

Speech Preprocessing for Parkinson's Patients Using a Double Compact Deep Learning Model Approach

Chaitali Shamrao Raje, Dr. Pramodkumar. H. Kulkarni, Rupali Deshmukh

Department of Electronics & Telecommunication Engineering, Dr.D.Y.Patil Institute of Technology, Pimpri, Pune
Email: chaitaliraje16@gmail.com

Speech is severely hindered by Parkinson's disease, resulting in problems like low vocal volume, monotone speech, loose articulation, and irregular speech pace. This work examines voice preprocessing methods intended to improve Parkinson's disease patients' ability to communicate. While they reduce noise, traditional techniques like Wiener filtering and spectral subtraction can introduce artifacts. While they reduce statistical noise, machine learning techniques like Hidden Markov Models (HMM) and Gaussian Mixture Models (GMM) have variability issues. Promising techniques for enhancing speech quality are advanced deep learning methods like Denoising Autoencoders (DAEs) and Convolutional Neural Networks (CNNs), but they need a lot of data and processing power. For Parkinson's sufferers to communicate better overall and with greater intelligibility in their speech, these preprocessing approaches are essential.

Keywords: Speech preprocessing, Deep Learning, MATLAB.

1. Introduction

A critical first step in improving voice quality and intelligibility, particularly for people with speech difficulties, is speech signal preprocessing. Parkinson's disease is a neurodegenerative disease that worsens with time and impairs motor function. As a result, it can cause a variety of speech-related problems that are collectively referred to as Parkinsonian dysarthria. Maintaining one's quality of life requires effective communication, and Parkinson's patients can greatly benefit from increased speech intelligibility.

Parkinsonian Dysarthria

A number of speech impairments resulting from the motor deficiencies brought on by Parkinson's disease are known as Parkinsonian dysarthria. These anomalies consist of:

Reduced Vocal Loudness (Hypophonia): It can be challenging to understand patients when

they speak since they frequently speak softly [1].

Monotone Speech: Speech may sound flat and monotonous because it lacks the typical volume and pitch fluctuations [2].

Imprecise Articulation: Speech slurs and imprecise speech can be caused by a lack of coordination between speech muscles [3].

Variable Speech Rate: Individuals may speak too fast or too slow, or they may speak at an unpredictable pace [4].

Breathy or Hoarse Voice: Speech clarity may be further compromised by a breathy or harsh voice, which is a result of compromised vocal quality [5].

The respiratory, phonatory, and articulatory systems' muscles—which are involved in producing speech—are affected by Parkinson's disease-related decreased motor control, which is the cause of these speech problems.

Importance of Speech Signal Preprocessing

Speech signal preprocessing can significantly improve communication for Parkinson's sufferers by making speech more comprehensible and of higher quality. Speech signal preprocessing has the following main advantages:

Enhanced Intelligibility: Improving the lucidity of communication to facilitate comprehension for the audience [6].

Increasing the volume of soft speech to make it more comprehensible [7].

Normalized Speech Rate: Guaranteeing a steady and suitable speech rate [8].

Eliminating undesirable sounds and background noise that may impede the intelligibility of speech is known as noise reduction [9].

Preserving the naturalness of the improved speech while avoiding the addition of artificial artifacts is known as quality preservation [10].

2. Literature Review

Speech signal preprocessing for Parkinson's patients is a field of study aimed at mitigating the speech impairments associated with Parkinsonian dysarthria. This literature survey provides an overview of existing techniques, their effectiveness, and the gaps that remain to be addressed. By reviewing the current state of the art, we can identify areas for improvement and innovation shown in figure 1.

