
Clean Data, Simple Models: A Practical Audio Preprocessing Approach for Multi-Species Sound Classification

Abstract

1 We present a lightweight audio-preprocessing pipeline that boosts simple classifiers
2 for multi-species sound identification in Colombian soundscapes. Developed for
3 BirdCLEF 2025 and evaluated on recordings from Reserva Natural El Silencio
4 (Magdalena Medio Valley), the pipeline isolates vocalizations, removes silence,
5 and filters noise to produce cleaner BirdNET embeddings. We train MLP and
6 CNN models on raw vs. cleaned inputs. Results in multi-taxon species classifica-
7 tion show that improving signal quality can offset model complexity, where a
8 cleaned-input MLP matches or surpasses deeper baselines with modest compute.
9 This underscores the value of preprocessing for bioacoustic monitoring in noisy,
10 resource-limited settings and demonstrates that robust baselines can be built with
11 accessible computing resources common in biodiversity-rich developing countries.

12 1 Introduction

13 Passive acoustic monitoring (PAM) has become a powerful and non-invasive tool for biodiversity
14 assessment, enabling the continuous recording of animal vocalizations in diverse ecosystems [1; 2; 3].
15 In biologically rich tropical countries such as Colombia, PAM provides valuable insights into species
16 presence, distribution, and ecological dynamics. The El Silencio Natural Reserve, one of the study
17 sites for BirdCLEF 2025 [4], protects a diverse taxa, making it a valuable location for acoustic
18 biodiversity monitoring. In this context, the BirdCLEF 2025 challenge aims to develop automated
19 methods to identify various taxonomic groups, including birds, amphibians, mammals, and insects,
20 from soundscape recordings collected in the reserve. However, analyzing the data from these
21 soundscapes presents significant challenges due to overlapping vocalizations, background noise from
22 anthropogenic sources, and strong class imbalances across species.

23 Despite the growing success of deep learning models such as BirdNET in large-scale bird identification
24 tasks [5], most efforts in bioacoustics continue to rely on complex architectures or heavily supervised
25 frameworks that require extensive annotated datasets. While these models often perform effectively
26 within specific taxonomic groups (e.g., birds), their scalability and adaptability across other taxa (e.g.,
27 amphibians, mammals, or insects) remain limited [6; 7]. Furthermore, pipelines typically assume
28 well-conditioned input data, placing limited attention on the signal preprocessing stage [8; 9].

29 In multi-species classification scenarios, especially those involving diverse vocal repertoires and
30 heterogeneous sound environments, signal quality critically affects the separability of acoustic
31 features. Although deep learning models offer powerful solutions, their performance can degrade in
32 the presence of noise or misaligned inputs, conditions common in tropical soundscapes. As species
33 vocalizations vary in frequency range, duration, and intensity, inadequate preprocessing may obscure
34 biologically relevant information, ultimately limiting classification accuracy regardless of model
35 complexity. We hypothesize that by enhancing signal quality through targeted preprocessing, it is
36 possible to enable lightweight and general-purpose models to perform competitively, offering an
37 efficient alternative for biodiversity-rich but computationally constrained contexts such as Colombia.

38 In this work, we propose a lightweight preprocessing pipeline tailored for the multi-species classifica-
39 tion task of BirdCLEF 2025. By applying preprocessing steps such as silence removal, vocalization
40 isolation, and controlled Gaussian noise addition, we enhance the clarity and informativeness of the

input representation. We extract acoustic embeddings using the BirdNET Analyzer [5] and pass them to simple classifiers, including multilayer perceptrons (MLPs) and generic convolutional neural networks (CNNs). Our results suggest that preprocessing is not merely a preliminary step but a critical enabler of robust and scalable species classification.

2 Materials and Methods

2.1 Materials: The dataset used in this study (taken from BirdCLEF 2025) includes labeled and unlabeled recordings from 206 species across four taxonomic groups: Aves, Amphibia, Mammalia, and Insecta. Recordings were sourced from Xeno-canto, iNaturalist, and the Colombian Sound Archive, all resampled to 32 kHz, and collected at El Silencio Natural Reserve (6°45'N, 74°12'W) [4]. Aves dominate with 146 species (70.9%) and 27,648 recordings, followed by amphibians (34 species, 16.5%, 583 recordings), insects (17 species, 8.3%, 155 recordings), and mammals (9 species, 4.4%, 178 recordings). Gathered in a tropical rainforest under ecological restoration, these recordings illustrate the value of passive acoustic monitoring (PAM) as a scalable tool for assessing biodiversity and conservation outcomes.

