



Improving Bird Classification with Primary Color Additives

Ezhini Rasendiran R¹, Chandresh Kumar Maurya²

¹Department of Metallurgical Engineering and Materials Science, Indian Institute of Technology Indore, India

²Department of Computer Science & Engineering, Indian Institute of Technology Indore, India

mems210005019@alum.iiti.ac.in, chandresh@iiti.ac.in

Abstract

We address the problem of classifying bird species using their song recordings, a challenging task due to environmental noise, overlapping vocalizations, and missing labels. Existing models struggle with low-SNR or multi-species recordings. We hypothesize that birds can be classified by visualizing their pitch pattern, speed, and repetition—collectively called **motifs**. Deep learning models applied to spectrogram images help, but similar motifs across species cause confusion. To mitigate this, we embed frequency information into spectrograms using primary color additives. This enhances species distinction, improving classification accuracy. Our experiments show that the proposed approach achieves statistically significant gains over models without colorization and surpasses the BirdCLEF 2024 winner, improving F1 by **7.3%**, ROC-AUC by **6.2%**, and CMAP by **6.6%**. These results show the effectiveness of incorporating frequency information via colorization.

Index Terms: Audio Classification, Bird Classification, BirdCLEF-2024, EfficientNet

1. Introduction

Audio classification is the process of categorizing audio recordings into predefined classes based on their acoustic characteristics [1]. This technique is increasingly being used in biodiversity conservation efforts, particularly for the monitoring of wildlife. By using advanced machine learning models and audio processing techniques, researchers can automatically identify bird calls, insect sounds, and other wildlife vocalizations [2, 3]. Such methods are crucial for tracking the presence, abundance, and behavior of species in their natural habitats. Conventional biodiversity monitoring methods often involve manual observation and data collection, which are time-consuming, labor-intensive, and prone to human error [4]. The adoption of audio classification enables more efficient and accurate monitoring, providing real-time insights into ecosystem health. For instance, automated systems can analyze large datasets from bioacoustic Passive Acoustic Monitoring (PAM) machines deployed in remote areas, making it feasible to monitor biodiversity over extended periods. By automating the identification of species through their vocalizations, aids researchers in making informed decisions to protect ecosystems and combat threats like habitat loss and climate change. This technology application is a step toward achieving global sustainability goals by ensuring the preservation of our planet’s rich biodiversity [5].

Deep learning breakthroughs have revolutionized bioacoustic audio classification by enabling the extraction of deep abstract features inherent in raw data [6]; however, overlapping species sounds, a wide array of diverse vocal motifs, and striking similarities between species’ calls continue to complicate

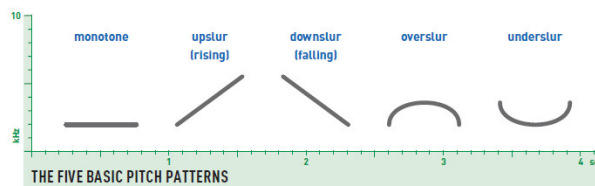


Figure 1: The five basic pitch pattern. Fig. taken from [9] with permission.

robust classification.

Unsupervised source separation, such as the MixIT approach [7], has been employed to disentangle overlapping bird vocalizations in complex soundscapes. This method improves classification performance by isolating individual calls and enabling cleaner extraction of acoustic features. It effectively reduces interference from various sound motifs of different classes, leading to a more precise species identification [8]. However, over-separation can, on some occasions, remove essential contextual cues, diminishing the detection probability of the most prominent species. In these instances, over-separation may also fragment longer vocalizations into isolated notes that resemble calls from other species and result in misclassification.

We therefore introduce a novel feature engineering method that embeds frequency information directly into the input mel spectrogram. This enhancement enables the model to capture frequency variations within similar vocal motifs across different species. By emphasizing these underlying differences, the model is better equipped to learn and distinguish between species calls. Our experiments revealed that this strategy significantly improves the robustness of the classification. Statistical analysis confirmed that our engineered features led to meaningful performance gains compared to model without this enhancement. Our main contributions are:

- Propose a novel idea of embedding frequency information into the mel spectrogram for solving the similar motif problem in the grey-scale mel spectrogram.
- We demonstrate through empirical study that our proposed approach is effective in handling similar motif pattern and outperforms the BirdCLEF 2024 winner model.