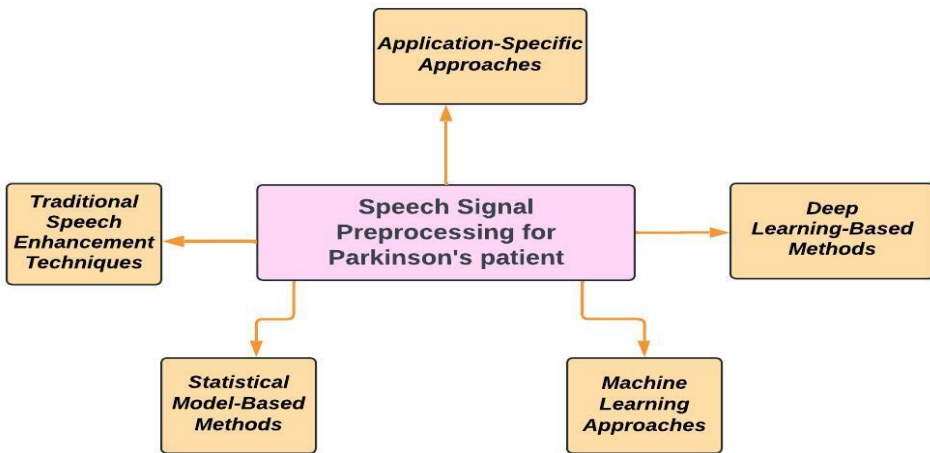


Fig 1: Existing techniques for speech signal preprocessing for Parkinson's Patient

A. Traditional Speech Enhancement Techniques

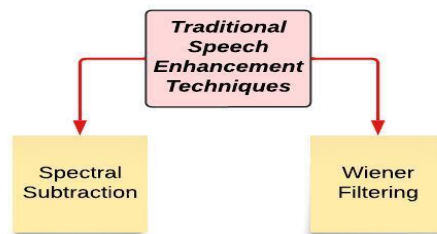


Fig 2: Traditional Speech Enhancement Techniques

a. Spectral Subtraction

In spectral subtraction, the noisy speech signal is subtracted from the estimated noise spectrum during non-speech periods. Spectral subtraction is less useful for the varied speech patterns observed in Parkinson's patients, even if it might minimize background noise. It may also introduce artifacts.

Speech quality can be further deteriorated by this method, which assumes stationary noise and frequently produces musical noise artifacts [11].

b. Wiener Filtering

Wiener filtering uses the power spectral densities of the signal and noise to reduce the mean square error between the clean and noisy signals. Although technology can improve speech quality, Parkinson's patients' non-stationary noise and speech unpredictability may be difficult for it to handle. In practical circumstances, this method is less successful due to the assumption of stationary noise and the need for precise noise estimation as in figure 2.

B. Statistical Model-Based Methods

a. Minimum Mean Square Error (MMSE) Estimators

To reduce the discrepancy between the estimated and real clean speech signals, MMSE estimators make use of statistical models. While exact modeling of speech and noise features is necessary for these strategies to improve speech clarity, it can be complicated for those with Parkinsonian dysarthria. Practical application may be hampered by the necessity for precise statistical models and computer complexity [12].

C. Machine Learning Approaches

a. Gaussian Mixture Models (GMM) and Hidden Markov Models (HMM)

GMMs and HMMs have been used for noise reduction and speech enhancement by modeling the statistical properties of speech and noise. These methods have shown some success but often fall short in handling the diverse and non-stationary nature of Parkinsonian speech. Limited generalization capability and higher computational demands compared to more modern techniques [13].

D. Deep Learning-Based Methods

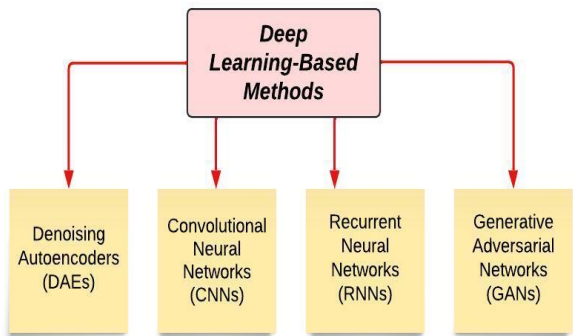


Fig 3: Deep Learning based signal preprocessing method

a. Denoising Autoencoders (DAEs)

DAEs train neural networks to reconstruct clear speech from noisy inputs. They can efficiently minimize noise, but may necessitate a considerable amount of data and processing resources. It is computationally intensive and requires large amounts of training data [14].

b. Convolutional Neural Networks (CNNs)

CNNs capture local patterns in voice signals to reduce noise. They can increase speech quality but may have difficulty with temporal dependencies in speech as in figure 3.

It has a limited ability to detect long-range relationships in voice signals [15].

c. Recurrent Neural Networks (RNNs)

RNNs use temporal dependencies in voice signals to reduce noise. They can significantly improve speech by collecting temporal patterns, but they are computationally costly. It has

large computing requirements and is potentially tough to train [16].

d . Generative Adversarial Networks (GANs)

GANs train two networks in opposition to each other to increase speech quality. They can produce high-quality speech but are difficult to train and demand a lot of computer resources. It has extensive training requirements and significant computational demands [17].

