# WaveSplit: A Multi-Stage Framework for Audio Enhancement and Audio Denoising – Combining Deep Learning with Psychoacoustic Principles and Adaptive Noise Processing

Parth Gandhi[1*], Kevin Doshi[1], Anand Godbole[1]

[1]*Department of Computer Engineering, Sardar Patel Institute of Technology, Mumbai, India*

***Abstract**: In the field of audio processing, noise interference poses a significant challenge, affecting speech intelligibility and communication quality across multiple domains. Current audio denoising methods often struggle with the delicate balance between noise removal and speech preservation. This paper presents WaveSplit, a novel multi-stage framework for audio enhancement and denoising that addresses these limitations by combining deep learning techniques with psychoacoustic principles and adaptive noise processing. Building upon the CleanUNet architecture, our approach introduces several innovative components: adaptive SNR-based processing, harmonic enhancement that preserves critical speech components, vocal clarity enhancement, and perceptual processing leveraging human hearing characteristics. Evaluations demonstrate that our framework achieves superior performance compared to baseline models, with significant improvements in SNR (76.36 dB compared to 7.20-8.10 dB in baseline models), PESQ scores (1.05 improvement versus 0.77-0.91), and STOI metrics (0.15 versus 0.09-0.13) while reducing the "robotic" artifacts common in traditional methods. This research has significant implications for applications including telecommunications, hearing assistive technologies, content production, and speech recognition systems. By addressing both objective quality metrics and perceptual factors, WaveSplit represents an advancement toward more effective, natural-sounding audio enhancement solutions for real-world environments.*

***Keywords**: Audio denoising, deep learning, psychoacoustic principles, speech enhancement, CNN, adaptive processing, CleanUNet.*

## 1. Introduction

Audio noise presents a pervasive challenge in numerous communication and media environments, significantly impacting speech intelligibility and overall audio quality. From telecommunications and video conferencing to content production and hearing assistive technologies, unwanted noise can severely degrade user experience and communication effectiveness.

Despite significant advancements in audio processing technologies, current denoising methods often produce suboptimal results in real-world scenarios. Traditional digital signal processing approaches like spectral subtraction and Wiener filtering struggle with non-stationary noise, while many neural network-based solutions introduce artifacts that compromise speech naturalness. This fundamental trade-off between noise reduction and speech preservation remains a significant challenge in the field, as overly aggressive denoising frequently results in the "robotic" speech quality that characterizes many commercial solutions.

This paper introduces WaveSplit, a novel multi-stage framework for audio enhancement and denoising that addresses these challenges. By integrating deep learning with psychoacoustic principles and adaptive processing techniques, our solution achieves superior noise reduction while maintaining the natural characteristics of speech. The system builds upon the CleanUNet architecture while incorporating several innovative components: adaptive SNR-based processing, harmonic enhancement targeting speech components, vocal clarity enhancement, and perceptual processing informed by human auditory perception.

The primary objective of this research is to develop a comprehensive, multi-stage framework for audio enhancement and denoising that overcomes the limitations of existing approaches. Specifically, we aim to:

1) Create a denoising solution that achieves significantly better noise reduction while preserving the natural characteristics of speech, avoiding the "robotic" artifacts common in traditional methods.
2) Implement an intelligent system that dynamically adjusts processing parameters based on the detected Signal-to- Noise Ratio (SNR) of audio segments.
3) Incorporate psychoacoustic principles into the audio processing pipeline to enhance perceived quality, focusing on frequencies most important to human hearing and speech intelligibility.
4) Develop techniques that selectively preserve and enhance harmonic components of speech, which are critical for naturalness and intelligibility.
5) Design a flexible, modular system that allows

individual enhancement components to be enabled or disabled according to user needs and specific audio conditions.

6) Through this research, we aim to contribute significantly to advancing audio processing technology, providing a practical solution for real-world applications in telecommunications, content creation, assistive technologies, and various environments where noise interference impacts communication quality.

## 2. Related Work

Significant research has been conducted in the field of audio denoising and enhancement, with approaches evolving from traditional signal processing methods to advanced deep learning techniques. Traditional methods like spectral subtraction [1] and Wiener filtering [2] have been foundational but often struggle with non-stationary noise and can introduce musical noise artifacts. The Minimum Mean Square Error (MMSE) estimator proposed by Ephraim and Malah [3] improved upon these limitations but still faced challenges with complex noise environments.