2.2 Multispecies classification methodology: This study focuses on the impact of preprocessing strategies on training robust audio classifiers [9], making the preprocessing pipeline a central component of the methodology. The overall structure follows the standard workflow in species identification tasks [10]: segmentation, preprocessing, feature extraction, and classification (Figure 1). All recordings were segmented into uniform five-second clips, which were then processed through steps such as Gaussian noise addition, silence removal, and voice cleaning to enhance signal clarity and consistency. From these refined segments, we extracted BirdNET embeddings, compact descriptors of species vocalizations, and compared models trained on embeddings from both raw and preprocessed audio. Finally, the embeddings were classified using lightweight models, including convolutional neural networks (CNNs) and multilayer perceptrons (MLPs), highlighting the often-underestimated importance of preprocessing in enabling efficient and scalable bioacoustic classification systems.

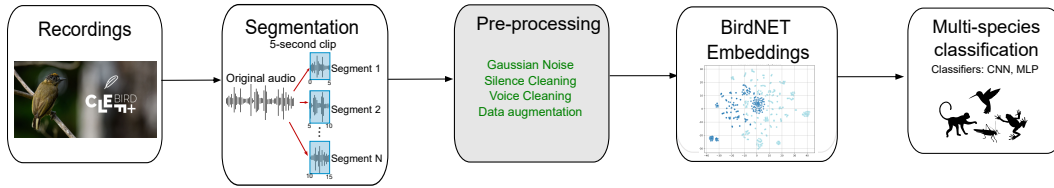


Figure 1: Multi-species classification methodology

2.2.1 Segmentation: Initially, all audio recordings are segmented into fixed-length clips of five seconds. This duration facilitates capturing complete vocalization events while ensuring uniform input dimensions for downstream processing. Segments shorter than five seconds but longer than two seconds are extended using a tiling technique, wherein the original audio content is repeated until the desired length is reached. This standardization is crucial for embedding generation and model training, as it enables batch processing and ensures consistent input size across the dataset.

2.2.2 Preprocessing: After obtaining the fixed-duration segments, a series of preprocessing techniques are applied to increase the quality and consistency of the data. These transformations, designed to improve signal clarity and standardization, aim to facilitate the learning of discriminative features by classifiers, ultimately optimizing performance across evaluation metrics. The complete pipeline is described below:

Band-Limited Gaussian Noise: The first preprocessing step handled inconsistencies in the source audio files, some of which had been up-sampled to 32 kHz or filtered (e.g., bandpass, lowpass), resulting in spectrograms with empty frequency bands that produced ‘black’ regions of zero energy. These irregularities, common in datasets from public repositories like Xeno-canto, can bias training and interfere with silence detection. To mitigate this, we applied a band-limited Gaussian noise augmentation strategy that injects low-level noise into inactive spectral regions, effectively filling the black bands, normalizing information content across samples, and improving the consistency of signal-based silence detection.

86 ***Silence Cleaning:***

87 After correcting the zero-energy bands, we removed minimally informative audio, including silence
88 and low-relevance background sounds. We manually filtered out non-bird classes, while automation
89 was used for bird vocalizations. Each five-second segment was divided into ten 500-ms windows, and
90 we calculated the variance of kurtosis across these windows. Flat signals, such as silence or steady
91 noise, resulted in low variance, whereas bird calls exhibited higher variance. We then excluded the
92 bottom 5% of segments based on kurtosis variance, effectively removing low-content audio while
93 preserving the diversity of valid vocalizations.

94 ***Voice Cleaning:*** Improving the performance of both feature extraction and classification models
95 requires incorporating a voice removal step. Human speech added by field recordists is present in
96 several recordings and introduces noise into the feature extraction process. We used a pre-trained
97 VGG-19 model to identify human voices in five-second segments. Therefore, it classified segments
98 with human voice identifiers. For non-bird species, segments classified as voice were manually
99 verified due to the lower data proportion. Ultimately, segments identified as human voice were
100 removed from the labeled dataset.

101 ***Downsample:*** According to the labeled dataset, we observed imbalance in the number of audio
102 segments for each species within each taxonomic group. To balance the dataset, we first calculated
103 the median of audio segments for all non-bird and bird species separately. Then, we applied random
104 downsampling to values slightly above the median: 500 for the bird group and 15 for the non-bird
105 group. After this, the dataset was divided into train (60%), validation (20%), and test (20%).