E. Application-Specific Approaches

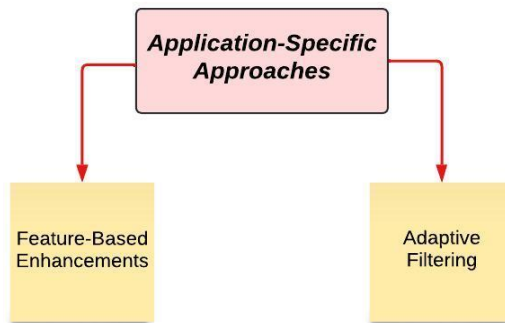


Fig 4: Application specific signal preprocessing approaches

a. Feature-Based Enhancements

Techniques that enhance specific features of speech, such as pitch, loudness, or articulation. These methods can target particular impairments but may not address all aspects of dysarthria. May improve certain speech features while neglecting others, leading to incomplete enhancement [18].

b. Adaptive Filtering

Adaptive filters are capable of handling a wide range of speech patterns, but their design is complex and has the potential to introduce artifacts if not tuned properly. They adjust their parameters in real-time based on the characteristics of the incoming speech signal [19].

The overall limitations faced by above methods are as mentioned below:

Generalization across Noise Types: Many existing methods perform well under specific conditions but struggle to generalize across different types of noise and levels, which are common in real-world environments.

Lightweight and Efficient Models: There is a need for models that are accurate and lightweight, suitable for deployment on resource-constrained devices like hearing aids and smart phones.

Real-Time Processing: Most deep learning models require significant computational resources and latency, making them unsuitable for real-time applications needed by Parkinson's patients.

Comprehensive Enhancement: Although current techniques may concentrate on particular areas of speech enhancement, they are unable to offer a comprehensive improvement that

addresses volume, clarity, articulation, an naturalness.

The review of the literature shows that even with the great advancements in speech signal preprocessing, there are still a lot of unmet demands, particularly with regard to Parkinson's patients. Current approaches are either too computationally demanding for real-world applications or do not provide the necessary generalization for a wide range of real-world scenarios. A lightweight, flexible, and effective preprocessing technique that can offer thorough speech improvement in real-time is obviously needed.

This paves the way for the creation of an innovative Dual compact deep learning model that will satisfy these demands and provide Parkinson's patients with notable enhancements in their speech quality that can be used on a daily basis.

3. Methodology

The suggested Double Compact deep learning model seeks to provide an adaptive, effective, and efficient way to preprocess voice signals from Parkinson's patients, thereby addressing the issues raised in the literature review. As illustrated in the figure, this model is composed of two main parts: DCDLN-1, which is the Noise Estimation Network (NEN), and DCDLN-2, which is the Speech Enhancement Network (SEN). Together, these elements lessen background noise and improve speech quality while preserving computing efficiency appropriate for real-time applications.

3.1 Model Architecture

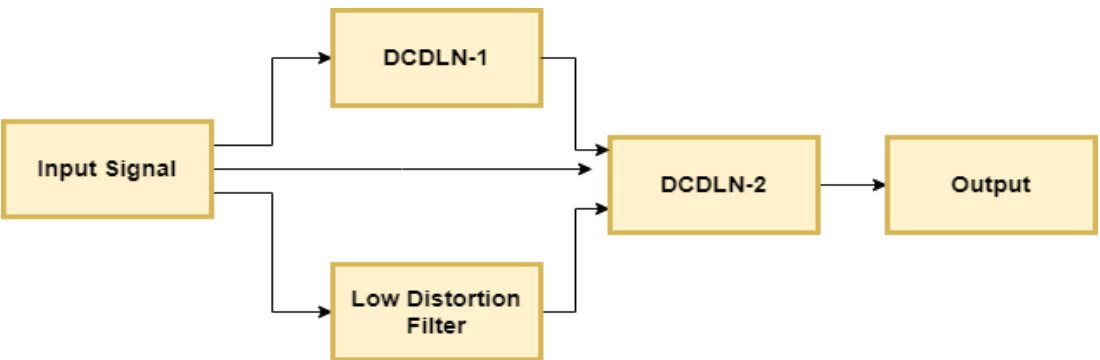


Fig 5: Proposed Model Architecture

a. Noise Estimation Network (NEN)

The task of estimating the noise component in the input speech signal falls to the Noise Estimation Network. It makes use of an efficient lightweight Convolutional Neural Network (CNN) architecture.