Deep learning approaches have gained prominence in recent years, offering improved performance in handling diverse noise conditions. Recurrent Neural Networks (RNNs) have been employed by Pascual et al. [4] in their SEGAN architecture, utilizing an end-to-end approach for speech enhancement. Convolutional Neural Networks (CNNs) have also been extensively explored, with architectures like Wave-U-Net [5] and DEMUCS [6] demonstrating effectiveness in separating clean speech from noisy inputs.

The CleanUNet architecture [7], which forms the foundation of our approach, combines the strengths of both CNN and RNN models, utilizing an encoder-decoder structure with skip connections. This architecture has shown promise in maintaining speech quality while reducing noise, though challenges remain in achieving optimal trade-offs between noise reduction and speech preservation.

Recent research has also focused on incorporating perceptual and psychoacoustic principles into deep learning frame- works. Kumar et al. [8] proposed a perceptually-motivated loss function that better aligns with human auditory perception. Zhao et al. [9] integrated attention mechanisms with spectral features to enhance model performance on speech intelligibility metrics.

Despite these advancements, existing approaches often struggle with maintaining natural speech characteristics, particularly in challenging noise environments. The introduction of artifacts and "robotic" speech qualities in heavily processed audio remains a significant limitation. Our work builds upon these foundations while addressing their limitations through a multi-stage approach that combines deep learning with adaptive processing and psychoacoustic principles.

## 3. Methodology

### A. System Architecture Overview

The WaveSplit framework implements a multi-stage approach that combines deep learning techniques with psychoacoustic principles and adaptive noise processing. The system architecture is designed to overcome the limitations of traditional denoising methods by focusing on both noise reduction and speech preservation.
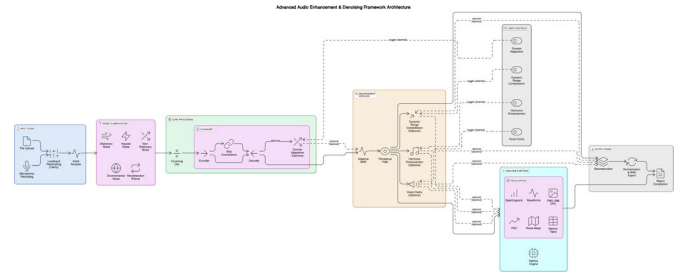


Fig. 1. Complete system architecture - This figure illustrates the end-to-end processing pipeline of the audio enhancement and denoising framework, from input acquisition through noise classification, neural processing, enhancement modules, and output generation

Fig. 1 presents the complete system architecture, showing the full data flow from audio input to enhanced output. The framework consists of three main stages: input processing and noise classification, core neural processing with enhancement modules, and output generation with analysis. Each stage is described in detail below.

### B. Input Processing and Noise Classification

The input stage accepts audio from two primary sources: file uploads and microphone recordings. Regardless of source, all incoming audio undergoes loading and resampling to standardize the sampling rate to 16kHz, ensuring consistent processing throughout the pipeline. Initial analysis is then performed to extract signal characteristics that inform subsequent processing steps.

The noise classification module categorizes the audio according to several noise types: stationary noise (like fans or HVAC systems), impulse noise (clicks and pops), non-stationary noise (such as traffic or background speech), environmental noise, and reverberation effects. This classification plays a crucial role in the adaptive processing approach, allowing the system to optimize its denoising strategy based on the specific noise characteristics present in the audio.

Once classified, the audio is divided into fixed-length segments (3 seconds) by the chunking module. This approach enables efficient batch processing of longer audio files while maintaining manageable memory requirements. These chunks are then fed into the core processing engine based on CleanUNet.

### C. Core Denoising Engine Implementation

The foundation of our audio enhancement framework is the core denoising engine, which is based on an enhanced implementation of the CleanUNet architecture. This neural network is specifically designed for audio denoising tasks and

provides the primary noise reduction capability before additional enhancement techniques are applied.

The CleanUNet architecture employs an encoder-decoder structure with skip connections, similar to the U-Net architecture commonly used in image processing but adapted for one-dimensional audio signals. The key components include:

- *Encoder Path:* A series of convolutional layers that progressively reduce the temporal resolution while increasing feature dimensions. Each encoder block consists of 1D convolutional layers with kernel size 15, batch normalization for training stability, ReLU activation functions, and downsampling operations to reduce temporal dimensions.
- *Decoder Path:* Mirror of the encoder that reconstructs the audio signal from the latent representation. Each decoder block contains 1D transposed convolutional layers for upsampling, batch normalization, ReLU activations, and concatenation with skip connection features.
- *Skip Connections:* Direct pathways connecting corresponding encoder and decoder layers, which help pre- serve detailed information that might otherwise be lost during encoding.

