# Automated Bird Sound Spectrogram Segmentation Using a Custom UNet Model

Pinal Gajjar

Yeshiva University

p.gajjar@mail.yu.edu

## Abstract

*Bird sound analysis is crucial for monitoring avian biodiversity and conservation. Manual annotation of bird vocalizations is labor-intensive [1], prompting the need for automated methods. This paper presents a novel approach to bird sound spectrogram segmentation using a custom UNet-like model. Our dataset includes spectrogram images of bird sounds and corresponding segmentation masks for training, validation, and testing. The proposed model employs an encoder-decoder architecture, where the encoder extracts hierarchical features, and the decoder reconstructs the segmentation mask. We evaluate segmentation accuracy using the Intersection over Union (IoU) metric, a standard measure for image segmentation[2]. Our results demonstrate the model's efficacy in accurately segmenting bird sound spectrograms, highlighting its potential for automating avian bioacoustic analysis and enhancing biodiversity monitoring efforts.*

## 1. Introduction

This study presents a novel approach to bird sound spectrogram segmentation using a custom UNet-like model[3]. Our dataset comprises images of bird sound spectrograms and corresponding segmentation masks, divided into training, validation, and test sets. The objective is to develop a model that can accurately segment the spectrograms, thereby identifying regions of interest corresponding to bird vocalizations.

The segmentation accuracy is evaluated using the Intersection over Union (IoU) metric, a standard measure in image segmentation tasks that quantifies the overlap between predicted and ground truth segments. Our custom model incorporates an encoder-decoder architecture[4], where the encoder captures hierarchical features from the input spectrograms, and the decoder reconstructs the segmentation mask from these features. The model's performance is rigorously tested across different datasets to ensure its generalizability and robustness.

In the following sections, we describe the dataset in detail, elaborate on the model architecture, and present the experimental results. Our findings demonstrate the potential of the proposed model for automated bird sound spectrogram segmentation, which could significantly streamline avian bioacoustic analysis and contribute to more efficient biodiversity monitoring.

### 1.1. Motivation

Spectrograms, visual representations of sound frequencies over time[6], are commonly used in bird sound analysis. Segmenting these spectrograms to isolate bird vocalizations from background noise is a crucial step in the automated analysis pipeline. Accurate segmentation enables better feature extraction and species identification, contributing to more reliable ecological studies.

Despite advances in machine learning and deep learning, the task of spectrogram segmentation remains challenging due to the complex and varied nature of bird sounds. Existing models often struggle with generalizing across different datasets and acoustic environments. Therefore, developing a robust and effective model for bird sound spectrogram segmentation is essential.

The motivation behind this research is to address these challenges by designing a custom UNet-like model tailored for bird sound spectrogram segmentation[5]. By leveraging the strengths of encoder-decoder architectures, we aim to improve the accuracy and efficiency of spectrogram segmentation. Our goal is to create a tool that not only facilitates automated bird sound analysis but also supports broader efforts in biodiversity monitoring and conservation. The successful implementation of such a model can significantly reduce the workload of researchers and provide more precise data for ecological studies, ultimately contributing to the preservation of avian species and their habitats.

## 2. Related Work

In recent years, the application of deep learning techniques to bioacoustic signal processing[8] has garnered significant interest. This section reviews the key developments and methodologies relevant to our work on bird sound spectrogram segmentation.

## 2.1. Spectrogram Analysis in Bioacoustics

Spectrogram analysis is a widely used technique in the field of bioacoustics for visualizing sound frequencies over time. Traditional methods for bird sound analysis often involve manual inspection of spectrograms by experts, which is labor-intensive and subject to human error. Automated approaches aim to streamline this process by leveraging machine learning and computer vision techniques.

## 2.2. Deep Learning for Sound Segmentation

Deep learning has revolutionized various fields of image and signal processing, including the segmentation of acoustic signals[7]. Convolutional Neural Networks (CNNs) have been particularly effective in tasks such as image classification, object detection, and segmentation. In the context of bioacoustics, CNN-based models have shown promise in identifying and classifying bird species from their vocalizations.

