitLab continuous integra-

<sup>4</sup>[https://manual.audacityteam.org/man/label\\_tracks.html](https://manual.audacityteam.org/man/label_tracks.html)

<sup>5</sup><https://labelstud.io/>

<sup>6</sup><https://pyinstaller.org/en/stable/>

<sup>7</sup><https://github.com/kahst/BirdNET-Analyzer>

<sup>8</sup><https://github.com/opencontainers/image-spec>

<sup>9</sup><https://github.com/GoogleContainerTools/kaniko>

<sup>10</sup><https://hub.docker.com/repository/docker/bengt/birdnet-annotator>

<sup>11</sup><https://huggingface.co/spaces/Bengt0/BirdNET-Annotator>

<sup>12</sup><https://huggingface.co/docs/hub/spaces-gpus>

tion<sup>13</sup>. It is based on Ubuntu 22.04<sup>14</sup>, installs Python 3.11<sup>15</sup>. Python runtime packages include the BirdNET-Analyzer wrapper `birdnetlib`<sup>16</sup>, the full-stack web toolkit `Gradio`<sup>17</sup>, the audio processing library `librosa`<sup>18</sup>, the interactive plotting library `Plotly`<sup>19</sup>, the audio resampling library `resampy`<sup>20</sup>, and the deep learning framework `TensorFlow`<sup>21</sup>.

## Results

Figure 1 shows the finished application from a user’s perspective. We made sure to handle all possible implications of interacting with the application so that it can be used in any order of the possible interaction steps. The intended interaction flow starts at the top where the user specifies the dataset to work on either by selecting from a dropdown of predefined locations or by entering the coordinates of the recording locations. This is necessary, as the BirdNET system filters its outputs using the eBird database<sup>22</sup> of local bird species. Next, the users select a recording of the dataset to work from their local machine, which in turn gets uploaded to the server. Note that this interaction scheme allows users to quickly iterate through the audio files of one dataset, which matches the workflow we identified in co-design with the taxonomy experts. The users can configure settings for the inference using the BirdNET Analyzer before a click on the “Detect” triggers the automated analysis. After about 1 minute, the system responds with a list of automatic detections, which can be selected for further inspection. Upon selecting one detection on the left side table, the right side updates to display a spectrogram, an audio player, an interactive label pull-down selection element and a link to Xeno-Canto. Using these, the user can visualize the sound and play it back to detect the bird call themselves. For reference, the user can open the relevant Xeno-Canto page right in the browser they are using the BirdNET-Annotator. The user can select another label from the searchable dropdown or enter a free-text label of their own. Once the user is happy with a label, they can click “Confirm” to add it to the output file. This process can be repeated until the high-confidence detections on this file are exhausted. Finally, the label track file can be downloaded to the user’s client computer by clicking “Download”. The downloaded file can be imported into Audacity for fine-tuning and completing the labels.

We evaluated our resulting application using a quantitative user study. To gather data, one of our taxonomy experts (f) labelled some ( $N = 46$ ) audio files, using Audacity either with or without our BirdNET-Annotator application as a preliminary stage. We made sure to randomize the files

so that difficulty rising and falling with the day and night cycle cannot skew results. The annotation effort took place over the course of several days, so that form-of-the-day effects should be eliminated. We found that our expert could label the files significantly faster when also using BirdNET-Annotator (M=00:10:48, SD=00:05:01) than with Audacity alone (M=00:06:52, SD=00:01:49),  $p = 0.001$ . This means that users of BirdNET-Annotator can save about 4 minutes (00:03:55) or about 35 % (36, 3 %) on an average file. This validates the overwhelmingly positive response that we got from the taxonomy experts, who reported being much faster when using BirdNET-Annotator.

## Conclusion and Future Work

In this paper, we introduced BirdNET-Annotator, a web-based tool that employs large, pretrained artificial neural models to semiautomatically annotate bird calls and bird songs in audio data files provided by the user. We have discussed how we developed, built, and hosted BirdNET-Annotator as well as how we evaluated it. BirdNET-Annotator has proven to be easily accessible, usable by taxonomists, and to improve user performance as measured in labelling time. In our exploratory effort to provide the field of annotating audio datasets with gold-standard labels, we have found a field that is on the one hand essentially necessary for machine learning algorithms to be even applicable. In our exchanges with the taxonomy experts, we have but also found a field that is in dire need for more automation of the tedious parts of the knowledge work, which can be provided by machine learning as we demonstrated by instance of our tool.