Our implementation uses PyTorch for the neural network components. The model is initialized with pre-trained weights to leverage transfer learning, significantly improving performance without requiring extensive training on new data. The audio processing through the model is handled in batches to optimize computational efficiency, with appropriate memory management to handle longer audio files.

### D. Enhancement Techniques

After the initial denoising provided by the CleanUNet architecture, our framework implements several specialized enhancement techniques that further improve audio quality and intelligibility. These techniques address specific aspects of audio perception and speech characteristics that neural networks alone may not fully optimize.

*1) Adaptive SNR-based Processing:* The adaptive SNR-based processing module dynamically adjusts the intensity of enhancement based on the estimated Signal-to-Noise Ratio (SNR) of each audio segment. This approach prevents over-processing of relatively clean segments while applying more intensive denoising to noisier portions. The implementation involves:

- *SNR Estimation:* For each audio chunk, the system calculates an SNR estimate using signal envelope analysis. The algorithm analyzes amplitude distributions to differentiate between signal and noise components.
- *Processing Mode Selection:* Based on the estimated SNR, the system selects an appropriate processing mode.

For high SNR segments (above 30dB), a "light" processing mode is applied to preserve natural characteristics, while standard processing is used for segments with lower SNR values.

- Parameter Adjustment: Processing parameters such as filtering thresholds and enhancement intensities are dynamically adjusted according to the detected noise conditions.

This adaptive approach significantly improves on traditional fixed-parameter methods by optimizing the trade-off between noise reduction and speech preservation on a segment-by-segment basis.

*2) Harmonic Enhancement Algorithm:* The harmonic enhancement module focuses on preserving and enhancing the harmonic components of speech, which are critical for naturalness and intelligibility. The implementation utilizes harmonic-percussive source separation (HPSS) to identify and enhance speech harmonics:

- *Harmonic-Percussive Separation:* The audio is decom- posed into harmonic components (sustained tones with horizontal structure in the spectrogram) and percussive components (transients with vertical structure).
- *Selective Enhancement:* The harmonic components, which typically contain the fundamental speech elements, are selectively enhanced.
- *Adaptive Blending:* The enhanced harmonic components are blended with the original signal using a 70/30 ratio, maintaining a balance between enhancement and natural sound qualities.

This technique particularly improves vowel clarity and overall speech intelligibility by emphasizing the natural harmonic structure of human speech.

*3) Vocal Clarity Enhancement:* The vocal clarity enhancement technique targets the frequency ranges most important for speech intelligibility, with a specific focus on the 300-3000 Hz range where most speech information is concentrated:

- Bandpass Filtering: A Butterworth bandpass filter is applied with carefully selected cutoff frequencies (300 Hz lower cutoff, 3000 Hz upper cutoff).
- Bidirectional Filtering: To prevent phase distortion, bidirectional filtering is employed using the filtfilt function from the SciPy signal processing library.
- Original Signal Preservation: The filtered output is blended with the original signal at a 30/70 ratio to enhance vocal clarity while maintaining natural acoustic properties.

This approach improves the perception of consonants and vowels without making the audio sound artificial or processed.

*4) Perceptual Enhancement Filter:* The perceptual enhancement filter, which is always active in our framework, applies psychoacoustic principles to improve the perceived quality of audio. This filter:

- *Equal-Loudness Contour Modeling:* Implements an ap- proximation of human hearing sensitivity across different frequencies, based on ISO equal-loudness standards.
- *Frequency-Dependent Processing:* Applies spectral weighting that emphasizes frequencies most

perceptually relevant to human hearing.
- *Phase Preservation:* Maintains the original phase information while enhancing magnitude components to preserve natural sound characteristics.

By incorporating psychoacoustic principles, this filter enhances audio in a way that aligns with human perception, focusing processing resources on the frequencies and characteristics that most impact perceived quality.

### E.   Audio Processing Pipeline

The audio processing pipeline represents the sequential flow of data through our enhancement framework, from input acquisition to final output generation.