2. Background and Problem Statement

2.1. Bird Sound Visualization

Birds can be identified by visualizing their sounds [9]. Key aspects include *pitch pattern*, *speed*, *repetition*, *pauses*, and *tone quality*. Spectrogram symbols offer a simplified artistic repre-

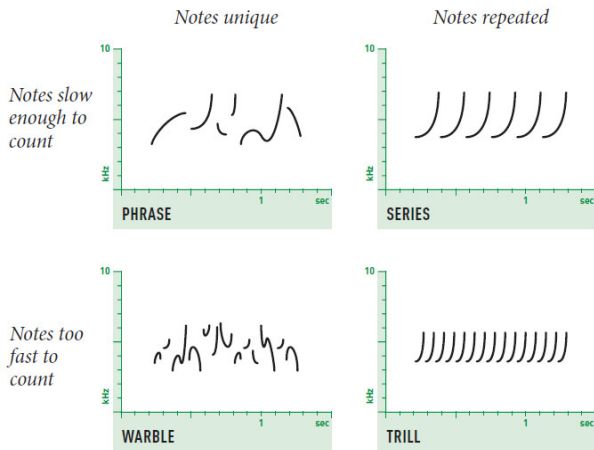


Figure 2: Repetition and speed categorization: Phrases, series, warbles, and trill. Fig. taken from [9] with permission

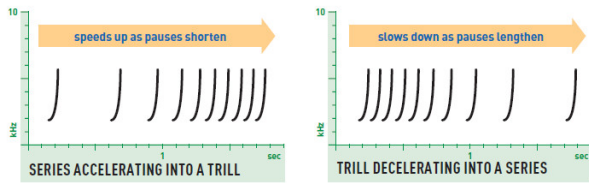


Figure 3: Series acceleration into a trill and trill decelerating into a series. Fig. taken from [9] with permission

sensation of bird sounds, emphasizing patterns over detail, making them useful for humans but not for computers. For experiments, real spectrograms are used to classify bird species. Unlike music, where exact pitch matters, bird sound identification focuses on pitch variations over time. This approach enhances species recognition by capturing distinctive auditory patterns while leveraging deep learning for accurate classification. All bird sounds can be classified into five subcategories (or three combinations thereof) as shown in Fig. 1.

Monotone sounds remain at a constant pitch and appear as horizontal lines on the spectrogram. **Upslurred** sounds increase in pitch, showing an upward tilt. **Downslurred** sounds decrease in pitch, showing a downward tilt. **Overslurred** sounds rise and then fall in pitch, with the highest point occurring in the middle. **Underslurred** sounds fall and then rise in pitch, with the lowest point occurring in the middle.

Repetition (aka motif) and speed relate to the fact that the birds sing the same note multiple times and the pace at which it happens, respectively. Together, there are four basic patterns of repetition and speed: **phrases**, **series**, **warbles** and **trills**. Phrases and series are slower sounds, where individual notes are distinct enough to count. Phrases contain unique notes that are not repeated, while series consist of a single note repeated multiple times. Warbles and trills are faster versions of phrases and series, with notes occurring too rapidly to count (typically faster than about eight notes per second). These motifs are intricately combined to form bird songs, which vary in its pitch and pace. Temporal variations include acceleration and deceleration (shown in Fig. 3). Such dynamic patterns highlight the complexity of avian acoustics, making bird calls invaluable for species identification and behavioral studies. Bird species audio

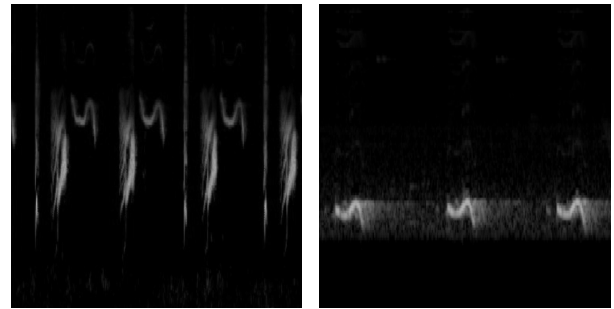


Figure 4: Grayscale mel spectrogram two bird species sharing some motif patterns

recognition ultimately comes down to identifying these unique motifs in a recording.

