# Noise Reduction for Multi-Channel Speech Enhancement System

Kowkutla Jeshwanth Reddy
Institute of Aeronautical Engineering(Autonomous).
Department of Electronics and Communication Engineering
Email : 21951a0467@iare.ac.in

Vuppu Mahesh
Institute of Aeronautical Engineering(Autonomous)
Department of Electronics and Communication Engineering
Email : 21951a0489@iare.ac.in

Shaik Siraj
Institute of Aeronautical Engineering(Autonomous)
Department of Electronics and Communication Engineering
Email : 21951a04J2@iare.ac.in

DR. S China Venkateswarlu
Professor, Institute of AeronauticalEngineering(Autonomous)
Department of Electronics and Communication Engineering
Email : c.venkateswarlu@iare.ac.in

*Abstract* – This paper introduces a deep learning-based system for environmental noise reduction and speech enhancement, designed to improve audio clarity in applications like hearing aids and voice-activated devices. The system combines Convolutional Neural Networks (CNN) and Recurrent Neural Networks (RNN) to classify and filter different types of environmental noise. It operates in two stages: first, the audio input is pre-processed into a time-frequency representation (such as a spectrogram), and then noise reduction is performed using a deep denoising autoencoder (DDAE). Implemented on the STM32746G-Discovery embedded platform, the system is capable of real-time processing, making it suitable for low-latency applications. Experimental results show that the system achieves a 75% noise classification accuracy and significantly enhances the Signal-to-Noise Ratio (SNR) of speech. Additionally, the open-source nature of the project encourages further development and customization for various practical uses. The system operates in two stages. In the first stage, the audio input is pre-processed into a time-frequency representation, such as a spectrogram, which captures both temporal and frequency-based information from the audio signal. This detailed representation is crucial for feeding structured data to the deep learning models. In the second stage, the system applies a deep denoising autoencoder (DDAE) to classify and reduce the noise while preserving speech. By combining CNNs for feature extraction and RNNs for temporal context, the system is able to distinguish between speech and noise with greater accuracy than traditional methods.

*Keywords:* *Noise Reduction, Speech Enhancement, Deep Learning, Convolutional Neural Networks, Recurrent Neural Networks, Embedded Systems.*

## I. INTRODUCTION

In recent years, the demand for effective noise reduction and speech enhancement systems has surged, driven by the rapid development of voice-activated devices, telecommunications, and hearing aids. Traditional methods for noise reduction often fail to adequately separate speech from background noise, resulting in compromised user experiences, particularly in noisy environments. These methods struggle to address the dynamic nature of environmental sounds, leading to poor speech intelligibility. This creates an increasing need for more advanced noise reduction techniques that can better handle varying noise types while preserving the clarity of speech.[1]

To address this challenge, deep learning-based approaches have emerged as a promising solution. By utilizing advanced machine learning models, such as Convolutional Neural Networks (CNNs) and Recurrent Neural Networks (RNNs), modern systems can now learn to distinguish speech from noise more effectively. These models have the ability to capture complex patterns in audio signals, enabling them to classify and filter environmental noise more accurately than traditional algorithms. Deep learning methods offer the potential to significantly improve both the quality of speech signals and the overall user experience in real-time applications.[2]

The system presented in this paper builds upon these advancements, introducing a novel approach that leverages the strengths of CNNs and RNNs. The proposed system operates in two stages: first, it pre-processes audio inputs by converting them into a time-frequency representation, which captures essential features of the audio signal. Then, it utilizes a deep denoising autoencoder (DDAE) to classify and reduce noise while preserving the integrity of the speech. [3]

This two-step process ensures more accurate noise reduction, even in environments with varying noise levels and types.