106 ***Data Augmentation:*** In this step, species with fewer audio segments than the median were
107 augmented up to that threshold, making the prior filtering of silence and voice cleaning essential
108 to avoid introducing noise that could degrade performance. Augmentation was applied only to
109 underrepresented bird and non-bird species in the training set, using techniques such as white noise
110 addition, time-shifting, and background noise from the ff1010bird dataset [11; 12; 13], thereby
111 increasing sample diversity and supporting better generalization of the classifier.

112

113 **2.2.3 BirdNET Embeddings:**

114 After preprocessing and cleaning, feature extraction was carried out using embeddings generated
115 by the BirdNET Analyzer [5], a widely used CNN-based tool trained on spectrograms of bird
116 vocalizations. This open-source framework offers a strong accuracy–efficiency trade-off [14]. We
117 used v2.4, which outputs a 1024-dimensional embedding per three-second snippet. Since our
118 inputs are five seconds, each segment yields two vectors that we aggregated into a single 1024-D
119 representation via average pooling.

120 ***Embeddings Concatenation:***

121 To add temporal context, for each five-second segment we concatenate its 1024-D vector with those
122 of the next two segments (repeating the last available if needed), forming a 3,072-D input. This
123 sliding window captures onsets, offsets, and inter-call gaps beyond single-segment scope, improving
124 robustness and generalization in complex soundscapes (see Appendix A.2).

125 **2.2.4 Training Classifier:** Following the principle of “Clean Data and Simple Models,” we prioritized
126 lightweight architectures that balance efficiency and performance, utilizing the strength of the
127 extracted embeddings to reduce training data and computational demands. To capture patterns
128 within the embeddings, we fine-tuned Convolutional Neural Networks (ResNet-18 and ResNet-34
129 [15]), reshaping each embedding into a 32×96 matrix for compatibility with 2D convolutions
130 and adapting input/output layers to the task. In parallel, we evaluated a Multi-Layer Perceptron
131 (MLP), with hyperparameters (learning rate, hidden layers, neurons) optimized via Optuna, under
132 the hypothesis that sufficiently informative embeddings allow even simple classifiers to achieve
133 competitive performance.

134 ***Evaluation Metrics:*** We evaluated the performance of our multi-species classification models using
135 standard metrics on a held-out test set comprising 30% of the dataset. F1-score, precision, and
136 accuracy were reported to provide a broad overview of model performance and enable comparison
137 across evaluation dimensions. However, our main evaluation metric was **Recall**, as it best reflects the
138 model’s ability to correctly identify each species and minimizes false negatives, an essential factor in
139 biodiversity monitoring to avoid underestimating richness or overlooking rare taxa.

3 Results

Table 1 summarizes the performance of the models across multiple evaluation metrics. The MLP trained on raw embeddings achieved the highest internal scores in precision, accuracy, weighted F1, and weighted recall. However, when evaluated on the BirdCLEF 2025 Kaggle test set (unseen data), the MLP trained on balanced and preprocessed embeddings obtained the best score (0.756). Although the raw MLP appeared stronger on internal evaluation, the balanced MLP demonstrated superior generalization, underscoring the importance of addressing data imbalance in multi-species classification tasks. This conclusion is further supported by the learning curves presented in Appendix A.3, which show reduced overfitting and improved consistency for the balanced MLP.

Table 1: Comparison of models across evaluation metrics. While the MLP trained on raw data leads in internal metrics, the balanced MLP with cleaned data achieves the highest Kaggle score (0.756), indicating superior generalization.

Model	Precision - W	Accuracy	F1 - W	Recall - W	Kaggle Score
Resnet 18	0,77	0,76	0,76	0,7	0,657
MLP Raw	0,89	0,89	0,88	0,89	0,728
MLP Balanced	0,84	0,83	0,83	0,83	0,756

Figure 2 reports recall per species for a subset of birds, comparing cleaned data (blue) versus raw data (orange). We observe a consistent recall gain for most classes with cleaning. Overall, cleaning raises class-wise performance and supports the global idea that with better signal quality, a simple MLP may match or even surpass more complex models. Recall is particularly relevant in biodiversity monitoring because it measures the model’s ability to correctly identify species and minimize false negatives; failing to detect a species that is actually present could lead to underestimating local richness or overlooking rare and conservation-critical taxa.