2.2. Problem Statement

Given a dataset $\{X_i, Y_i\}_{i=1}^n$ with $X_i = [x_1, x_2, \dots] \in \mathcal{X}$ containing multiple instances x_j (recordings chopped in fix-size windows as discussed in §3), and weak labels $Y_i \in \{0, 1\} \in \mathcal{Y}$ at the recording labels. Here, $Y_i = 1$ if any of the instance x_i is positive and $Y_i = 0$ if all the instances are negative. Our goal is to build a multi-class multi-label classifier $f : \mathcal{X} \rightarrow \mathcal{Y}$.

3. Methodology

This section describes our approach to solving the audio classification of bird sounds. We hypothesize that using the mel spectrogram alone as an audio feature and feeding it to the deep learning model is insufficient. The reason is that the grayscale mel spectrogram (being a single channel) loses the frequency information (when used as an image fed to a convolution neural network (CNN)) which is a critical part of characteristic sound identification. The second issue in using the grayscale mel spectrogram to classify different bird sounds is that different birds appear to have similar motifs (series, trills, warbles, etc.) at different frequencies. For example, in Fig. 4, two different birds (Blyth's Reed Warbler and Asian Koel) share a motif pattern. The third issue is that external sound sources like horns and birds may appear similar (in terms of motif) at the same or different frequencies. If we somehow can embed the frequency information into the mel spectrogram, we hope to solve this problem. Next, we discuss such an approach and split our discussion into the following: (1) acoustic event detection, (2) feature engineering mel spectrogram, and (3) frequency information embedding by primary color additives. Finally, we discuss the model architecture.

3.1. Acoustic event detection

We utilize BirdCLEF 2024¹ dataset having 182 species for classification. In this, all audio recordings are *weakly labeled*, that is, a label is present at the recording level instead of the duration level. Its duration can span anywhere between 3 seconds to 30 minutes. First audio is denoised to segregate the instances of bird sound activity, and a high pass filter is applied with a cutoff frequency of 300Hz. We find that acoustic events under

¹<https://www.kaggle.com/competitions/birdclef-2024>

this threshold do not constitute any significant activity. This results in a recording that is highly populated by target acoustic events. Denoising and high pass filter reduces the energy levels of irrelevant sounds. Next, the energy is calculated for each frame as the sum of the squared absolute values of the samples in that frame. The timings of descending energy peak value points above mean energy are calculated (using `find_peaks` package in `scipy.signal` library). Then a 5-second window is wrapped around the energy peaks making complete enclosure over the bird motif constituting an acoustic event. Then a multiple event level example grouping takes place with the condition that the following lower peak energy instances should not share more than 50% time with the preceding acoustic events (to prevent overlapping acoustic events). Overall, we choose 5 events with 5 seconds duration each (sometimes it could be less number of events depending on the presence of the significant acoustic event). The probability of finding the labeled audio event among 5 sound events above mean energy is assumed as 1 in a focal recording which means intentionally capturing sounds from a specific area of interest compared to a PAM where the device records the entire sound without a focus. Thus, our process is guaranteed to have primary labels sound among the 5 mined sound events.

3.2. Feature engineering

Next, we create a mel spectrogram of detected acoustic events. As mentioned at the beginning of this section, mel spectrograms are single-channel and represent the pitch shift or time-stretching phenomena for different acoustic events in the recording. Computer vision models like Resnet [10], Deformable Convolutional Networks (DCNs) [11], etc. are translational invariance (in a sense) and hence will predict the same label for different acoustic events if motifs are shared. Alternatively, the model assigns a high probability to a secondary label as well in the case of shared motifs making the classification less accurate. Thus, the output can not be relied upon completely and we need to encode sound frequencies in some way during the learning mechanism. To this end, the mel spectrogram is normalized in the range (0,1) followed by scaling on log-scale and then again normalized within the range (0,1).