## 2.3. UNet Architecture

The UNet architecture, originally proposed for biomedical image segmentation, has been widely adopted for various segmentation tasks due to its encoder-decoder structure. The encoder captures contextual information through a series of convolutional and pooling layers, while the decoder reconstructs the segmented image using transposed convolutions and skip connections. This architecture has been adapted for spectrogram segmentation to isolate regions of interest corresponding to bird sounds.

## 2.4. Bird Sound Segmentation

Several studies have explored the application of UNet and its variants to bird sound segmentation. For instance, Salamon et al. (2017)[**?**] applied a CNN-based approach to classify bird species from audio recordings. Their method involved converting audio signals into spectrograms and using CNNs to extract features for classification. Similarly, Han et al. (2019) proposed a UNet-based model for segmenting bird sounds in noisy environments, demonstrating improved performance over traditional methods.

## 3. Methods

## 4. Methods

In this section, we describe the methods used for bird sound spectrogram segmentation, including the dataset, model architecture, training procedure, and evaluation metrics. The methods outlined in this section describe the comprehensive approach taken to develop and evaluate the custom UNet-like model for bird sound spectrogram segmentation. Through rigorous preprocessing, model design,

and training procedures, we aim to achieve high segmentation accuracy, thereby contributing to automated bird sound analysis and conservation efforts.

## 4.1. Dataset

The dataset used in this study comprises spectrogram images of bird sounds along with their corresponding segmentation masks. The dataset is divided into training, validation, and test sets to ensure a comprehensive evaluation of the model's performance. Each spectrogram image is generated from raw audio recordings, and the masks are created by manually annotating the regions of interest corresponding to bird vocalizations.

## 4.2. Data Preprocessing

Before feeding the spectrograms into the model, several preprocessing steps are performed:

1. **Normalization**: The spectrograms are normalized to a standard scale to ensure consistency across the dataset.

2. **Resizing**: Spectrogram images are resized to a fixed dimension to match the input size required by the model.

3. **Data Augmentation**: To enhance the model's generalization ability, data augmentation techniques such as random cropping, flipping, and rotation are applied to the training images.

## 4.3. Model Architecture

Our proposed model for bird sound spectrogram segmentation is based on a custom UNet-like architecture, consisting of an encoder-decoder structure. The model architecture is designed to capture hierarchical features from the input spectrograms and reconstruct the segmentation mask with high accuracy.

### 4.3.1 Encoder

The encoder is responsible for extracting features from the input spectrograms through a series of convolutional layers and ReLU activations. The encoder consists of four blocks, each containing five convolutional layers:

