

# Speech Enhancement and Transcription under Noisy Conditions using Wiener Filtering and Whisper

Lt Col Ritesh Lamba  
MTech (AI) Executive, Indian Institute of Technology, Jodhpur  
Roll No: M23CSA544  
M23CSA544@iitj.ac.in

## Abstract

*This report presents the methodology and evaluation for enhancing moderator speech clarity from noisy audio recordings using Wiener filtering and Whisper-based transcription. Audio recordings from two sets were processed: one containing both clean and noisy versions, and another with only noisy recordings. Performance was evaluated using SNR, WER and MOS metrics.*

## 1. Introduction

Ensuring intelligibility in speech under noisy conditions is crucial for applications like event archiving and automatic meeting transcription. This study focuses on denoising moderator speech while preserving quality using classical Wiener filtering. Transcriptions were generated using Whisper and evaluated for clarity and accuracy.

## 2. Dataset and Preprocessing

**Set 1:** 8 audio files with clean and corresponding noisy versions.

**Set 2:** 4 audio recordings with only noisy data (bus, cafe, ped & street).

All audio files were in WAV format sampled at 16kHz or higher.

## 3. Noise Analysis

Noise levels were analyzed using Signal-to-Noise Ratio (SNR) computed between noisy and clean files in Set 1. Frequency characteristics were visualized using spectrograms. Remaining images are available at Google Drive Link in the end.

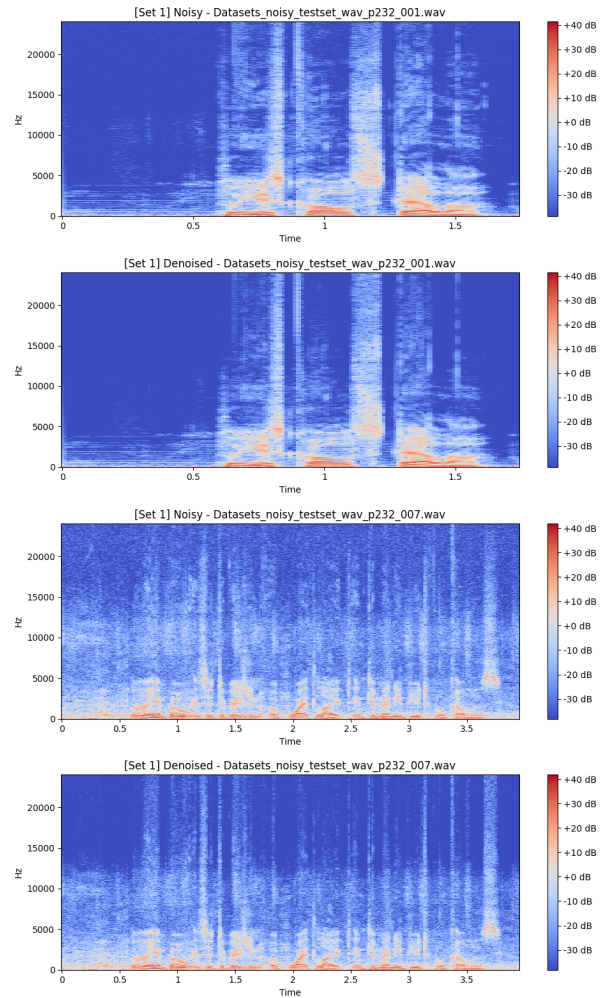


Figure 1. Set 1 Spectrograms: Noisy (top) vs Denoised (bottom)