3.3. Frequency information embedding by primary color additives

As mentioned, mel spectrogram is missing the frequency information when fed to the deep learning model. To solve this issue, we proceed as follows. The pixel representation of any mel frequency bin in an image is in an RGB channel color ratio of (1:1:1). To discriminate pixels of mel frequency bins based on their frequency information variation, the mel spectrogram is divided into 3 equal regions between lowest frequency (f_{min}) and highest frequency (f_{max}) initialized during mel spectrogram creation. Therefore the number of mel frequency bins in a region will be $n_{bins} = total_bins/3$ where $total_bins$ is a parameter set during mel spectrogram creation. The first region of the mel frequency bins whose frequency range starts from f_{min} Hz its primary color pixels will be mapped to `color_array` $RG(1-t, t)$ where t is defined as

$$t = \frac{\text{Index of mel frequency bin (counted from } f_{min})}{n_{bins}} \quad (1)$$

As the index of the bins increases, the red color channel linearly decreases from pure red, while the green color channel linearly increases, both colors transitioning in equal amounts

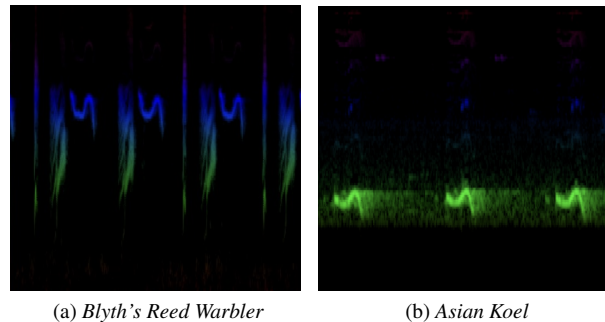


Figure 5: Mel spectrogram of two bird species sharing some motif patterns. Motifs now are distinguishable by a deep learning model because of colorization.

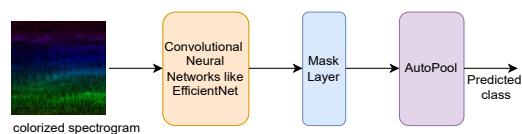


Figure 6: Model architecture

and simultaneously in every bin upwards. This primary color channel addition gives a distinct secondary color to all the mel frequency bins in a region. Finally, pixel values of the mel frequency bins are multiplied with the color array `color_array*pixel.value`. The same operation is performed on the mel frequency bins of other regions with the color arrays `color_array GB(1-t, t)` and `BR(1-t, t)`, respectively. As a result, we get a colorized mel spectrogram which looks like Fig. 5. Note that colorization may be seen as an *approximation* to the frequency information encoding in the mel spectrogram. As such, a colorized mel spectrogram may help distinguish two different bird species sharing a motif pattern which we show in the empirical section.

3.4. Model Architecture

For audio classification, we use the EfficientNetB0 architecture which is a type of CNN [12] for learning followed by the AutoPool layer [13] as shown in Fig. 6. AutoPool is a pooling mechanism designed for weakly labeled audio classification task scenarios involving multiple-instance learning (MIL). Unlike traditional pooling methods such as max or average pooling, AutoPool introduces a trainable pooling function that learns how to aggregate instance-level predictions into recording-level predictions during training. This flexibility allows the model to adaptively balance between max pooling (highlighting dominant signals) and average pooling (capturing broader patterns) based on the dataset characteristics. AutoPool is well-suited for our task because learning happens from weakly labeled data. Whereas predictions happen at a recording level during testing. AutoPool layer pools the logits of 5 instances per recording and then passes it to the sigmoid activation function for multi-class multi-label prediction. For recordings with less than five instances, zero image channels are used to fill the remaining instances, and before AutoPool, binary mask is applied on the zero image channels. Auto Pooled probability aggregation for

multiple instances learning is given by (2).

$$\hat{P}^\alpha(Y|X) = \sum_{x \in X} \hat{p}(Y|x) \left(\frac{\exp(\alpha \cdot \hat{p}(Y|x))}{\sum_{z \in X} \exp(\alpha \cdot \hat{p}(Y|z))} \right) \quad (2)$$

where $\hat{P}(Y|x)$ is the probability of class Y given instance x and α is a scalar parameter learned along with the model parameters during training, which controls the pooling behavior. In particular, $\alpha = 0$ equals the unweighted mean, $\alpha = 1$ equals soft-max pooling, and $\alpha \rightarrow \infty$ is a max operator. Note that the auto pooling is done over each class separately for handling *multi-label* problem. However, we do not vary α for each class which we leave for future study.