Although adding our tool, the BirdNET-Annotator, as a new stage to the annotation workflow complicates it considerably, we could show that the increased automation is worth it in terms of time saved. Still, integrating the workflow into one tool can foreseeably improve the workflow even more, because that would avoid the import, and export steps of the label track file. We explored building the labelling functionality as in Audacity right into the spectrogram plot of BirdNET-Annotator, but that did not work because Gradio does not expose functionality necessary for writing annotations and listening for user interactions with them. As a workaround, one could implement the functionality oneself, but that seems daunting because of the complexity of the Plotly element we used for this plot. Instead, it seems to us that the application needs to be reimplemented in a framework which supports this central requirement. We explored the possibility of rewriting BirdNET-Annotator in the plotly-aware full-stack web framework `Dash`<sup>23</sup> and could demonstrate that combining the spectrogram plot with the interactive labelling component does indeed work as expected. Given the potentially immense impact of better, AI-assisted tools for the taxonomy domain, this seems like a worthwhile endeavor.

<sup>13</sup><https://docs.gitlab.com/ee/ci/>

<sup>14</sup><https://ubuntu.com/>

<sup>15</sup><https://www.python.org/>

<sup>16</sup><https://pypi.org/project/birdnetlib/>

<sup>17</sup><https://pypi.org/project/gradio/>

<sup>18</sup><https://pypi.org/project/librosa/>

<sup>19</sup><https://pypi.org/project/plotly/>

<sup>20</sup><https://pypi.org/project/resampy/>

<sup>21</sup><https://pypi.org/project/tensorflow/>

<sup>22</sup><https://ebird.org/>

<sup>23</sup><https://plotly.com/dash/>

## References

- Cobos, M.; Ahrens, J.; Kowalczyk, K.; and Politis, A. 2022. An overview of machine learning and other data-based methods for spatial audio capture, processing, and reproduction. *EURASIP Journal on Audio, Speech, and Music Processing*, 2022(1): 1–21.
- Eichinski, P.; Alexander, C.; Roe, P.; Parsons, S.; and Fuller, S. 2022. A convolutional neural network bird species recognizer built from little data by iteratively training, detecting, and labeling. *Frontiers in Ecology and Evolution*, 10: 133.
- Gibbons, A.; Donohue, I.; Gorman, C.; King, E.; and Parnell, A. 2023. NEAL: an open-source tool for audio annotation. *PeerJ*, 11: e15913.
- Kahl, S.; Wood, C. M.; Eibl, M.; and Klinck, H. 2021. Bird-NET: A deep learning solution for avian diversity monitoring. *Ecological Informatics*, 61: 101236.
- Li, B.; Burgoyne, J. A.; and Fujinaga, I. 2006. Extending Audacity for Audio Annotation. In *ISMIR*, 379–380.
- Mitchell, V.; Ross, T.; May, A.; Sims, R.; and Parker, C. J. 2015. Empirical investigation of the impact of using co-design methods when generating proposals for sustainable travel solutions.
- Potamitis, I. 2016. Deep learning for detection of bird vocalisations. *arXiv preprint arXiv:1609.08408*.
- Serrano, L. 2021. *Grokking Machine Learning*. Simon and Schuster.
- Sugai, L. S. M.; Silva, T. S. F.; Ribeiro, J., José Wagner; and Llusia, D. 2018. Terrestrial Passive Acoustic Monitoring: Review and Perspectives. *BioScience*, 69(1): 15–25.
- Xie, J.; Zhong, Y.; Zhang, J.; Liu, S.; Ding, C.; and Triantafyllopoulos, A. 2023. A review of automatic recognition technology for bird vocalizations in the deep learning era. *Ecological Informatics*, 73: 101927.

Location

Specify the location of the dataset you are working on.

