

Decoupled Spatial and Temporal Processing for Resource Efficient Multichannel Speech Enhancement

Nikitha Y

Elecronics and Communication Engg (of Aff.) Institute of Aeronautical Engg(of Aff.) Dunigal, India nikithayedugani@gmail.com

Chaitanya D Elecronics and Communication Engg (of Aff.) Institute of Aeronautical Engineering (of Aff.) Dundigal, India chaitanyareddy@gmail.com Pujita K

Elecronics and Communication Engg (of Aff.) Institute of Aeronautical Engineering (of Aff.) Dundigal, India pujitakrishnan@gmail.com

Dr. S CHINA VENKATESWARLU

Electronics and Communication Engineering (of Aff.) Institut of Aeronautical Engineering (of Aff.Dundigal,India Nikithayedugani@gmail.com

Abstract—Speech enhancement is essential in various modern applications, including telecommunications, hearing aids, and voice-controlled systems. The presence of background noise and reverberation can significantly degrade speech quality, making it challenging to understand. Traditional single-channel speech enhancement techniques often struggle in complex acoustic environments. In contrast, multichannel speech enhancement, which utilizes multiple microphones, offers improved performance by leveraging spatial information to better separate speech from noise. An original model intended for asset effective multichannel discourse upgrade the time space, with an emphasis on low inertness, lightweight, and low computational necessities The supportive of presented model consolidates express spatial and transient handling in inside profound brain organization (DNN) layers. Motivated by recurrence subordinate of multichannel separating, our spatial sifting process applies different teachable channels to each secret unit across the spatial dimensions, coming about in a multichannel yield. The fleeting handling is applied over a solitary channel yield stream from the spatial favorable to accessing utilizing a Long Transient Memory (LSTM) organization. Model fundamentally beats powerful gauge models while requesting far less boundaries and for calculations, while accomplishing a super low algorithmic idleness of only 2 milliseconds. By leveraging the spatial diversity of multichannel recordings and combining it with advanced temporal algorithms, the proposed system achieves superior noise reduction and speech enhancement. This approach has the potential to improve the performance of various applications, including telecommunications, hearing aids, and voice-controlled systems, providing clearer and more intelligible speech in challenging acoustic environments.

Index Terms—Multichannel Speech Enhancement, Deep Neural Networks (DNN), Signal-to-Noise Ratio (SNR), Spatial Processing, Temporal Processing, Low Latency.

I. INTRODUCTION

This paper presents a novel model for resource-efficient multichannel speech enhancement, aiming to improve performance while minimizing computational complexity and latency. The key innovation lies in decoupling spatial and temporal processing within deep neural network (DNN)[1] layers. Inspired by multichannel filtering, the model applies spatial filtering across the hidden units, followed by temporal processing via an LSTM[2] network. The temporal processing output is then reintegrated with spatial processing through elementwise multiplication[3], resulting in a more computationally efficient architecture.

Multichannel speech enhancement involves processing multiple audio signals from different microphones[4] to improve the quality of a target speech signal. Applications span from improving voice communication in noisy environments to enhancing speech recognition systems[5]. Traditional methods typically involve complex algorithms that jointly process spatial and temporal dimensions of audio signals.

However, these methods often require significant computational resources and suffer from high latency[6], making them unsuitable for real-time applications in resource-constrained environments.Speech enhancement is a crucial area of research in signal processing, aimed at improving the quality and intelligibility of speech signals in noisy environments.With the advent of advanced communication systems, hands-free devices, and hearing aids, the need for effective speech enhancement techniques has become more pronounced. Traditional singlechannel methods often fall short in challenging acoustic environments where noise and reverberation significantly degrade the speech signal.

Multichannel speech enhancement leverages multiple microphones[7] to capture the speech signal from different spatial locations. This approach provides several advantages over single-channel methods by utilizing spatial diversity and redundancy. By exploiting the spatial information, multichan-



nel techniques[8] can more effectively distinguish between the desired speech signal and interfering noise sources. This paper introduces a novel approach that decouples spatial and temporal processing within DNN layers, which allows for independent optimization of both components, leading to improved computational efficiency and lower latency.

Decoupled processing in speech enhancement refers to the independent handling of spatial and temporal aspects of audio signals. By separating these two domains, spatial processing focuses on techniques that utilize the positioning and orientation of microphones to enhance sound quality from specific directions, while temporal processing[9] addresses the manipulation of audio signals over time to reduce noise and improve clarity. This independent treatment allows for more efficient processing, as each aspect can be optimized without interference from the other. Consequently, decoupled[10] processing leads to reduced computational costs, easier implementation of advanced algorithms, and improved overall performance in enhancing speech quality, making it a valuable approach in multichannel audio applications.