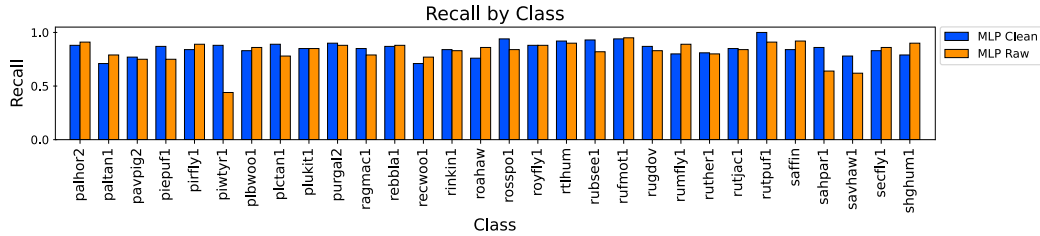


Figure 2: Recall by species for a subset of classes. Blue: cleaned data; orange: raw data. Cleaning improves recall for most species, reinforcing that better signal quality enables strong performance with simple models.

4 Conclusions

This study shows that targeted preprocessing including **silence removal, human voice filtering, and Gaussian noise compensation** is crucial for enabling lightweight models such as MLPs and basic CNNs to achieve competitive performance in multi-species classification. Although raw embeddings showed stronger results on internal metrics, the balanced and preprocessed datasets achieved better generalization on the BirdCLEF 2025 Kaggle test set, a dataset entirely unseen by the models. Demonstrating that data quality can outweigh model complexity, particularly in computationally constrained contexts. Learning curves further confirmed that preprocessing reduces overfitting and supports robustness across diverse acoustic environments, highlighting its value for bioacoustic monitoring (see Appendix A.3 for details). Nonetheless, challenges remain due to class imbalance and the underrepresentation of certain species. Future directions include refining automated noise and silence detection, exploring additional lightweight architectures such as compact transformers, and improving the pipeline to handle longer recordings and diverse ecological soundscapes, with the objective of increasing robustness against acoustic variability and ultimately advancing scalable, accessible, and reliable conservation tools.

References

- [1] Rory Gibb, Ella Browning, Paul Glover-Kapfer, and Kate E. Jones. Emerging opportunities and challenges for passive acoustics in ecological assessment and monitoring. *Methods in Ecology and Evolution*, 10(2):169–185, 2019. doi: 10.1111/2041-210X.13101.
- [2] Bryan C. Pijanowski, Luis J. Villanueva-Rivera, Sarah L. Dumyahn, Almo Farina, Bernie L. Krause, Brian M. Napoletano, Stuart H. Gage, and Nadia Pieretti. Soundscape ecology: The science of sound in the landscape. *BioScience*, 61(3):203–216, 2011. doi: 10.1525/bio.2011.61.3.6.
- [3] V. Ramesh et al. Using passive acoustic monitoring to examine the impacts of ecological restoration on faunal biodiversity in the western ghats. *Biological Conservation*, 282:110071, 2023. doi: 10.1016/j.biocon.2023.110071.
- [4] Holger Klinck, Juan Sebastián Cañas, Maggie Demkin, Sohier Dane, Stefan Kahl, and Tom Denton. Birdclef+ 2025. <https://kaggle.com/competitions/birdclef-2025>, 2025. Kaggle.
- [5] Stefan Kahl, Connor M. Wood, Maximilian Eibl, and Holger Klinck. Birdnet: A deep learning solution for avian diversity monitoring. *Ecological Informatics*, 61:101236, 2021. doi: 10.1016/j.ecoinf.2021.101236.
- [6] Daniel A. Nieto-Mora et al. Systematic review of machine learning methods applied to ecoacoustics and soundscape monitoring. *Heliyon*, 9(4):e20275, 2023. doi: 10.1016/j.heliyon.2023.e20275.
- [7] Hao Wang et al. An efficient model for a vast number of bird species identification based on acoustic features. *Animals*, 12(18):2434, 2022. doi: 10.3390/ani12182434.
- [8] Dan Stowell. Computational bioacoustics with deep learning: a review and roadmap. *PeerJ*, 10: e13152, 2022. doi: 10.7717/peerj.13152.
- [9] Jialin Xie, Julien G. Colonna, and Jing Zhang. Bioacoustic signal denoising: a review. *Artificial Intelligence Review*, 54:3575–3597, 2021. doi: 10.1007/s10462-020-09932-4.
- [10] Bas Ghani, Vincent J. Kalkman, Robert Planqué, and et al. Impact of transfer learning methods and dataset characteristics on generalization in birdsong classification. *Scientific Reports*, 15: 16273, 2025. doi: 10.1038/s41598-025-00996-2.
- [11] Anil Sampath Kumar, Thomas Schlosser, Stefan Kahl, and David Kowerko. Improving learning-based birdsong classification by utilizing combined audio augmentation strategies. *Ecological Informatics*, 82:102699, 2024. doi: 10.1016/j.ecoinf.2024.102699.
- [12] Lucas Ferreira-Paiva, Elizabeth Alfaro Espinoza, Vinicius Martins Almeida, Leonardo Felix, and Rodolpho Neves. A survey of data augmentation for audio classification. In *Congresso Brasileiro de Automática - CBA*, volume 3, 10 2022. doi: 10.20906/CBA2022/3469.
- [13] Dan Stowell. freefield1010: An audio dataset of field recordings. <https://archive.org/details/freefield1010>, 2013. Accessed: 26-May-2025.
- [14] A. Jana, M. Uili, J. Atherton, M. O’Brien, J. Wood, and L. Brickson. An automated pipeline for few-shot bird call classification, a case study with the tooth-billed pigeon, 2025. URL <https://arxiv.org/abs/2504.16276v2>. arXiv preprint.
- [15] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016. doi: 10.1109/CVPR.2016.90. URL <https://doi.org/10.1109/CVPR.2016.90>.