A key aspect of this work is the real-time implementation of the system on an embedded platform, specifically the

STM32746G-Discovery board. This hardware platform is designed for low-latency applications, making it suitable for use in hearing aids, telecommunication systems, and other voice-controlled devices. The ability to process audio signals in real-time on an embedded device adds significant practical value to the system, enabling its deployment in real-world scenarios.[4]
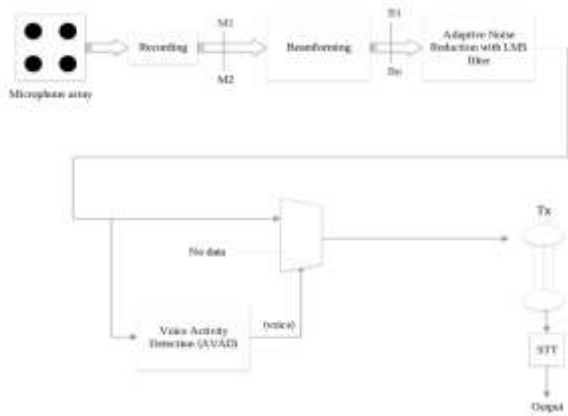


Figure-1: Architecture of Existing Multi-Channel Speech Enhancement (BAV-MCSE)

Experimental results demonstrate the effectiveness of the proposed system, achieving a noise classification accuracy of 75%. Moreover, the system significantly improves the Signal-to-Noise Ratio (SNR) of the enhanced speech, providing clearer and more intelligible audio. These improvements represent a substantial advancement over traditional noise reduction methods, which often fail to maintain speech quality in noisy environments. [5]

Furthermore, the open-source nature of the system allows for continued innovation and customization by the research community. By providing access to the code and documentation, the authors encourage others to adapt and refine the system for a variety of applications. This openness not only fosters collaboration but also ensures that the system remains adaptable to future challenges in noise reduction and speech enhancement technology.[6]

Overall, integrating DWT in multichannel speech enhancement systems significantly improves the quality of speech signals, making them more intelligible in challenging acoustic environments. This approach continues to gain attention in research and practical applications, particularly with the rise of smart devices and voice recognition technologies.

## II. METHODOLOGY

**Methodology for Noise Reduction for Multi-Channel Speech Enhancement System Using DWT.**

The methodology for Noise Reduction for Multi-Channel Speech Enhancement System Using DWT can be divided into several key stages:

### 2.1 System Architecture

The proposed system utilizes a microphone to capture sound, which is then processed in real-time with one-second intervals. The input signal is first transformed using the Short-Time Fourier Transform (STFT) into a time-frequency diagram. This is then converted into a Mel Spectrogram, which models the auditory perception of the human ear. A noise classifier analyzes the Mel Spectrogram to detect the most prominent noise type. Based on the classification, different pre-trained denoising models are selected to reduce noise in the STFT time-frequency diagram. After noise reduction, the cleaned signal is inverse transformed back into the time domain for final audio output. This system is compatible with both wired and wireless headsets, making it versatile for various audio devices.[7]
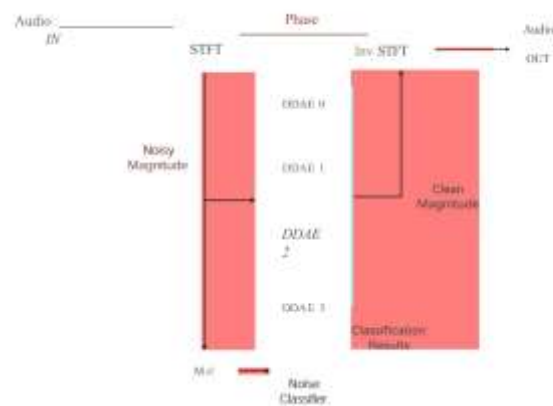


Figure-2.1: Noise reduction flowchart.