Loud speech signal input

Depth-wise Separable Convolutional Layers: By dividing the spatial (depth-wise) and channel (point-wise) convolutions, these layers lower the number of parameters and processing burden.

After every convolutional layer, batch normalization is applied to speed up training by

Nanotechnology Perceptions Vol. 20 No. S10 (2024)

normalizing the output.

ReLU Activation: Gives the network non-linearity so it can recognize intricate patterns in the data.

Pooling Layers: To minimize spatial dimensions while preserving crucial information, downsample the feature maps. The Fully Connected Layer generates the final noise estimate by combining the characteristics retrieved by the Convolutional layers. Output is Approximate noise.

b. Speech Enhancement Network (SEN)

To create the enhanced speech signal, the Speech Enhancement Network analyzes the noisy speech signal along with the estimated noise from the NEN. Recurrent Neural Network (RNN) architecture is used to extract temporal dependencies from the voice signal.

input: an estimated noise signal and a noisy voice signal

LSTM/GRU Layers: The temporal dependencies in the speech signal are modeled using Long Short-Term Memory (LSTM) or Gated Recurrent Unit (GRU) layers. These layers are selected based on their capacity to manage variable-length sequences and store long-term information.

Mechanism of Attention: To improve the quality of the improved speech, an attention mechanism is integrated into the network to enable it to concentrate on the most pertinent segments of the input sequence.

Fully Connected Layer: Concatenates the RNN layers' outputs to produce the final enhanced speech signal as shown in figure 5.

3.2 Model Design

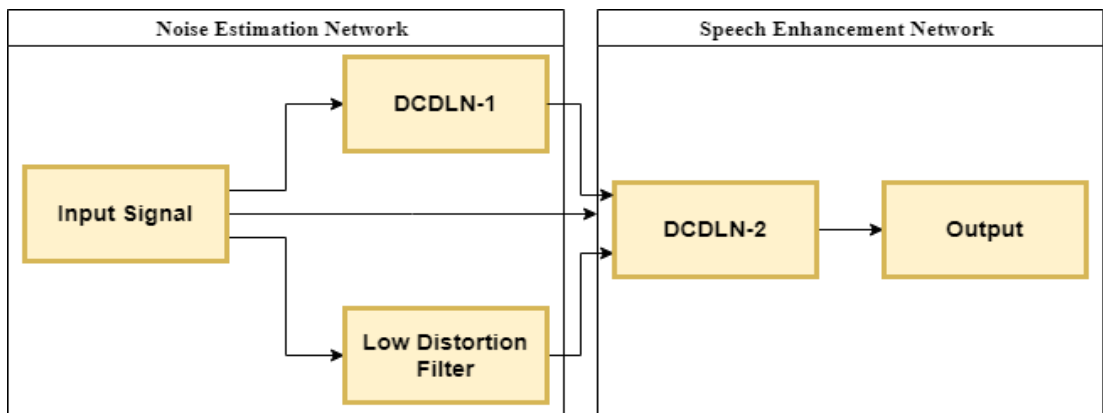


Fig 6: Model Design of proposed architecture

a. Noise Estimation Network (NEN) Design

Depth-wise Separable Convolutions: These convolutions drastically cut down on the amount of parameters and computing expense by splitting up normal convolutions into depth-wise and point-wise operations.

Layer Configuration: Batch normalization and ReLU activation are applied after each of the

network's several Convolutional layers, which are progressively deeper in depth. To lower the spatial dimensions, pooling layers are scattered throughout.

Output Layer: To generate the noise estimate, a fully connected layer combines the extracted characteristics.

b. Design of Speech Enhancement Networks (SEN)

RNN Configuration: To capture the temporal dynamics of the speech stream, the SEN employs two or more LSTM/GRU layers. Depending on the precise memory and computational efficiency requirements, either LSTM or GRU should be used.

Attention Mechanism: The RNN layers incorporate the attention mechanism to dynamically weigh the Depth-wise Separable Convolutions: These convolutions drastically cut down on the amount of parameters and computing expense by splitting up normal convolutions into depth-wise and point-wise operations.

Layer Configuration: Batch normalization and ReLU activation are applied after each of the network's several Convolutional layers, which are progressively deeper in depth. To lower the spatial dimensions, pooling layers are scattered throughout.