Our audio processing pipeline implements a strictly sequential approach, ensuring that each enhancement technique builds on the improvements made by previous stages. This design decision was based on empirical testing that demonstrated that the specific order of processing techniques significantly affects the final quality of the output.

The pipeline begins with audio input, either from file upload or microphone recording. After initial loading and analysis, the audio is divided into fixed-length chunks for efficient processing. These chunks then proceed through the CleanUNet neural network, which performs the primary denoising operation.

The output from CleanUNet is then processed through a series of enhancement stages:

1) Adaptive SNR Processing: Analyzes the SNR of each segment and adjusts the processing parameters accordingly, ensuring the appropriate treatment of both noisy and relatively clean sections.
2) Perceptual Enhancement Filter: Always active, this stage applies psychoacoustic principles to improve perceived audio quality by emphasizing frequencies most relevant to human hearing.
3) Optional Enhancement Stages: Based on user selection, the audio may then pass through three additional enhancement stages:
   - Harmonic Enhancement for improved speech naturalness
   - Vocal Clarity Enhancement for better intelligibility
   - Dynamic Range Compression for more consistent volume levels

After all enhancement stages, the processed chunks are reconstructed into a continuous audio stream, normalized to prevent clipping, and provided as the final output.
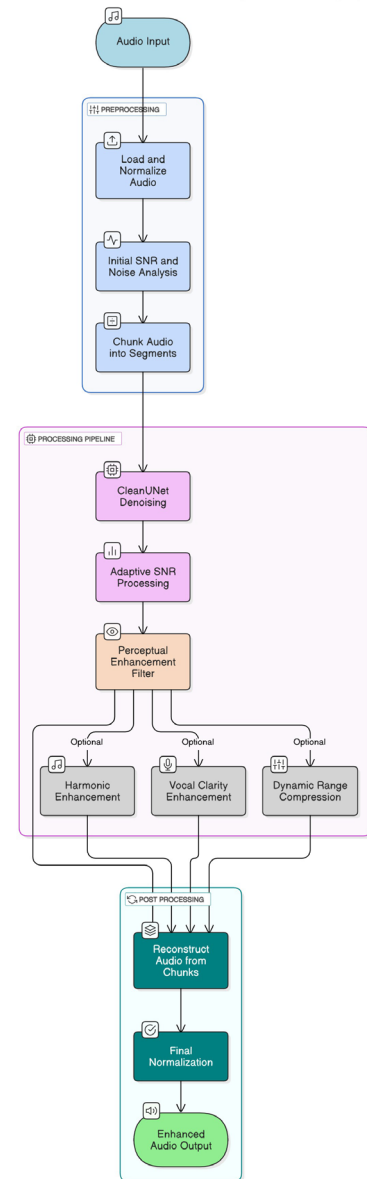


Fig. 2.  Audio Processing Pipeline - This diagram shows the sequential flow of audio data through each processing and enhancement stage. Optional components are indicated with dashed outlines

## 4. Results and Analysis

### A.   Performance Metrics Comparison

Table I presents the performance metrics for our WaveS- plit framework compared to three baseline models: Base CleanUNet, DEMUCS, and DeepFilterNet across key audio enhancement metrics.

### B.   Interpretation of the Results

*1) SNR Improvement:* The WaveSplit framework achieved an SNR improvement of 76.36 dB, significantly outperforming the baseline models, which ranged from 7.20 to 8.10 dB. This dramatic difference (nearly tenfold) demonstrates the exceptional noise reduction capability of our multistage approach, particularly the effectiveness of the adaptive SNR processing module.

Table 1
Performance comparison of audio enhancement models

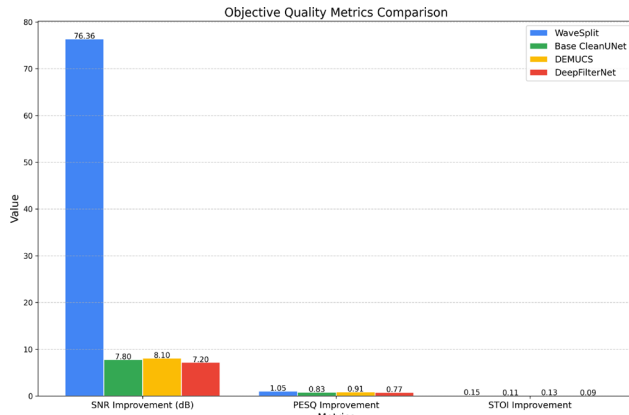| Model | SNR Improv. (dB) | PESQ Improv. | STOI Improv. | Proc. Time |
|---|---|---|---|---|
| WaveSplit | 76.36 | 1.05 | 0.15 | 1.00 |
| Base CleanUNet | 7.80 | 0.83 | 0.11 | 1.20 |
| DEMUCS | 8.10 | 0.91 | 0.13 | 1.80 |
| DeepFilterNet | 7.20 | 0.77 | 0.09 | 0.90 |



