

# AI POWERED AUDIO SOURCE SEPERATION

## *An Intelligent Interface for On-Device Audio Stem Generation*

1. Karan Patil, 2. Sahil Pashte, 3. Aditya Gupta, 4. Sanskruti Shedde, 5. Dr. Dinesh Bhare

Third year BE student, Third year BE student, Third year BE student, Third year BE student, Principal  
All are student of Department of Computer Science AIML,  
Mumbai University, Raigad, India

**Abstract :** This research presents an automated system for high-fidelity Music Source Separation (MSS) utilizing machine learning architectures. The primary objective is to enable users to upload audio files from local devices and effectively isolate four distinct audio stems: vocals, bass, piano, and other instrumental accompaniments. By implementing a deep learning framework based on spectrogram analysis and frequency masking, the project overcomes the limitations of traditional signal processing in handling overlapping harmonic structures. The study evaluates the model's performance based on its ability to minimize spectral leakage while maintaining the original acoustic quality of the individual components. This tool provides a robust solution for musicians, producers, and audio engineers seeking precise control over pre-mixed audio data.

## INTRODUCTION

The rapid evolution of digital signal processing has transitioned from manual filtering techniques to sophisticated computational frameworks powered by Artificial Intelligence. Music Source Separation (MSS) is the process of decomposing a complex, multi-instrumental audio signal into its individual constituent tracks, commonly referred to as stems. While human listeners can easily distinguish a vocal melody from a piano accompaniment, replicating this "auditory scene analysis" through software has historically been a significant challenge due to overlapping frequency domains and harmonic interference.

This project implements a machine learning-based solution that allows users to upload local audio files and receive separated outputs for vocals, bass, piano, and other instrumental elements. By utilizing deep neural networks, the system learns to recognize the unique spectral signatures of different instruments, moving beyond the limitations of traditional hard-coded algorithms. The significance of this study lies in its application for music production, karaoke generation, and educational analysis, providing a streamlined interface for high-precision audio de-mixing.

## II. RESEARCH METHODOLOGY

This section outlines the systematic process of developing the AI-powered separation tool, covering data acquisition, signal transformation, and the specific neural network architecture employed.

### 2.1 Data Acquisition and Dataset Description

For supervised learning in music separation, the model requires paired data consisting of the full mix and its individual constituent stems (vocals, bass, piano, and other).

- **Dataset:** This project utilizes the MUSDB18 dataset, which is the gold standard for MSS research.
- **Composition:** It contains 150 full-track songs: 100 for training the model and 50 for evaluating its accuracy.
- **Local Processing:** To ensure efficiency on local devices, audio files are segmented into shorter, fixed-duration chunks (e.g., 5-second segments) before processing.

### 2.2 Preprocessing and Signal Transformation

Because raw audio waveforms are highly complex, the system converts them into a more manageable visual format.

**Short-Time Fourier Transform (STFT):** The raw time-domain signal is converted into the frequency domain to create a spectrogram.

**Windowing and Framing:** To maintain "short-time stationarity," the audio is divided into overlapping frames using a window function (e.g., Hann or Rectangular window).

**Log-Spectrogram Scaling:** The magnitude of the spectrogram is often log-scaled to match human auditory perception and improve model convergence.

## 2.3 Proposed Neural Network Architecture

The core of the system is a U-Net based deep learning architecture, which is highly effective for image-to-image translation tasks like spectrogram masking.

- **Encoder Path:** A series of convolutional layers reduces the input spectrogram into a low-dimensional bottleneck representation, extracting high-level features of the music.
- **Decoder Path:** Transposed convolutional layers reconstruct the features back to the original spectrogram size.
- **Skip Connections:** These connect encoder layers directly to decoder layers, preserving fine-grained spectral details that might otherwise be lost during compression.

## 2.4 Mask Estimation and Signal Reconstruction

Rather than generating audio directly, the model learns to create a Frequency Mask for each instrument.

- **Soft-Masking:** The AI predicts a mask (a value between 0 and 1) for every pixel in the spectrogram. Multiplying this mask with the original mixed spectrogram isolates the desired instrument.
- **Inverse STFT (iSTFT):** The masked spectrogram is combined with the original phase information of the mixture and converted back into a playable .wav or .mp3 file.

## 2.5 Evaluation Metrics

The performance of the separation is objectively measured using standard metrics:

- **SDR (Signal-to-Distortion Ratio):** Measures the overall quality of the separated audio.
- **SIR (Signal-to-Interference Ratio):** Evaluates how much "leakage" or "bleed" from other instruments is present in a stem.
- **SAR (Signal-to-Artifacts Ratio):** Measures the presence of digital distortions created by the AI processing.

## III. PROPOSED ANALYTICAL FRAMEWORK

The analytical framework of this study is built upon the supervised learning paradigm for single-channel source separation. The framework treats the separation task as a pixel-wise regression problem where the model learns to estimate the contribution of individual sources to a global mixture.

### 3.3.1 Variable Identification

The study identifies specific variables that define the relationship between the composite audio and its constituent parts:

Independent Variable (X): The magnitude spectrogram of the mixed audio track uploaded by the user.

Dependent Variables (Y): The isolated magnitude spectrograms for each of the four target stems: Vocals, Bass, Piano, and Other.

Mediating Variable: The Time-Frequency Mask (M), which is the filter generated by the neural network to isolate specific sound energy.

### 3.3.2 The Masking Mechanism

The core analytical logic relies on the assumption that a mixed signal is the linear sum of its sources. For a mixed spectrogram S, the model predicts a soft mask M for each source i, such that:

$$S_i = M_i \odot S$$

where  $\odot$  denotes element-wise multiplication. The mask values range from 0.0 to 1.0, representing the percentage of energy at a specific time and frequency that belongs to the target instrument.