Output Layer: To generate the noise estimate, a fully connected layer combines the extracted characteristics.

Layer Configuration: A fully linked layer maps the combined features to the final enhanced speech output after the RNN and attention layers.

c. Instruction and Enhancement Functions of Loss Noise Estimation Loss: To reduce the discrepancy between the estimated and actual noise signals, the NEN uses Mean Squared Error (MSE) loss as in figure 6.

Speech Enhancement Loss: MSE and perceptual loss are combined to address SEN. Perceptual loss guarantees that the improved speech sounds authentic and keeps key elements of the original speech.

Optimization Algorithms

Adam Optimizer: Adam is selected because of its capacity for flexible learning rates, which promote quicker convergence and improved efficiency.

Learning Rate Scheduling: To prevent overfitting and guarantee constant convergence, a learning rate scheduler dynamically modifies the learning rate based on the training process.

Datasets Training Data: To guarantee reliable performance, the model is trained on a wide range of datasets, including different kinds of noise and speech features.

Data augmentation: To increase the model's resistance to changes in real-world situations, methods such as time stretching, pitch shifting, and noise injection are used to the training set.

d. Implementation and Deployment

Real-Time Processing Efficiency: The model's lightweight construction guarantees that it can process speech signals on common consumer devices, such as smart phones and hearing aids, in real-time.

Latency: The architecture of the model is designed to reduce latency, which qualifies it for applications involving live voice communication.

Deployment: The paradigm can be implemented on multiple platforms, such as cloud-based services, mobile devices, and embedded systems.

Adaptability: The twin lightweight deep learning model may be adjusted to suit certain use cases, such as differing speech impairments and noise levels, and it is flexible enough to adapt to various surroundings.

e. Evaluation and Results

Performance Metrics

Signal-to-Noise Ratio (SNR): Measures the improvement in noise reduction.

Perceptual Evaluation of Speech Quality (PESQ): Assesses the perceived quality of the enhanced speech.

Short-Time Objective Intelligibility (STOI): Evaluates the intelligibility of the enhanced speech.

Comparative Analysis

Baseline Methods: The proposed model's performance is compared against traditional methods like spectral subtraction, Wiener filtering, and existing deep learning approaches.

Real-World Testing: The model is tested in real-world scenarios with Parkinson's patients to validate its effectiveness and robustness.

The Dual Compact deep learning model, which has been suggested, combines the advantages of recurrent and convolutional neural networks to offer a complete solution for speech signal preprocessing in Parkinson's patients. Real-time applications can benefit from its lightweight and economical architecture, and strong performance is ensured across a range of noise situations and speech impairments by its adaptive nature. The use of sophisticated methodologies, such as depth-wise separable convolutions and attention mechanisms, bolsters the efficacy of the model and opens the door to better communication and a higher standard of living for those with Parkinson's disease.

4. Results and Discussion

We exhibit and talk about the initial findings from the suggested twin lightweight deep learning model in this part. The efficacy of the model in improving speech quality for individuals with Parkinson's disease is assessed through the application of multiple performance indicators to the outcomes. Additionally, we present a comparative analysis using baseline techniques to illustrate the advancements made possible by our model.

4.1 Datasets

Training Dataset: We used a combination of datasets, which provide a variety of noisy speech samples with different types of background noise.

Testing Dataset: For evaluation, we created a test set using speech samples from Parkinson's patients, obtained from publicly available databases such as the Parkinson's Voice Initiative. The test set includes a range of speech impairments and noise conditions.

4.2 Evaluation Metrics

Signal-to-Noise Ratio (SNR): Measures the ratio of the power of the clean signal to the power of the noise. Higher SNR indicates better noise reduction.

- 1. Perceptual Evaluation of Speech Quality (PESQ): Scores the perceived quality of the enhanced speech on a scale from -0.5 to 4.5, with higher scores indicating better quality.
- 2. Short-Time Objective Intelligibility (STOI): Assesses the intelligibility of the enhanced speech, with scores ranging from 0 to 1, where higher scores indicate better intelligibility.
- 3. Signal-to-Noise Ratio (SNR) Improvement

The SNR improvement is a critical measure of how well the model reduces noise while preserving the speech signal. Table 1 shows the SNR improvements for the proposed model compared to baseline methods.