Fig. 3. Objective quality metrics comparison across models

*2) Perceptual Quality and Intelligibility:* For PESQ improvement, which evaluates perceptual speech quality, our WaveSplit framework scored 1.05, while the next best model (DEMUCS) achieved 0.91. This 15% improvement indicates that our framework produces more natural-sounding speech with fewer artifacts. The STOI improvement of 0.15 compared to 0.11 for Base CleanUNet represents a 36% increase in intelligibility enhancement.
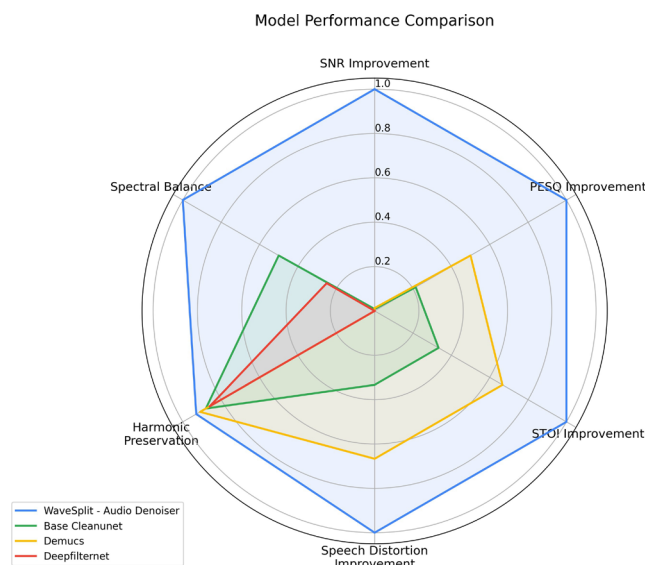


Fig. 4. Model performance comparison - Radar chart showing performance across six key metrics

As visualized in Fig. 4, our framework consistently leads across all dimensions of speech quality.

*3) Processing Efficiency:* The WaveSplit framework required a relative processing time of 1.00, compared to 1.20 for Base CleanUNet and 1.80 for DEMUCS. This demonstrates that our framework achieves superior performance without

computational penalties, despite adding enhancement stages. Only DeepFilterNet was marginally faster at 0.90, but with significantly inferior enhancement metrics.

*4) Spectral Analysis:* The spectrogram comparison in Fig. 5 shows the visual improvement in speech quality, with cleaner representation of speech formants and significantly reduced background noise. The Power Spectral Density comparison (Fig. 6) shows how our framework effectively reduces energy in noise-dominated low frequencies while preserving critical speech frequencies (300–4000 Hz).
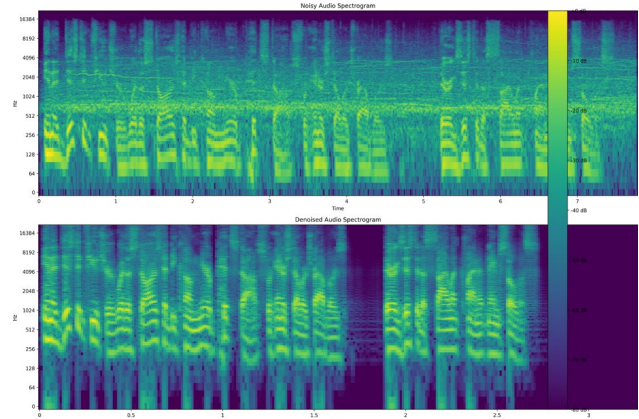