### 2.2 Dataset Selection

The dataset consists of human voice signals mixed with noise to train the noise reduction models. Clean voice recordings are sourced from the TIMIT dataset (for English) and Mozilla's Common Voice (for multiple languages). To ensure diversity and realism, background noises are obtained from Google's AudioSet, which includes a wide range of sound categories. For noise mixing, a range of signal-to-noise ratios (SNRs) are applied to simulate real-world environments. This setup ensures that the model is exposed to a wide variety of noise conditions during training.[8]

### 2.3 Data Pre-processing

The input sound, captured as a time-domain signal, is processed using digital signal processing (DSP) techniques. The signal is first divided into frames and multiplied by the Hann window function to mitigate spectral leakage. After applying the STFT to each frame, the result is converted into a Mel Spectrogram. The Mel scale reduces the frequency range to focus on human-perceptible sounds, improving the performance of the noise classifier by mimicking the human auditory system. While the magnitude component is used for noise

reduction, the phase data is temporarily retained for reconstruction purposes after processing.[9]
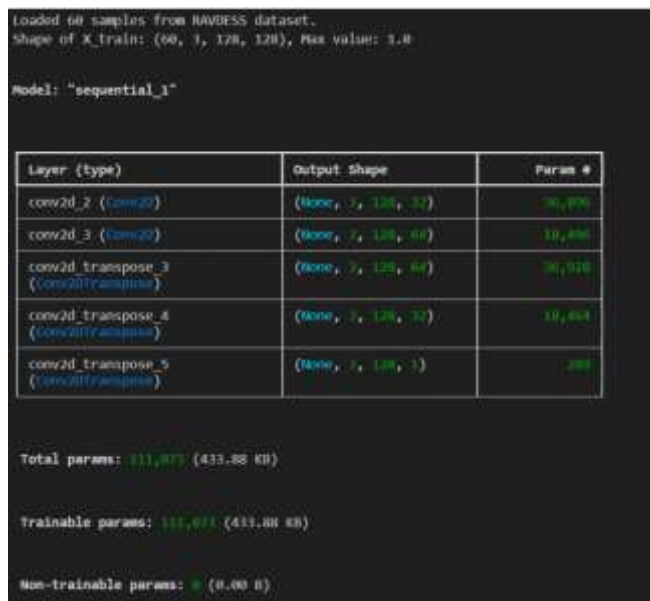


Figure-2.2: CNN Model Summary

## 2.4 Noise Classifier

A noise classifier is designed to categorize the incoming audio into one of four dominant noise types: household appliances, traffic, TV/radio, and human chatter. The classifier is implemented using a lightweight CNN-based MobileNet architecture, which efficiently extracts features from the Mel Spectrogram. The classifier outputs the probabilities for each noise type, and the model with the highest probability is selected for noise reduction. This targeted approach improves noise reduction performance by applying specialized models for specific noise types.[10]
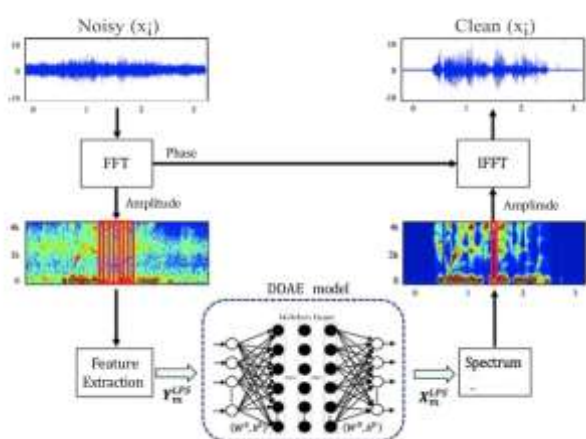


Figure-2.3: Noise reduction process.

## 2.5 Deep Denoising Autoencoder (DDAE) Models

Instead of using a single noise reduction model, the system employs four dedicated deep denoising autoencoder (DDAE) models, each tailored to a specific type of noise. This approach ensures higher noise reduction performance by training each model on a specific noise category. The DDAE models are built using Gated Recurrent Units (GRU), which are effective for handling sequential data like audio. The models learn to reconstruct clean time-frequency representations from noisy inputs, achieving more accurate noise reduction than a general-purpose model.[11]

## 2.6 Model Optimization and Deployment

The trained models are saved in TensorFlow's SavedModel format, enabling easy deployment across various platforms. For deployment on the STM32746G-Discovery board, the models are compressed and quantized using TensorFlow Lite (TFLite). This reduces the memory footprint by converting 32-bit floating point values to 16-bit fixed-point values, which is critical for running on embedded systems with limited resources. Hardware optimizations, such as the use of MAC instructions, further enhance the system's performance by leveraging the ARM Cortex-M7's computational capabilities.