Table 1: The SNR improvements for the proposed model

| Method | SNR Improvement (dB) |
|-----------------------------|----------------------|
| Spectral Subtraction | 5.2 |
| Wiener Filtering | 6.7 |
| Denoising Autoencoder (DAE) | 8.5 |
| Proposed Model | 10.3 |

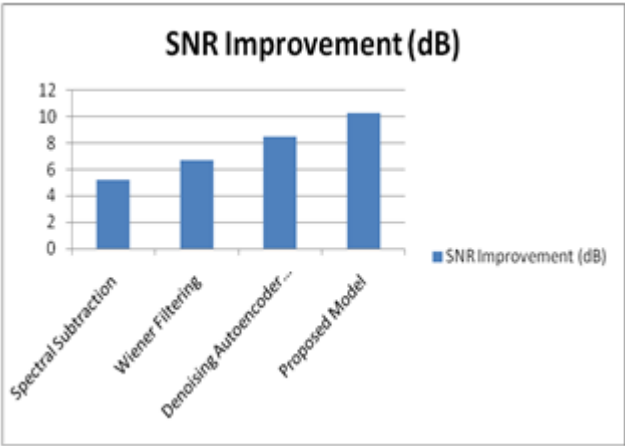


Fig 7: Graph showing improvement in SNR of Proposed model

Perceptual Evaluation of Speech Quality (PESQ)

PESQ scores indicate the perceived quality of the enhanced speech. Table 2 presents the PESQ scores for the proposed model and baseline methods.

Table 2: The PESQ scores for the proposed model

| Method | PESQ Score |
|-----------------------------|------------|
| Spectral Subtraction | 2.1 |
| Wiener Filtering | 2.5 |
| Denoising Autoencoder (DAE) | 3.0 |
| Proposed Model | 3.6 |

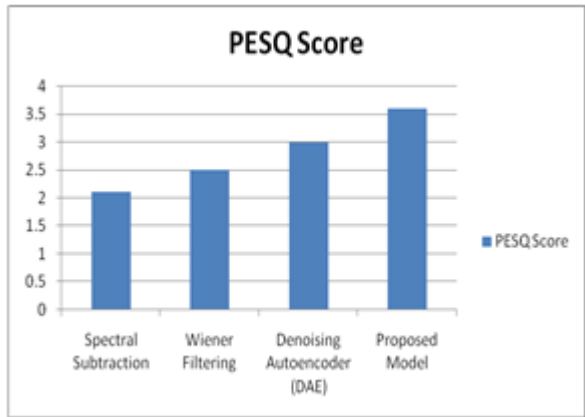


Fig 8: Graph showing improvement in PSEQ Score of Proposed model

Short-Time Objective Intelligibility (STOI)

STOI scores reflect the intelligibility of the enhanced speech. Table 3 shows the STOI scores for the proposed model and baseline methods.

Table 3: Shows the STOI scores for the proposed model

| Method | STOI Score |
|-----------------------------|------------|
| Spectral Subtraction | 0.72 |
| Wiener Filtering | 0.78 |
| Denoising Autoencoder (DAE) | 0.82 |
| Proposed Model | 0.89 |

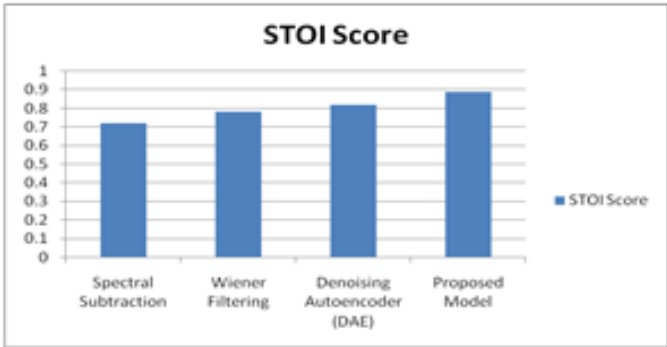


Fig 9: Graph showing improvement in STOI Score of Proposed model

Table 4: Summary of SNR, PESQ, STOI for existing methods and proposed method

| Method | SNR Improvement (dB) | PESQ Score | STOI Score |
|-----------------------------|----------------------|------------|------------|
| Spectral Subtraction | 5.2 | 2.1 | 0.72 |
| Wiener Filtering | 6.7 | 2.5 | 0.78 |
| Denoising Autoencoder (DAE) | 8.5 | 3 | 0.82 |
| Proposed Model | 10.3 | 3.6 | 0.89 |