Preset

Fernando de Noronha, Brazil

Latitude

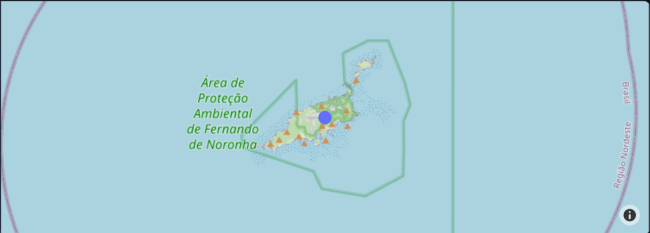
Latitude of the dataset.

-3.853825

Longitude

Longitude of the dataset.

-32.418656



Recording

Specify a recording to annotate.

0:00 / 1:00

Filename:

IMEIO1SET17\_20170922\_180000-0-100.wav

Date / Time:

2017-09-22 18:00:00

Inference

Specify how to automatically detect birds.

Confidence

Required confidence

0.01

Sensitivity

Sensitivity of detector

1

Overlap

Detection Window Overlap

0

Detect

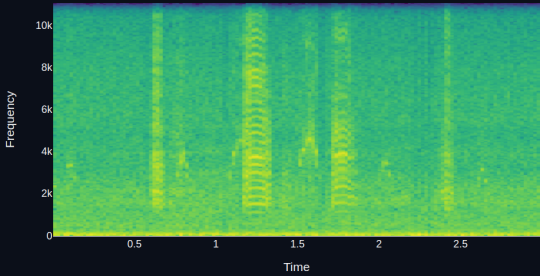
Detections

Select a detection for editing.

Common Name	Scientific Name	Start	End	Confidence
Masked Booby	Sula dactylatra	24	27	0.9958509206771851
Masked Booby	Sula dactylatra	21	24	0.9828821420669556
Masked Booby	Sula dactylatra	36	39	0.9823676347732544
Masked Booby	Sula dactylatra	9	12	0.9794644713401794
Masked Booby	Sula dactylatra	3	6	0.9661660194396973
Masked Booby	Sula dactylatra	0	3	0.8754146099090576
Masked Booby	Sula dactylatra	57	60	0.8594850301742554
Masked Booby	Sula dactylatra	18	21	0.8308070302009583
Masked Booby	Sula dactylatra	48	51	0.7499291896820068
Masked Booby	Sula dactylatra	6	9	0.6222408413887024
Great Crested Grebe	Podiceps cristatus	3	6	0.5404766798019409
Masked Booby	Sula dactylatra	33	36	0.402085542678833
Great Crested Grebe	Podiceps cristatus	15	18	0.2885759770870209
Masked Booby	Sula dactylatra	12	15	0.21965019404888153
Egyptian Goose	Alopochen aegyptiaca	36	39	0.1321738213300705
Ruddy Turnstone	Arenaria interpres	42	45	0.1316462904214859
Eurasian Wigeon	Mareca penelope	42	45	0.08782588690519333
Masked Booby	Sula dactylatra	15	18	0.07446610927581787
Great Crested Grebe	Podiceps cristatus	9	12	0.04598791152238846
Black-bellied Plover	Pluvialis squatarola	15	18	0.04422108829021454
Green-winged Teal	Anas crecca	3	6	0.04235334321856499
Eurasian Wigeon	Mareca penelope	15	18	0.040283240377902985
Black-bellied Plover	Pluvialis squatarola	42	45	0.03654693067073822
White-bellied Bustard	Eupodotis senegalensis	0	3	0.027458513155579567
Egyptian Goose	Alopochen aegyptiaca	24	27	0.02666369639337063
Great Crested Grebe	Podiceps cristatus	21	24	0.025333596393465996
Masked Booby	Sula dactylatra	27	30	0.02151246927678585
Masked Booby	Sula dactylatra	42	45	0.02027459815144539
Egyptian Goose	Alopochen aegyptiaca	9	12	0.01940160058438779

Detection

Edit the label of a detection.



0:00 / 0:03

Confidence

Confidence of automatic detection.

0.9958509

Name

Masked Booby · Sula dactylatra

[Xeno](#)  
[Canto](#)

Confirm

Start	End	Title
24	27	Masked Booby · Sula dactylatra
21	24	Masked Booby · Sula dactylatra
36	39	Masked Booby · Sula dactylatra
3	6	Great Crested Grebe · Podiceps cristatus

IMEIO1SET17\_20170922\_180000.tx...

Download

Figure 1: A screenshot of the BirdNET-Annotator application.