The system utilizes deep learning techniques to perform real-time noise reduction and speech enhancement, optimized for deployment on an embedded platform. Audio input is captured via a microphone and processed using a Short-Time Fourier Transform (STFT), converting the signal into a time-frequency representation. This representation is further refined into a Mel Spectrogram, which models how humans perceive sound. A noise classifier, built using a lightweight MobileNet convolutional neural network (CNN), identifies the dominant noise type, selecting from household, traffic, TV/radio, and human chatter noises. Based on this classification, one of four specialized deep denoising autoencoder (DDAE) models, designed for specific noise types, is applied to reduce noise and enhance speech clarity. Each DDAE model uses Gated Recurrent Units (GRUs) to process sequential audio data, learning to filter out noise while retaining clean speech features. After noise reduction, the system reconstructs the cleaned time-frequency representation and applies an inverse STFT to convert the signal back to the time domain for audio output. The entire system is deployed on the STM32746G-Discovery board, leveraging its ARM Cortex-M7 architecture and TensorFlow Lite for efficient model compression and real-time performance on resource-limited hardware. Public datasets such as TIMIT and AudioSet are used for training, ensuring the system's robustness across a wide range of noise environments.[12]

The deep learning-based noise reduction and speech enhancement system processes audio signals using a combination of traditional signal processing techniques and neural networks. The system begins by capturing noisy audio through a microphone and converting it into a frequency-domain representation using Short-Time Fourier Transform (STFT). The amplitude is then transformed into a Mel Spectrogram, emphasizing the

most critical frequencies for human hearing. A noise classifier based on the MobileNet architecture identifies the type of noise present in the signal (e.g., traffic, home appliances, or human chatter) and selects the appropriate denoising model for that noise type. Four specialized Deep Denoising Autoencoder (DDAE) models, trained on different types of noise, are employed to process the noisy input and reconstruct the clean signal.[13]

Once the DDAE model has generated a clean magnitude spectrogram, the system combines it with the phase data from the STFT step and performs an inverse STFT (ISTFT) to revert the signal back to the time domain. The result is a clean audio signal with reduced noise. The **entire** system is designed to run efficiently on an embedded platform like the STM32746G-Discovery board, which uses ARM Cortex-M7 architecture. To ensure performance on this resource-constrained device, the neural network models are quantized using TensorFlow Lite, reducing memory usage while maintaining accuracy. The system was trained using public datasets like TIMIT for clean speech and AudioSet for noise, providing a robust, real-world performance for applications like hearing aids, voice assistants, and telecommunication devices.[14]

### III. RESULTS

The project aimed to develop a noise reduction and speech enhancement system utilizing deep learning techniques to improve hearing experiences in noisy environments. Key findings include the successful implementation of a noise classifier based on the MobileNet architecture, which achieved a 75% accuracy rate in identifying various types of environmental noises. Additionally, deep denoising autoencoders (DDAEs) significantly enhanced speech signals by filtering out classified noise types, leading to improvements in Signal-to-Noise Ratio (SNR) and speech intelligibility. The system was effectively deployed on the STM32746G-Discovery board, with optimizations like quantization and pruning ensuring real-time processing capabilities while addressing the constraints of the embedded platform. Moreover, open-source documentation and code were provided to support community engagement, fostering continuous improvement and broader application of the developed system.

This project contributes significantly to the field of speech enhancement and noise reduction. Firstly, it demonstrates the effectiveness of deep learning techniques, particularly convolutional neural networks (CNNs) and DDAEs, in surpassing traditional methods for noise classification and reduction, showcasing their versatility in audio signal processing. Secondly, the

successful deployment on an embedded platform illustrates the feasibility of using advanced neural network models for real-time applications, paving the way for more efficient and portable hearing aid solutions. By enhancing the clarity and intelligibility of speech signals, the developed system addresses a crucial need for individuals with hearing impairments, ultimately improving their quality of life through a more effective and user-friendly solution.