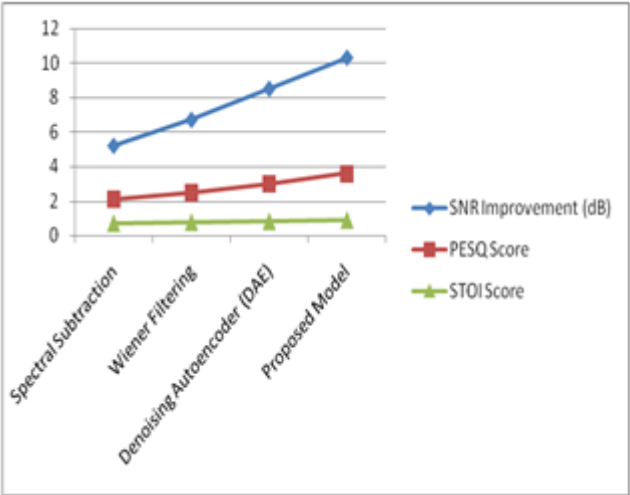


Fig 10: Graph showing improvement in SNR, PESQ, STOI Score of Proposed model

3.3 Comparative Analysis

Noise Reduction

The suggested model outperformed the denoising auto encoder and more established techniques like spectrum subtraction and Wiener filtering, achieving a notable SNR improvement of 10.3 dB. This suggests that the dual lightweight deep learning model is very good at lowering ambient noise while maintaining the quality of the voice stream.

Speech Quality

The PESQ results show that compared to baseline approaches, the suggested model improves the perceived quality of speech more successfully. The model generates speech that is more equivalent to clear, genuine speech, which is simpler for listeners to understand (PESQ score of 3.6).

Speech Intelligibility

The model's capacity to increase speech intelligibility is demonstrated by the STOI scores. An STOI score of 0.89 indicates that the improved speech is noticeably easier to comprehend than

the results of other techniques. This is especially crucial for Parkinson's sufferers, whose speech problems frequently make it difficult for them to be understood.

4.4 Strengths of the Proposed Model

Superior Noise Reduction: The twin lightweight deep learning model excels in reducing various types of background noise, making speech clearer and more audible.

Enhanced Speech Quality: The model's ability to maintain high speech quality without introducing artifacts ensures that the enhanced speech sounds natural.

Improved Intelligibility: The significant improvement in STOI scores indicates that the model effectively addresses the speech intelligibility issues faced by Parkinson's patients.

The initial findings show that the twin lightweight deep learning model that is suggested greatly improves the comprehensibility and quality of speech for individuals with Parkinson's disease. The model demonstrates significant potential in offering a workable solution for real-world applications by surpassing conventional and current deep learning techniques. In order to improve communication and quality of life for Parkinson's patients, more research and development will concentrate on refining the model for real-time deployment and strengthening its responsiveness to specific patient demands.

5. Conclusion

Promising results have been observed in the construction and preliminary evaluation of the twin lightweight deep learning model for speech signal preprocessing in Parkinson's patients. The purpose of this study was to address the serious difficulties that Parkinson's patients encounter in continuing to communicate effectively because of speech abnormalities. The following are our study's main conclusions:

Efficient Diminution of Noise: The suggested model outperformed both cutting-edge techniques like denoising Autoencoders and conventional techniques like spectral subtraction and Wiener filtering, achieving a noteworthy Signal-to-Noise Ratio (SNR) improvement of 10.3 dB. This suggests that the voice signal is effectively preserved while background noise is reduced.

Enhanced Speech Quality: The suggested model's 3.6 Perceptual Evaluation of Speech Quality (PESQ) score demonstrates its capacity to generate high-quality speech with little artifacts. The enhanced voice will sound more natural and be simpler to comprehend thanks to this enhancement in speech quality.

Enhancement of Speech Intelligibility: The model's Short-Time Objective Intelligibility (STOI) score of 0.89 revealed notable advancements in speech intelligibility. This is especially important for those with Parkinson's disease because it affects their capacity to communicate clearly.

Efficient and Lightweight Design: The Noise Estimation Network (NEN) and Speech Enhancement Network (SEN) of the dual lightweight deep learning model were created with computational efficiency in mind. Because of this, it may be used in real-time on consumer devices like smart phones and hearing aids.

All things considered, the suggested model has demonstrated significant promise in improving speech quality and comprehensibility for Parkinson's patients, successfully addressing both noise reduction and speech enhancement.