3.4.1. Loss function

The model in Fig. 6 is optimized for binary cross-entropy loss function (3).

$$\mathcal{L}_{BCE} = -\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^C [y_{ij} \log \hat{y}_{ij} + (1 - y_{ij}) \log (1 - \hat{y}_{ij})] \quad (3)$$

where N is the sample size and C is the number of classes.

4. Experiments

4.1. Dataset

We utilize BirdCLEF 2024 data [14] for training and validation of the proposed model. It has 24459 audio recording of 182 bird species of which we took only 23920 (possibly some duplicates). Due to the unavailability of the *hidden* test data of BirdCLEF competition, the test set is prepared from the given audio files. We filter files with both primary and secondary labels as the secondary labels are noisy and not reliable². We get 1873 files which are discarded. To be consistent with the test data used in the existing literature, we use recordings with primary labels only. The remaining files are split in the ratio of 80:20 for training and validation.

4.2. Baselines

We use the winner model³ from BirdCLEF 2024 as a baseline which is essentially an EfficientNET-B0 with some augmentations applied to the input such as horizontal cutmix, leveraging pseudo-labeled data, etc. Additionally, we removed the colorization from the mel spectrogram to perform an ablation study to see the effect of the color additives.

4.3. Training and Evaluation details

As discussed in §3.4, the model is trained following K-fold cross-validation (K=5 in our case) strategy. During validation, we feed top-5 acoustic event windows selected from the validation set following the procedure as mentioned in §3.1. Note that this evaluation is more realistic for PAM recordings (used in BirdCLEF hidden set evaluation) compared to just testing the *first* 5-second window (which may or may not contain the actual species in a PAM environment) for the presence of species (as done by the winners of the BirdCLEF). The model is optimized

²<https://www.kaggle.com/competitions/birdclef-2024/discussion/540969>

³<https://www.kaggle.com/competitions/birdclef-2024/discussion/512197>

Table 1: *Performance comparison of the proposed approach with the winner model from BirdCLEF 2024 and without colorization on the 5-fold validation set. * entry indicates a statistically significant difference from the model without colorization on the Wilcoxon signed-rank test (one-tailed) with $\alpha = 0.05$.*

Model	Macro-F1	Macro ROC-AUC	CMAP
Winner model BirdCLEF	0.6371	0.9220	0.6915
Proposed Model w/o Colorization	0.6676	0.9765	0.7217
Proposed Model w/ Colorization	0.6833*	0.9797*	0.7374*

with the initial learning rate of $3e^{-3}$ decaying to $1e^{-6}$ following cosine annealing LR scheduler. We set the batch size to 90 and epochs to 30. AdamW optimizer is used for the minimization of the loss function. The winner model is directly borrowed from the code provided by the winners⁴. We use the same data splits to train and test the winner model for a fair comparison.

4.4. Evaluation Metrics

Macro ROC-AUC, macro F1, and Class-averaged Mean Average Precision (CMAP) are used for the performance evaluation of the model. These metrics are best suited for multi-class, multi-label problems with imbalanced classes. In brief, CMAP is the mean of the per-class precision scores used in the previous BirdCLEF competitions [15]. However, BirdCLEF 2024 used macro ROC-AUC.

4.5. Results

As shown in Table 1, we can observe that the proposed approach outperforms the winner model on all metrics by 7.3% on F1, 6.2% on ROC-AUC, 6.6% on CMAP, respectively. Interestingly, we do not use any data augmentation during training/inference whereas the winner model does. To see the effect of the proposed color additives for frequency embedding, we perform an ablation study. As shown in the Table 1 (2nd row), we find that the colorization is effective in identifying bird species with similar motifs.