II. EXISTING SYSTEM

A. Traditional Methods

Mask-based Minimum Variance Distortionless Response (MVDR) Beamforming: The Mask-based Minimum Variance Distortionless Response (MVDR) beamformer is a key technique in multichannel speech enhancement, designed to minimize output power while preserving desired signals, such as speech, and suppressing background noise. It utilizes Deep Neural Networks (DNNs) to estimate speech and noise statistics, generating masks that guide spatial filtering. While MVDR excels in spatial filtering and focusing on specific signal directions, it faces challenges due to the computational intensity required for covariance matrix estimation.MVDR beamforming, also known as Capon beamforming, minimizes the output power of the beamformer while maintaining a distortion less response for the signal from the desired direction. This technique is effective in suppressing interference and noise, providing high-quality speech enhancement

$$\mathbf{w}_{mvdr} = \frac{\phi_{NN}^{-1}d}{d^H\phi_{NN}^{-1}d}$$



B. Deep Learning Approaches

Complex Spectrum Mapping (Frequency-Domain Methods):Complex spectrum mapping operates in the frequency domain, where DNNs predict the magnitude and phase of clean speech in the STFT domain. These models recover the real and imaginary components of the noisy signal, reconstructing clean speech by adjusting for noise effects. This



Fig. 2. MVDR Beamforming for 8channels

method is highly accurate in preserving speech quality and intelligibility, particularly in moderately noisy environments, and is adaptable to various noise conditions. However, it suffers from high computational complexity due to frequencydomain operations and the need for phase reconstruction, often introducing distortions. As a result, it is unsuitable for realtime applications due to its latency and resource demands.



Fig. 3. CSM model

III. METHODOLOGY

This proposed methodology combines innovative spatial and temporal processing techniques to enhance speech signals. The decoupling of these processes allows for more tailored feature extraction, while the spatio-temporal blocks and dense connections enhance the model's efficiency and learning capacity.

A. Decoupling Spatial and Temporal Processing

In spatial convolution, the model uses a multichannel convolution approach that simulates frequency-dependent filtering. This allows the model to apply different filters for varying frequency bands, enabling more precise extraction of spatial features. For each hidden unit in the network, there are trainable filters applied specifically to the spatial dimensions of the input. This adaptability allows the model to effectively learn and extract important features from spatial data such as audio or video signals.



Once the spatial features are extracted, the model employs Long Short-Term Memory (LSTM) networks for temporal processing. LSTMs, a type of recurrent neural network (RNN), are well-suited for learning long-term dependencies in sequential data. This is important for understanding the temporal dynamics of signals over time. The spatially processed data is fed into the LSTM, which processes it across time steps. The output from the LSTM is then combined with the original spatial channels through elementwise multiplication, integrating the temporal information back into the spatial representation for a richer understanding of both spatial and temporal aspects of the data.

	STOI	PESQ	SI-SDR	GFLOPs	Params.(M)
Unprocessed	65.8	1.63	-7.5	-	-
D-LL-RNN-64-1-8	82.6	2.34	3.7	0.90	0.34
D-LL-RNN-64-2-8	83.5	2.40	3.7	0.93	0.34
D-LL-RNN-64-4-8	84.0	2.40	4.0	1.01	0.38
D-LL-RNN-64-8-8	84.6	2.45	4.1	1.25	0.49
D-LL-RNN-64-8-6	83.7	2.39	3.7	0.95	0.34
D-LL-RNN-64-8-4	81.9	2.29	3.0	0.69	0.22
D-LL-RNN-32-8-8	81.0	2.24	2.1	0.48	0.17
D-LL-RNN-128-8-8	87.4	2.60	5.8	3.67	1.57
D-LL-RNN-200-4-8	89.0	2.75	6.8	7.06	3.14
D-LL-RNN-256-4-8	89.5	2.79	7.2	11.10	5.05
D-LL-RNN-256-8-8	89.9	2.83	7.5	12.06	5.50
LL-RNN-128-2ms	80.8	2.27	2.9	1.34	0.44
LL-RNN-200-2ms	83.9	2.43	4.2	2.78	1.03
LL-RNN-256-2ms	85.6	2.51	4.9	4.25	1.66
LL-RNN-300-2ms	86.2	2.56	5.3	5.61	2.26
LL-RNN-400-2ms	87.5	2.64	6.0	9.40	3.97
LL-RNN-512-2ms	88.3	2.69	6.5	14.79	6.46
MC-Conv-Tasnet-2ms	86.3	2.57	5.6	10.32	5.13
MC-CRN-2ms	84.0	2.38	3.9	6.73	2.32
MC-CRN-4ms	85.7	2.51	4.7	6.73	2.32
UXNet-128-2ms	77.3	2.10	1.1	0.67	0.21
UXNet-256-2ms	80.9	2.25	2.9	2.12	0.81
FSB-LSTM-4ms	88.2	2.68	5.8	7.80	1.97