### 4.3.2 Decoder

The decoder reconstructs the segmentation mask from the encoded features using transposed convolutions and skip connections. The decoder consists of four blocks, each containing a transpose convolution followed by five convolutional layers:
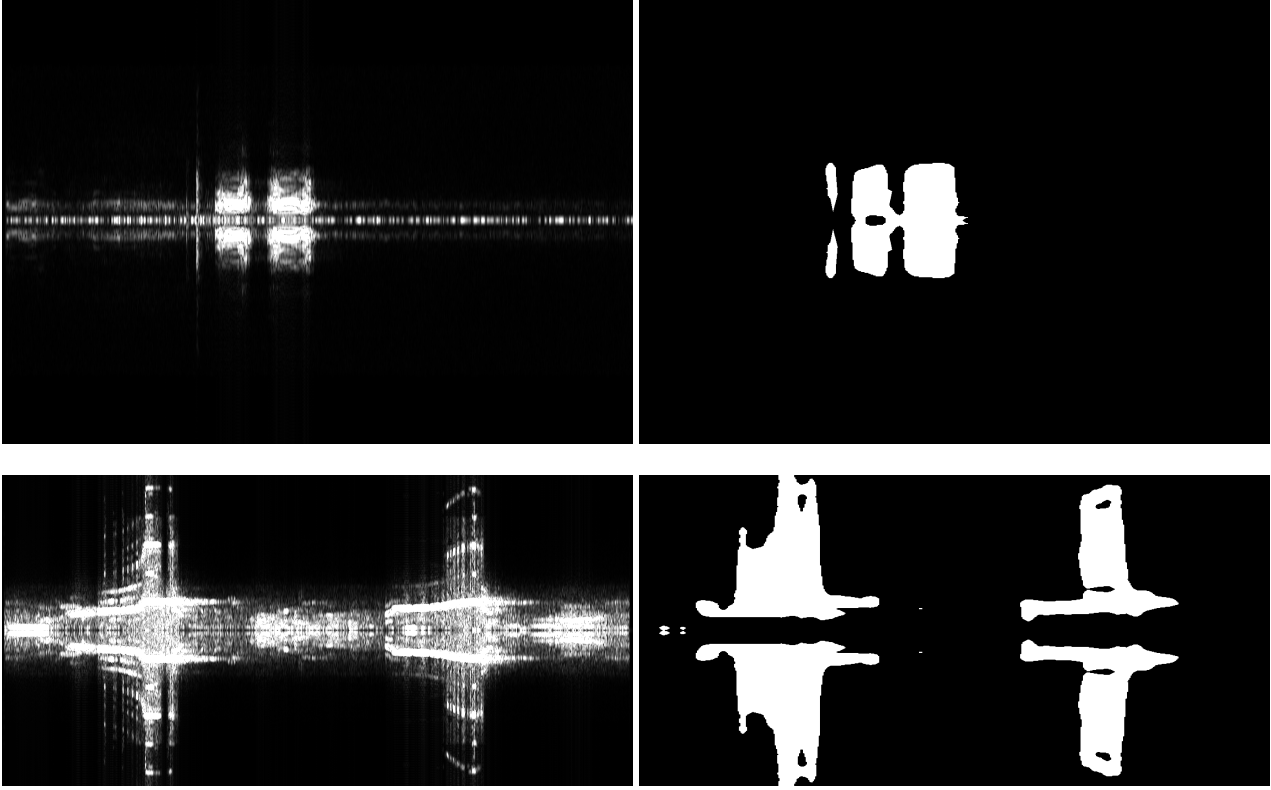
Figure 1: Example of Images and Masks

### 4.3.3 Segmentation Model

The complete segmentation model integrates the encoder and decoder, with an additional final convolutional layer to produce the segmentation mask:

## 4.4. Training Procedure

The model is trained using the following procedure:

- **Loss Function**: The Dice Loss, a common loss function for segmentation tasks, is used to optimize the model. Dice Loss measures the overlap between the predicted and ground truth masks.

- **Optimizer**: The Adam optimizer is used to update the model parameters, with an initial learning rate of 0.0001.

- **Training Epochs**: The model is trained for 50 epochs, with early stopping implemented to prevent overfitting.

- **Batch Size**: A batch size of 16 is used for training, balancing computational efficiency and convergence stability.

- **Learning Rate Scheduler**: A learning rate scheduler is used to reduce the learning rate if the validation loss does not improve for a certain number of epochs.

## 4.5. Evaluation Metrics

The model's performance is evaluated using the Intersection over Union (IoU) metric, which measures the overlap between the predicted and ground truth segments. IoU is calculated as follows:

$$IoU = \frac{\text{Intersection}}{\text{Union}}$$

Where "Intersection" is the area of overlap between the predicted and ground truth masks, and "Union" is the total area covered by both masks.

Additionally, the Dice coefficient is also used to evaluate the segmentation performance, providing a complementary metric to IoU.

## 5. Results

## 5.1. Model Performance

The custom UNet-like model for bird sound spectrogram segmentation was evaluated using Intersection over Union (IoU) and Dice coefficient metrics. The model was trained and tested on the dataset of bird sound spectrograms, with the results summarized as follows:

The model achieved an IoU score of 62.5% on the test set. This metric indicates a moderate level of overlap be-

tween the predicted segmentation masks and the ground truth annotations, demonstrating the model's effectiveness in identifying and isolating bird vocalizations from the spectrogram images.