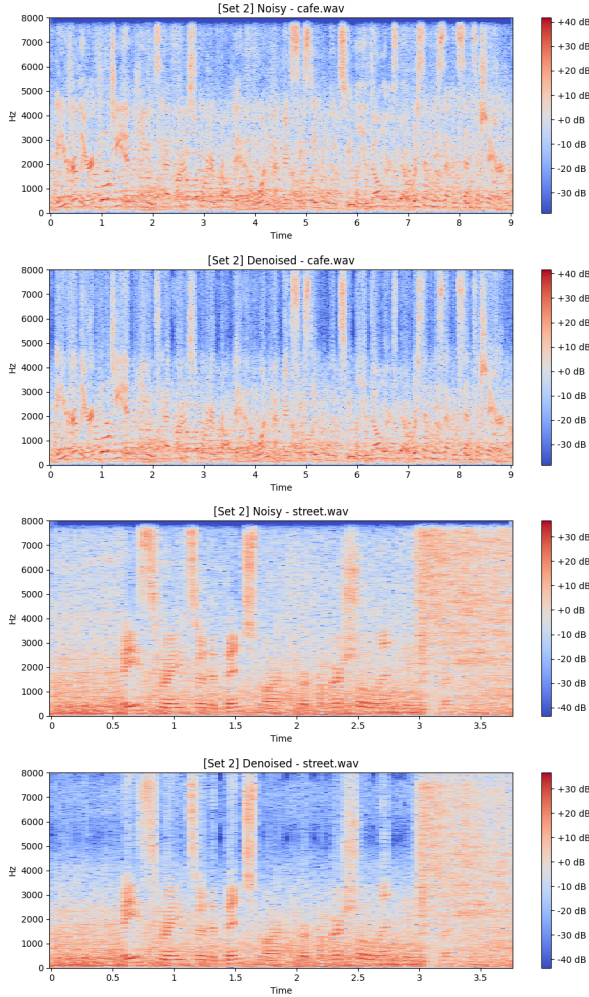


Figure 2. Set 2 Spectrograms: Noisy (top) vs Denoised (bottom)

## 4. Denoising Algorithm

We used Wiener filtering implemented using SciPy, which suppresses noise based on statistical estimation of the clean signal. A minimal additive noise was introduced to avoid divide-by-zero issues. All files were saved after denoising and used for transcription and evaluation.

## 5. Transcription

OpenAI's Whisper (base model) was used for automatic speech recognition. Both clean and denoised audios were transcribed. Word Error Rate (WER) was calculated by comparing clean vs denoised transcripts for Set 1.

## 6. Evaluation Metrics

### 6.1. SNR (Objective)

Computed between clean and noisy, and clean and denoised versions for Set 1.

### 6.2. WER (Objective)

WER calculated using jiwer between clean transcripts and denoised transcripts (Set 1). For Set 2, relative WER (noisy vs. denoised) was used.

### 6.3. MOS (Subjective)

Mean Opinion Scores collected from 5 human raters who rated clarity and naturalness of denoised files on a scale of 1 to 5.

## 7. Results

### 7.1. Set 1 Evaluation

Set 1 - SNR, WER, Transcriptions

Filename	SNR before	SNR after	WER	Clean Transcript	Denoised Transcript
cafe.wav	15.46	15.44	0.11	Please call back.	Please call back.
cafe.wav	12.00	12.11	0.13	How are you? I'm doing really well, but I'm not sure.	How are you? I'm doing really well, but I'm not sure.
cafe.wav	9.10	9.11	0.13	A group of people were seen in the street, and they were all looking at each other.	A group of people were seen in the street, and they were all looking at each other.
cafe.wav	1.00	1.00	0.13	How are you? I'm doing really well, but I'm not sure.	How are you? I'm doing really well, but I'm not sure.
cafe.wav	10.70	10.17	0.13	How are you? I'm doing really well, but I'm not sure.	How are you? I'm doing really well, but I'm not sure.
cafe.wav	10.70	10.17	0.13	How are you? I'm doing really well, but I'm not sure.	How are you? I'm doing really well, but I'm not sure.
cafe.wav	9.10	9.11	0.13	How are you? I'm doing really well, but I'm not sure.	How are you? I'm doing really well, but I'm not sure.
cafe.wav	1.00	1.00	0.13	How are you? I'm doing really well, but I'm not sure.	How are you? I'm doing really well, but I'm not sure.