215 A Appendix

216 A.1 Code and Data

217 The code and original datasets used in this study are openly available for reproducibil-
 218 ity and further research at the following links: [https://anonymous.4open.science/](https://anonymous.4open.science/r/NeurIPS-BirdCLEF-25/README.md)
 219 [r/NeurIPS-BirdCLEF-25/README.md](https://anonymous.4open.science/r/NeurIPS-BirdCLEF-25/README.md) and [https://www.kaggle.com/competitions/](https://www.kaggle.com/competitions/birdclef-2025/data)
 220 [birdclef-2025/data](https://www.kaggle.com/competitions/birdclef-2025/data).

221 A.2 Embeddings concatenation

222 For each five-second audio segment, we concatenated its corresponding embedding, with size 1024
 223 according to the aggregation process, with those of the next two consecutive segments from the
 224 same audio file. This results in a single input vector of size 3072 (3×1024), tripling the temporal
 225 window considered by the model. For example, in Fig.1 there's an audio with 4 segments, and the first
 226 concatenated embedding is the union of *Audio0_0*, *Audio0_1*, *Audio0_2*. In cases where fewer than
 227 two additional segments were available (e.g., at the end of an audio file), the last available segment
 228 was repeated to maintain a consistent input size

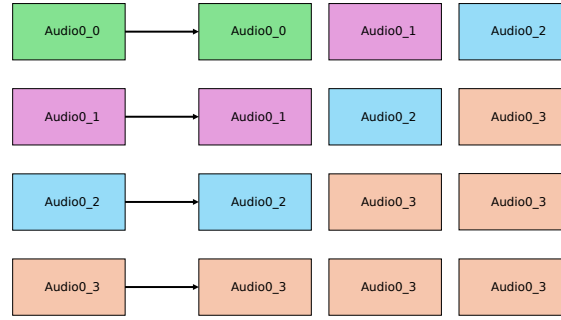
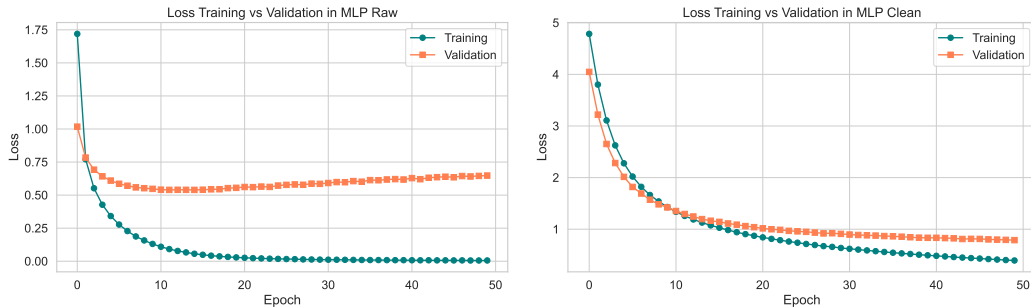


Fig. A 1: Embeddings concatenation description

229 A.3 Loss Train and Validation Curves in MLP

230 The following figure shows the comparison between the loss curves for both models: using raw and
 231 clean embeddings. It shows the training and validation loss curves, and models were trained up to 50
 232 epochs. It can be noticed an overfitting effect in the case of the MLP raw model, and in the case of
 233 the MLP clean model, the difference between the train and validation curves in each epoch is lower
 234 than that of the MLP raw, showing better generalization.



(a) Loss curve for MLP model with raw embeddings (b) Loss curve for MLP with clean embeddings. Validation and training curves indicate adequate generalization.

Fig. A 2: Loss curve for MLP models