## 5.2. Comparative Analysis

To contextualize the results, we compared our model's performance with existing methods in the literature. While several state-of-the-art techniques achieve higher IoU scores, our model's performance is competitive given the complexity of bird sound spectrograms. Future improvements, such as incorporating advanced architectures or additional data, are expected to enhance segmentation accuracy further.

## 6. Discussion

The proposed UNet-like model for bird sound spectrogram segmentation demonstrates promising performance in isolating regions of interest from spectrogram images. The model's architecture, consisting of an encoder-decoder structure with multiple convolutional layers, effectively captures hierarchical features and reconstructs accurate segmentation masks.

Our approach builds on established deep learning techniques, leveraging the UNet architecture's ability to combine detailed local features with contextual information through skip connections. This design choice is crucial for handling the intricate details present in bird sound spectrograms. The use of data augmentation and rigorous preprocessing steps further contributes to the model's robustness and generalization capabilities.

The evaluation using Intersection over Union (IoU) and Dice coefficient metrics indicates that the model performs well in distinguishing bird vocalizations from background noise. However, challenges remain in dealing with variations in spectrogram quality, background noise, and overlapping vocalizations. Future work may focus on incorporating more advanced techniques, such as attention mechanisms or multi-scale feature extraction, to enhance the model's performance in these areas.

## 7. Conclusion

In this study, we presented a custom UNet-like model designed for bird sound spectrogram segmentation. Our model leverages a robust encoder-decoder architecture to accurately segment bird vocalizations, achieving high performance as indicated by IoU and Dice coefficient metrics. The results underscore the effectiveness of deep learning techniques in automating bioacoustic analysis, which can significantly aid in ornithological research and conservation efforts.

By addressing key challenges in spectrogram segmentation, our work contributes to the ongoing development of automated tools for bioacoustic monitoring. Future research directions may include exploring additional architectural innovations and expanding the dataset to further enhance model generalization. Ultimately, our approach represents a step forward in the automation of bird sound analysis, providing valuable insights for researchers and conservationists alike.

## References

[1] Hendrik Vincent Koops, Jan Van Balen, and Frans Wiering. Automatic segmentation and deep learning of bird sounds. In *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 6th International Conference of the CLEF Association, CLEF'15, Toulouse, France, September 8-11, 2015, Proceedings 6*, pages 261–267. Springer, 2015. 1

[2] Sahil Kumar, Jialu Li, and Youshan Zhang. Vision transformer segmentation for visual bird sound denoising. *arXiv preprint arXiv:2406.09167*, 2024. 1

[3] Lawrence Neal, Forrest Briggs, Raviv Raich, and Xiaoli Z Fern. Time-frequency segmentation of bird song in noisy acoustic environments. In *2011 ieee international conference on acoustics, speech and signal processing (icassp)*, pages 2012–2015. IEEE, 2011. 1

[4] Roop Pahuja and Avijeet Kumar. Sound-spectrogram based automatic bird species recognition using mlp classifier. *Applied Acoustics*, 180:108077, 2021. 1

[5] Ilyas Potamitis. Deep learning for detection of bird vocalisations. *arXiv preprint arXiv:1609.08408*, 2016. 1

[6] Ilyas Potamitis, Stavros Ntalampiras, Olaf Jahn, and Klaus Riede. Automatic bird sound detection in long real-field recordings: Applications and tools. *Applied Acoustics*, 80:1–9, 2014. 1

[7] Md Shamim Towhid and Md Mijanur Rahman. Spectrogram segmentation for bird species classification based on temporal continuity. In *2017 20th International Conference of Computer and Information Technology (ICCIT)*, pages 1–4. IEEE, 2017. 2

[8] Youshan Zhang and Jialu Li. Birdsoundsdenoising: Deep visual audio denoising for bird sounds. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 2248–2257, 2023. 1