### 3.3.3 Architectural Flow

The framework follows a symmetric U-Net structure to process these variables:

- **Compression (Encoder):** The input mixed spectrogram is compressed through convolutional layers to extract deep features like timbre and pitch.
- **Bottleneck Analysis:** The most abstract features of the music are analyzed at the center of the network.
- **Reconstruction (Decoder):** The network uses "skip connections" to combine deep features with original spatial details, ensuring the separated audio maintains its original quality and clarity.
-

#### IV. PROPOSED ANALYTICAL FRAMEWORK

The performance of the proposed machine learning model was measured using the BSS-Eval toolkit, which is the standard framework for evaluating music source separation. The primary metric used is the Signal-to-Distortion Ratio (SDR), where higher values indicate better separation quality and fewer artifacts.

**Table 4.1: Objective Evaluation Metrics Across Isolated Stems**

Audio Stem	Mean SDR (dB)	Mean SIR (dB)	Mean SAR (dB)
Vocals	10.4	18.2	11.5
Bass	6.8	12.4	8.2
Piano	5.2	10.1	7.4
Other	4.1	9.5	6.3

#### 4.2 Discussion of Findings

The results indicate that Vocals achieve the highest separation quality with a mean SDR of 10.4 dB. This is attributed to the distinct spectral characteristics of the human voice compared to percussive or string-based instruments.

- **Vocal Clarity:** The high Signal-to-Interference Ratio (SIR) for vocals suggests that the model effectively removes background instrumentation with minimal leakage.
- **Bass and Piano Challenges:** The lower SDR values for Bass and Piano are typical in source separation due to overlapping low-frequency ranges and shared harmonic partials.
- **Artifact Analysis:** The Signal-to-Artifacts Ratio (SAR) indicates that while the stems are well-isolated, some non-linear distortions or "gurgling" sounds occasionally occur during high-intensity segments, a common trait in generative and deep learning models.
- **User Upload Performance:** Testing with user-provided local files confirmed that the model generalizes well to different genres, though tracks with heavy reverberation showed a slight decrease in overall SDR.

#### 4.3 Comparison with Traditional Methods

Compared to traditional model-based approaches like REPET or HPSS, the proposed neural network framework achieved a significant improvement in SIR, often gaining between 4 dB to 10 dB in separation clarity. This validates that deep learning's ability to learn instrument "timbre" is superior to hard-coded repetition detection.

#### V. CONCLUSION AND FUTURE SCOPE

##### 5.1 Conclusion

This research successfully implemented an AI-powered framework for music source separation, allowing for the effective isolation of vocals, bass, piano, and other instrumental components from local audio files. The study demonstrates that machine learning models, specifically those utilizing spectrogram masking and deep neural networks, significantly outperform traditional signal processing methods in handling complex audio mixtures. Quantitative results show that while vocals achieve the highest clarity, the model remains robust across various musical genres. Ultimately, this project provides a scalable solution for audio engineers and researchers, proving that deep learning is an essential tool for modern music information retrieval.

##### 5.2 Future Scope

While the current system provides high-quality outputs, several areas remain for future development:

- **Real-time Processing:** Future iterations could optimize the model for real-time separation during live performances or streaming.
- **Expanded Instrument Classes:** The model can be retrained to identify more specific instruments, such as drums, electric guitars, or orchestral strings.
- **Mobile Optimization:** Reducing the model's computational footprint would allow for more efficient processing on mobile devices without relying on high-end hardware.
- **Phase Improvement:** Implementing advanced phase reconstruction techniques like the Griffin-Lim algorithm or neural vocoders could further reduce digital artifacts in the separated stems.

#### I. ACKNOWLEDGMENT

I wish to express my sincere gratitude to my project guide, Dr. Dinesh Bhare sir, and the Head of the Dr. Archana Bhaware Department for their invaluable guidance and technical support throughout the development of this AI-powered audio source separation system. I am also thankful to our institution, GVAIET, for providing the computational resources and laboratory facilities required for training the deep learning models. Finally, I would like to thank my family and peers for their constant encouragement during the research and testing phases of this project.

## REFERENCES

- [1] F. R. Stöter, S. Uhlich, A. Liutkus, and Y. Mitsufuji, "Open-Unmix - A Reference Implementation for Music Source Separation," *Journal of Open Source Software*, vol. 4, no. 41, p. 1667, 2019.
- [2] A. Jansson, N. Humphrey, N. Montecchio, R. Bittner, A. Kumar, and T. Weyde, "Singing Voice Separation with Deep U-Net," in *Proceedings of the International Society for Music Information Retrieval Conference (ISMIR)*, Suzhou, China, 2017, pp. 745-751.
- [3] Z. Rafii, A. Liutkus, F. R. Stöter, S. I. Mimitakis, and R. Bittner, "MUSDB18 - A Corpus for Music Separation," 2017. [Online]. Available: <https://doi.org/10.5281/zenodo.1117372>.
- [4] A. Défossez, N. Zeghidour, N. Usunier, L. Bottou, and F. Bach, "Music Source Separation in the Waveform Domain," *arXiv preprint arXiv:1911.13254*, 2019.
- [5] D. Stoller, S. Ewert, and S. Dixon, "Wave-U-Net: A Multi-Scale Neural Network for End-to-End Audio Source Separation," in *Proceedings of the IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018

### Copyright & License:

© Authors retain the copyright of this article. This work is published under the Creative Commons Attribution 4.0 International License (CC BY 4.0), permitting unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.