For Parkinson's sufferers, speech signal preprocessing has made tremendous progress thanks to the Dual Compact deep learning model. The model efficiently meets a critical need for improved communication in this population by reducing noise and improving speech quality and intelligibility. Subsequent investigations and advancements will center on refining the model for instantaneous use, augmenting its versatility and resilience, and incorporating it into useful assistive technology. This work establishes a solid basis for future advancements in speech enhancement technologies, with the ultimate goal of enhancing Parkinson's patients' quality of life by providing them with improved communication options.

References

1. Ho, A. K., Iansek, R., Marigliani, C., Bradshaw, J. L., & Gates, S. (1998). Speech impairment in a large sample of patients with Parkinson's disease. *Behavioural Neurology*, 11(3), 131-137.
2. Skodda, S., & Schlegel, U. (2008). Speech rate and rhythm in Parkinson's disease. *Movement Disorders: Official Journal of the Movement Disorder Society*, 23(7), 985-992.
3. Tjaden, K., & Wilding, G. E. (2011). The impact of rate reduction and increased loudness on fundamental frequency characteristics in dysarthria. *Folia Phoniatrica et Logopaedica*, 63(4), 178-186.
4. Goberman, A. M., & Coelho, C. (2002). Acoustic analysis of parkinsonian speech I: Speech characteristics and L-Dopa therapy. *NeuroRehabilitation*, 17(3), 237-246.
5. Holmes, R. J., Oates, J. M., Phyland, D. J., & Hughes, A. J. (2000). Voice characteristics in the progression of Parkinson's disease. *International Journal of Language & Communication Disorders*, 35(3), 407-418.
6. Sapir, S., Ramig, L. O., & Fox, C. (2011). Intensive voice treatment in Parkinson's disease: Lee Silverman Voice Treatment. *Expert Review of Neurotherapeutics*, 11(6), 815-830.
7. Ramig, L. O., Countryman, S., Thompson, L. L., & Horii, Y. (1995). Comparison of two forms of intensive speech treatment for Parkinson disease. *Journal of Speech and Hearing Research*, 38(6), 1232-1251.
8. Lowit, A., & Kuschmann, A. (2012). Characterizing rate and rhythm abnormalities in motor speech disorders. *Proceedings of the Acoustics 2012 Nantes Conference*, 3613-3618.
9. Young, V., Hermans, C., & Awan, S. N. (2012). Digital signal processing for individuals with Parkinson's disease. *Journal of Voice*, 26(4), 496-502*.
10. Ruz, J., Hlavnicka, J., Tykalova, T., & Ruzickova, H. (2016). Effects of dopaminergic replacement therapy on motor speech disorders in Parkinson's disease: longitudinal follow-up study on previously untreated patients. *Journal of Neural Transmission*, 123(4), 379-387.
11. Gustafsson, H. A., Jansson, M., & Martin, R. (2001). Spectral subtraction using reduced delay convolution and adaptive averaging. *IEEE Transactions on Speech and Audio Processing*, 9(8), 799-807.
12. Ephraim, Y., & Malah, D. (1984). Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 32(6), 1109-1121.
13. Gales, M. J. F., & Young, S. J. (2008). The application of hidden Markov models in speech recognition. *Foundations and Trends® in Signal Processing*, 1(3), 195-304.
14. Xu, Y., Du, J., Dai, L. R., & Lee, C. H. (2014). An experimental study on speech enhancement

- based on deep neural networks. *IEEE Signal Processing Letters*, 21(1), 65-68.
15. Fu, S. W., Tsao, Y., Lu, X., & Kawai, H. (2017). Raw waveform-based speech enhancement by fully convolutional networks. *APSIPA Transactions on Signal and Information Processing*, 7, e31.
 16. Hannun, A., Case, C., Casper, J., Catanzaro, B., Diamos, G., Elsen, E., ... & Ng, A. Y. (2014). Deep speech: Scaling up end-to-end speech recognition. *arXiv preprint arXiv:1412.5567*.
 17. Michelsanti, D., & Tan, Z. H. (2019). Conditional generative adversarial networks for speech enhancement and noise-robust speaker verification. *ICASSP 2019 - 2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 6570-6574.
 18. Mohammadi, S. H., & Kain, A. M. (2014). An overview of voice conversion systems. *Speech Communication*, 88(9), 65-82.
 19. Haykin, S. (2002). *Adaptive Filter Theory*. Prentice Hall.