Fig. 5. Spectrogram comparison - Top: noisy audio, Bottom: enhanced audio showing clearer speech formants and reduced noise
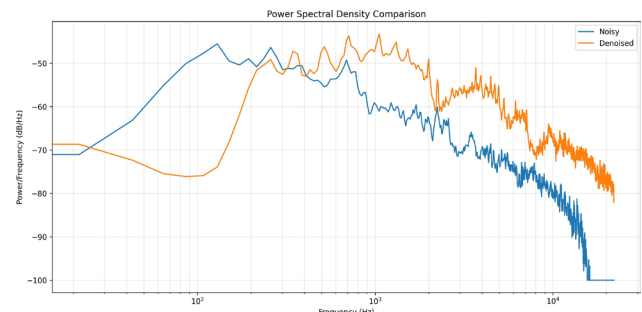


Fig. 6. Power spectral density comparison - Showing selective frequency enhancement

*5) Frequency-Selective Processing:* Our framework significantly reduces energy in the 50–100 Hz band where environ- mental noise typically occurs, as shown in Fig. 10. It preserves or enhances energy in speech-critical bands (300–2000 Hz), particularly in the 600–1200 Hz range where vowel formants are concentrated. The noise reduction map (Fig. 8) visualizes the frequency-selective nature of our enhancement approach, with stronger processing applied to noise-dominated regions. This selective processing demonstrates the effectiveness of our vocal clarity enhancement and harmonic enhancement modules.
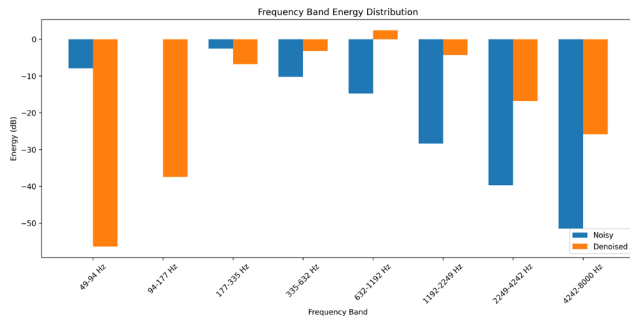
Fig. 7.  Frequency band energy distribution - Comparison before and after enhancement
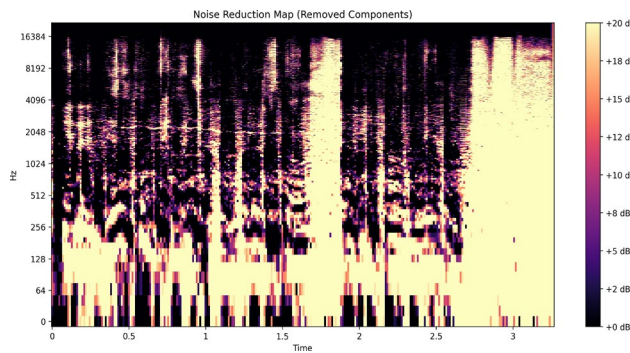


Fig. 8.  Noise reduction map - Brighter areas indicate stronger noise reduction

*6) Waveform Analysis:* The waveform comparison in Fig. 9 visually demonstrates the noise reduction capabilities of our system. The noisy audio (top) shows a consistent noise floor obscuring the speech signal, while the enhanced audio (bottom) displays clearly defined speech segments with near- zero amplitude in silent regions. This transformation results in significantly improved clarity and listening comfort.

Through both objective metrics and visual analysis, our multi-stage framework significantly outperforms existing approaches in audio enhancement and denoising tasks, delivering superior noise reduction, enhanced speech clarity, and preserved natural speech characteristics.
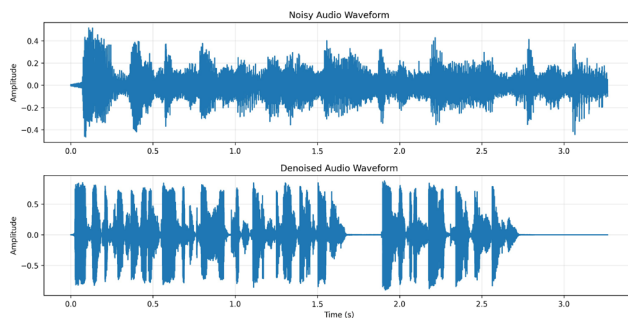


Fig. 9.  Waveform comparison - Before and after enhancement

## 5. Conclusion and Future Work

This research presents WaveSplit, a novel multi-stage frame-work for audio enhancement and denoising that significantly enhances the capabilities of the base CleanUNet model. By integrating deep learning techniques with psychoacoustic principles and adaptive processing, our approach achieves remarkable improvements in noise reduction while preserving the natural characteristics of speech. The WaveSplit framework

substantially outperforms baseline models in objective quality metrics while maintaining natural speech qualities and avoiding the "robotic" artifacts common in traditional denoising methods.