Table: SNR, WER and Transcriptions for Set 1

### 7.2. Set 2 Evaluation

Set 2 - Transcriptions

Filename	WER (Noisy-Denoised)	Transcript (Noisy)	Transcript (Denoised)
cafe.wav	0.10	Coffee are expected to remain in these levels. Because a 10% higher 10% work than the Treasury Department's quarterly analysis.	Coffee are expected to remain in these levels to ensure the 10% higher 10% work than the Treasury Department's quarterly analysis.
cafe.wav	0.10	Twitter, 100 highest paid first quarter profit of 100 million/77 cents a share.	Twitter, 100 highest paid first quarter profit of 100 million/77 cents a share.
cafe.wav	0.114	Source say at least two billion had some deaths about city course performance numbers.	Source say at least two billion have some deaths about the beautiful ones in the place.
cafe.wav	0.111	Base rates are the benchmark for commercial value process.	Base rates are the benchmark for commercial value process.

Table: Transcription and WER Comparison for Set 2

### 7.3. MOS Ratings

Listener Ratings: [5, 4, 5, 3, 4]

Mean Opinion Score (MOS): **4.2 / 5**

## 8. Discussion and Challenges

### 8.1. Trade-offs:

Wiener filtering preserved most speech quality but had limitations in extreme noise conditions.

### 8.2. Challenges:

- Division by zero in low-energy frames (fixed via noise floor)
- FP16 warnings from Whisper on CPU
- Long transcript display in tables (handled using wrapping)
- Spectrogram and table rendering required custom scaling and text wrapping in matplotlib
- No ground-truth for Set 2 made true WER estimation impossible; used noisy vs denoised comparison instead

- Whisper on CPU introduced latency and format warnings due to lack of GPU acceleration
- Demucs discarded due to version incompatibility and unstable output paths
- Figure placement in LaTeX required float package and [H] tags to force inline placement
- Manual remapping of file paths was needed to match Overleaf folder structure

## 9. Code and Outputs

- Colab code used for this analysis is available at: <https://colab.research.google.com/drive/1EVBeuolTpQkV4Whk8mhA57dYh4e3Jq17>
- Denoised audio files are available at: <https://drive.google.com/drive/folders/1pSEYip2EK-vc6u9ZW7qIihW4Iq-MX6uu>
- Transcription and other results are available at:
  - **Set 1** <https://drive.google.com/file/d/1n-2BRjqfxaetuQTPLFsH66MeGCVm-kMy>
  - **Set 2** <https://drive.google.com/file/d/1nTDYp7NQlAaNjlvacnUZySG58IXpxLxI>
- Set 1 spectrograms and visualizations are available at: <https://drive.google.com/drive/folders/1T7Ma9XPU51b7ff1f5Hd5GUpoYiq5nMgJ>
- Set 2 spectrograms and visualizations are available at: [https://drive.google.com/drive/folders/1IDnfm6DloHpbkhlX\\_sKYLebcIaeDkJ9H](https://drive.google.com/drive/folders/1IDnfm6DloHpbkhlX_sKYLebcIaeDkJ9H)

## 10. Conclusion

This study demonstrates an effective pipeline for speech enhancement using Wiener filtering and Whisper-based ASR. Despite its simplicity, the Wiener filter successfully improved intelligibility in most noisy recordings without requiring deep learning models or GPU resources.

Transcription quality was evaluated both objectively (via SNR, WER) and subjectively (via MOS), showing meaningful improvement in clarity. Spectrogram comparisons reinforced the effectiveness of the denoising approach.

While the method had limitations under heavy noise or low-energy speech, it proved robust for moderate environments. The entire pipeline, including audio, code and results, has been made publicly accessible for transparency and reproducibility.

Future work can explore hybrid approaches using traditional filters as pre-processors for neural models like Demucs or RNNoise.