5. Conclusion and Future Works

The present work studies bird classification problem through the lens of multiple instance learning. To tackle the problem of identifying bird species with similar motifs pattern, we propose the idea of color additives for frequency embedding. The empirical results reveal that the colorization is an effective method for classifying birds with similar motifs at different frequency bands. There are several open research directions to proceed. We hypothesize that the colorization will be more effective for classifying overlapping vocalization combined with the voice activity detection module for audio segmentation [16]. Another interesting study may be to see the effect of the α parameter in the multi-class multi-label problem.

6. Limitations

Our study is limited to identify primary labels. As such, our model can not classify unknown bird species present in the audio recordings.

⁴https://github.com/arpoyda/BirdCLEF_2024

7. References

- [1] S. Hershey, S. Chaudhuri, D. P. Ellis, J. F. Gemmeke, A. Jansen, R. C. Moore, M. Plakal, D. Platt, R. A. Saurous, B. Seybold *et al.*, “Cnn architectures for large-scale audio classification,” in *2017 IEEE international conference on acoustics, speech and signal processing (icassp)*. IEEE, 2017, pp. 131–135.
- [2] S. Kahl, C. M. Wood, M. Eibl, and H. Klinck, “Birdnet: A deep learning solution for avian diversity monitoring,” *Ecological Informatics*, vol. 61, p. 101236, 2021.
- [3] O. Sheikh, J. Doe, and J. Smith, “Bird whisperer: Leveraging large pre-trained acoustic model for bird call classification,” in *Proceedings of Interspeech*. ISCA, 2024. [Online]. Available: https://www.isca-archive.org/interspeech_2024/sheikh24_interspeech.html
- [4] R. D. Gregory, D. W. Gibbons, and P. F. Donald, “Bird census and survey techniques,” *Bird ecology and conservation*, pp. 17–56, 2004.
- [5] S. Chowfin and A. Leslie, “Using birds as bioindicators of forest restoration progress: A preliminary study,” *Trees, Forests and People*, vol. 3, p. 100048, 2021.
- [6] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [7] S. Wisdom, E. Tzinis, H. Erdogan, R. Weiss, K. Wilson, and J. Hershey, “Unsupervised sound separation using mixture invariant training,” *Advances in neural information processing systems*, vol. 33, pp. 3846–3857, 2020.
- [8] T. Denton, S. Wisdom, and J. R. Hershey, “Improving bird classification with unsupervised sound separation,” in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 636–640.
- [9] E. Birding, “Visualizing sound,” 2024, accessed: 2024-12-28. [Online]. Available: <http://earbirding.com/blog/specs/visualizing-sound>
- [10] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770–778.
- [11] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, “Deformable convolutional networks,” in *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2017, pp. 764–773.
- [12] M. Tan and Q. V. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *Proceedings of the 36th International Conference on Machine Learning (ICML)*. PMLR, 2019, pp. 6105–6114.
- [13] B. McFee, J. Salamon, and J. P. Bello, “Adaptive pooling operators for weakly labeled sound event detection,” *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, vol. 26, no. 11, p. 2180–2193, Nov. 2018. [Online]. Available: <https://doi.org/10.1109/TASLP.2018.2858559>
- [14] S. Kahl, T. Denton, H. Klinck, V. Ramesh, V. Joshi, M. Srivathsa, A. Anand, C. Arvind, H. Cp, S. Sawant *et al.*, “Overview of birdclef 2024: Acoustic identification of under-studied bird species in the western ghats,” in *CLEF 2024-25th Conference and Labs of the Evaluation Forum*, 2024, pp. 1948–1957.
- [15] S. Kahl, F.-R. Stöter, H. Goëau, H. Glotin, R. Planque, W.-P. Vellinga, and A. Joly, “Overview of birdclef 2019: large-scale bird recognition in soundscapes,” in *CLEF 2019-conference and labs of the evaluation forum*, vol. 2380, no. 256. CEUR, 2019.
- [16] N. Gu, K. Lee, M. Basha, S. K. Ram, G. You, and R. H. Hahnloser, “Positive transfer of the whisper speech transformer to human and animal voice activity detection,” in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2024, pp. 7505–7509.