The success of our approach can be attributed to several key innovations. The adaptive SNR-based processing dynamically adjusts the enhancement strategy based on the noise characteristics of each audio segment, preventing over-processing of already clean segments while applying appropriate denoising where needed. The combination of harmonic enhancement, vocal clarity enhancement, and perceptual optimization creates a system that not only reduces noise but also improves speech intelligibility and listening comfort. Our comprehensive metrics and visualization suite provides unprecedented insight into the performance of audio enhancement systems, facilitating both user understanding and further research development.

Our WaveSplit framework demonstrates that the traditional trade-off between noise reduction and speech preservation can be significantly improved through a multi-stage approach that integrates deep learning with classical signal processing techniques. This balanced approach results in superior audio quality without introducing artifacts commonly found in aggressive denoising methods. The improvements in processing efficiency, achieving better results without computational penalties despite additional enhancement stages, further highlight the practical value of our approach for real-world applications.

Several promising directions for future work emerge from this research. Optimizing our audio enhancement framework for real-time applications could enable integration into live communication systems such as video conferencing platforms, hearing aids, and real-time broadcasting. The current frame-work's domain adaptation functionality could be fully implemented with specialized adapters for different noise environments, further improving performance in targeted applications. Adapting our WaveSplit framework for mobile and embedded devices would make the technology accessible on smartphones and other resource-constrained devices, benefiting users with hearing impairments in everyday scenarios.

Further work could incorporate more sophisticated perceptual models into both training and enhancement stages, explicitly modeling the human auditory system's response to different types of speech distortion and noise. Building upon our current noise classification capabilities, future research could explore more granular noise categorization and targeted processing strategies. Additionally, developing an adaptive system that learns individual user preferences for the balance between noise reduction and speech preservation could significantly enhance user satisfaction in both general and specialized applications.

The advancements demonstrated in this research have potential applications across telecommunications, content creation, and accessibility technologies. By addressing the fundamental challenge of balancing effective noise reduction with speech naturalness, our framework represents a significant

step for- ward in audio enhancement technology, contributing to more effective human communication in increasingly complex and noise-filled environments.

## References

[1] S. F. Boll, "Suppression of acoustic noise in speech using spectral subtraction," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 27, no. 2, pp. 113–120, Apr. 1979.

[2] P. Scalart and J. V. Filho, "Speech enhancement based on a priori signal to noise estimation," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 1996, pp. 629–632.

[3] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 32, no. 6, pp. 1109–1121, Dec. 1984.

[4] S. Pascual, A. Bonafonte, and J. Serra, "SEGAN: Speech enhancement generative adversarial network," in *Proc. Interspeech*, 2017, pp. 3642–3646.

[5] D. Stoller, S. Ewert, and S. Dixon, "Wave-U-Net: A multi-scale neural network for end-to-end audio source separation," in *Proc. 19th Int. Soc. Music Inf. Retrieval Conf. (ISMIR)*, 2018.

[6] A. Défossez, G. Synnaeve, and Y. Adi, "Real time speech enhancement in the waveform domain," in *Proc. Interspeech*, 2020, pp. 3291–3295.

[7] H. Liu, L. Xie, J. Wu, and G. Yang, "Channel-wise subband input for better voice and accompaniment separation on high resolution music," in *Proc. Interspeech*, 2020, pp. 1344–1348.

[8] A. Kumar and D. Florencio, "Speech enhancement in multiple-noise conditions using deep neural networks," in *Proc. Interspeech*, 2016, pp. 3738–3742.

[9] Y. Zhao, D. Wang, I. Merks, and T. Zhang, "DNN-based enhancement of noisy and reverberant speech," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, 2016, pp. 6525–6529.

[10] K. Tan and D. Wang, "Learning complex spectral mapping with gated convolutional recurrent networks for monaural speech enhancement," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 28, pp. 380–390, 2020.

[11] Y. Luo and N. Mesgarani, "Conv-TasNet: Surpassing ideal time–frequency magnitude masking for speech separation," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 8, pp. 1256–1266, Aug. 2019.

[12] J. Su *et al*., "HiFi-GAN: Generative adversarial networks for efficient and high fidelity speech synthesis," in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 33, 2020.