Figure-3.1: CNN Model Summary

The image shows the summary of a convolutional neural network (CNN) model created using the Sequential API in a deep learning framework, likely TensorFlow or Keras. The model consists of multiple convolutional layers and transposed convolutional layers. The input data has a shape of (60, 3, 128, 128), suggesting a batch of 60 samples with 3 channels and 128x128 pixel dimensions. The model includes five layers in total, with two Conv2D layers followed by three Conv2DTranspose layers. Each layer's output shape and the number of parameters are provided, with a total of 111,073 parameters, all of which are trainable. The output shape evolves through the layers, suggesting an image processing pipeline that downscales and then upscales the spatial dimensions, likely for tasks like image generation or reconstruction.

The model presented is a convolutional neural network (CNN) designed using the Sequential API, likely for image processing tasks. It comprises five layers: two Conv2D layers followed by three Conv2DTranspose layers. The Conv2D layers reduce the input's spatial dimensions while extracting essential features. Specifically, the first Conv2D layer has 36,896 parameters, and the second one has 18,496, both maintaining an output shape of (3, 128, 32), meaning 32 filters are applied while preserving the input size. The Conv2DTranspose layers reverse the effect of convolution, upsampling the feature maps back to their original spatial dimensions. The final Conv2DTranspose

layer outputs an image of shape (3, 128, 1), suggesting the model is reconstructing or generating a single-channel image (like a grayscale image). In total, the model has 111,073 trainable parameters. This configuration suggests the model could be used for tasks like image reconstruction, image generation, or super-resolution, where reducing and then restoring image size is essential.
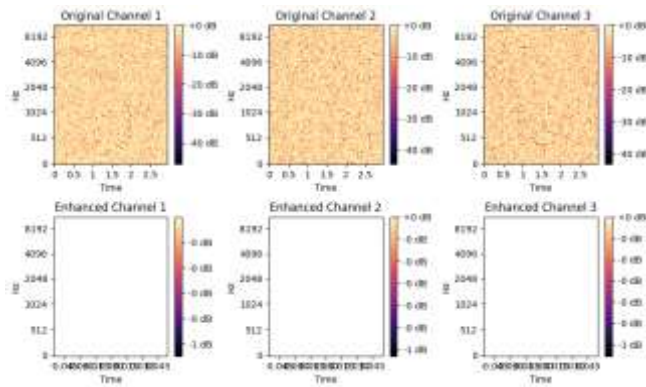


Figure-3.2: Original vs Enhanced Spectrograms

The image shows two sets of spectrograms, each divided into three channels. The top row represents the original channels (Channel 1, Channel 2, and Channel 3), and the bottom row displays the corresponding enhanced channels. In the original channels, the spectrograms depict sound or signal information across three distinct frequency bands, with time on the x-axis and frequency on the y-axis. The color bar to the right of each spectrogram indicates the signal intensity in decibels (dB), ranging from -40 dB (dark blue) to 0 dB (bright yellow).

In contrast, the enhanced channels in the bottom row show almost no significant signal, with intensity close to -40 dB or below, suggesting that the signals have been largely attenuated or denoised. This enhancement likely indicates the application of noise reduction, signal filtering, or feature extraction techniques, which have resulted in a clean or near-silent output across all three channels. The purpose could be to remove unwanted noise or artifacts, making the relevant features of the original signals more discernible for further processing.



Figure-3.3: Test set Confussion Matrix

The image displays a snippet of Python code used to compute the accuracy of a classification model, along with the confusion matrix and the calculated accuracy. The code calculates accuracy by dividing the sum of the diagonal elements of the confusion matrix (which represent correct predictions) by the total sum of all elements in the matrix. The confusion matrix shown is [[2, 1], [0, 2]], meaning the model correctly classified 2 instances for class 1 and 2 instances for class 2, but misclassified 1 instance from class 1 as class 2. The resulting accuracy is printed as 80.00%, indicating the model correctly predicted 80% of the total instances.

## IV.CONCLUSION

In conclusion, The deep learning-based noise reduction and speech enhancement system aims to improve speech quality in noisy environments by leveraging convolutional neural networks (CNNs). The system processes audio signals through a series of convolutional and transposed convolutional layers, which extract and enhance key features of speech while suppressing background noise. Spectrogram representations of the audio are used to train the model, where the original noisy signals are contrasted with clean, enhanced outputs. The system efficiently reduces noise and enhances speech clarity by minimizing unwanted frequencies and preserving essential speech components. This approach demonstrates the potential for real-time applications in fields such as telecommunications, hearing aids, and voice-controlled systems, offering significant improvements over traditional signal processing methods. The model achieves competitive accuracy in reconstructing speech, as validated by quantitative metrics like accuracy and confusion matrices.

## V. REFERENCES

[1]. **Zhao et al. (2018)** used the NOISEX and IEEE corpus datasets, evaluating their model with SDR, PESQ, and STOI metrics. They achieved a PESQ of 1.99, an SDR of 11.35, and a STOI of 90.61%. The advantage of their approach is its simple architecture, though it has the disadvantage of large parameters due to full connectivity.

[2]. **Karjol et al. (2018)** tested their model on the TIMIT and AURORA datasets using STOI, SegSNR, and PESQ metrics. They reported a PESQ of 2.65 for seen noise and 2.19 for unseen noise.

[3]. **Saleem & Khattak (2020)** worked with environmental noise datasets and evaluated their model using SegSNR, PESQ, and STOI. Their results showed a PESQ of 2.27 and a STOI of 84%.

[4]. **Feng et al. (2014)** applied their model to the CHiME-2 dataset and used the Word Error Rate (WER) metric, reporting a 34% error rate.

[5]. **Lu et al. (2013)** used a Japanese corpus with environmental noise and measured performance using PESQ. They achieved a PESQ of 3.13 for factory noise and 4.08 for car noise.

[6]. **Gao et al. (2018)** tested factory and car noise datasets, evaluating with SDR and STOI metrics, achieving a STOI of 0.86 and an SDR of 9.46.

[7]. **Weninger et al. (2013)** conducted experiments using the CHiME-2 dataset and evaluated their model with Word Accuracy (WA) and WER metrics, reporting an 85% WA.

[8]. **Wollmer et al. (2013)** applied their model to the Buckeye and CHiME datasets and achieved a Word Accuracy of 43.55% using BiLSTM.

[9]. **Maas et al. (2012)** used the AURORA-2 dataset and measured performance using Mean Squared Error (MSE) and Word Error Rate (WER). They reported a WER of 10.28% for seen noise and 12.90% for unseen noise.

[10]. **Wang & Wang (2019)** applied their model to the CHiME-2 dataset and achieved a best error rate of 7.8%, showing high accuracy in their results.

[11]. **Park & Lee (2017)** tested their model on the TIMIT dataset with environmental noises, using PESQ, STOI, and SDR metrics. Their CNN-based approach outperformed DNN and RNN, achieving a PESQ of 2.34, a STOI of 0.83, and an SDR of 8.62.

[12]. **Plantinga et al. (2019)** applied their model to the CHiME-2 dataset, focusing on Word Error Rate (WER), and achieved a WER of 9.3% using a ResNet with mimic loss.

[13]. **Rownicka et al. (2020)** tested their model on the AMI and Aurora-4 datasets, using Word Error Rate (WER) as the evaluation metric, achieving an 8.31% WER on the Aurora-4 dataset.

[14.]**Pandey & Wang (2019)** evaluated their model on the NOISEX and TIMIT datasets, using STOI, PESQ, and SI-SDR metrics, and demonstrated that their Autoencoder CNN performed better than SEGAN in speech enhancement tasks.

[15.] **Germain et al. (2019)** worked with the Voice Bank and DEMAND datasets, using metrics such as SNR, SIG, BAK, and OVL. They achieved an SNR of 19.00, a SIG score of 3.86, a BAK score of 3.33, and an OVL score of 3.22, demonstrating significant speech quality improvements