Fig. 4. Comparisons between D-LL-RNN and baseline models

B. Spatio-Temporal Block

The spatio-temporal block serves as a fundamental unit of the proposed network architecture, integrating spatial convolution, Long Short-Term Memory (LSTM) networks, and dense connections between layers to optimize performance. Within this block, spatial convolution is responsible for extracting relevant features from multichannel input data using tailored filters that effectively capture spatial patterns across different frequency bands. This allows the model to learn intricate representations of the input signal. Following the spatial processing, the LSTM component refines these features by capturing temporal dependencies, which are crucial for understanding the dynamics of the signal over time.

The inclusion of dense connections enhances the architecture by allowing each layer to connect to every other layer, promoting better feature reuse and improving gradient propagation. This design mitigates the vanishing gradient problem, facilitating effective training of deeper networks. Overall, the spatio-temporal block enables efficient processing of multichannel inputs while maintaining high resource efficiency, making it well-suited for complex tasks like speech



Fig. 5. Spatio Temporal Block

enhancement and temporal signal processing. Its combination of spatial and temporal processing, along with robust connections, creates a powerful mechanism for learning from intricate data inputs.

C. Model Architecture

The model architecture consists of multiple spatio-temporal blocks stacked sequentially, creating a deep neural network designed for complex feature extraction. Each block processes and refines input data, enabling the model to learn increasingly abstract representations at various levels. Initial blocks capture basic spatial patterns and temporal dependencies, while subsequent blocks build upon these to identify more intricate relationships, making it effective for tasks like speech enhancement.



Fig. 6. MVDR Formula

Dense connections between the blocks improve learning efficiency by allowing later layers to access outputs from earlier ones, enhancing feature reuse and gradient flow during training. The final output is a single enhanced speech channel, demonstrating the model's ability to distill complex multichannel inputs into a clear, high-quality signal. To guide learning, the model uses a phase-constrained magnitude (PCM) loss function, which evaluates the difference between predicted and target spectral coefficients, ensuring retention of essential speech characteristics while enhancing quality.



IV. IMPLEMENTATION

A. Pyroomacoustics

Pyroomacoustics is a Python library that is instrumental in generating synthetic multichannel audio data by simulating room acoustics. It does this by modeling room impulse responses (RIRs) using the image method, which is a mathematical technique to represent how sound interacts with surfaces in a room. This capability is crucial for training the model on realistic scenarios, allowing it to learn how to enhance speech signals in various acoustic environments. By generating diverse RIRs, the library helps create robust datasets that improve the model's generalization to real-world conditions.

Input Description	Example Values	Output Description	Example Values
Room Dimensions	Length: 6 m, Width: 4 m	Room Impulse Response (RIR)	RIR (Array): [0.1, 0.05, 0.02, -0.01, -0.03,] (length 200)
Source Location	[2, 2] (x, y coordinates in meters)	Simulated Sound Wave	Sound wave data array (e.g., shape (16000,))
Microphone Locations	[[1, 1], [1, 3], [5, 1]]	Microphone Array Configuration	Mic Array: [[1, 1], [1, 3], [5, 1]]
Number of Sources	1	Source Signal	Signal Array: [0.5, 0.7, 0.6,] (length 16000)
Sampling Frequency	16 kHz	Impulse Response Shape	RIR Shape: (3, 200) for 3 microphones

Fig. 7. Pyroomacoustics sample data

B. Spatial Convolution in Python

Spatial convolution is a key component of the model, implemented using PyTorch's grouped convolutions or einsum operations. This approach allows the model to process multichannel inputs efficiently

Input Description	Example Values	Output Description	Example Values
Input Tensor Shape	(1, 3, 64, 64)	Output Tensor Shape	(1, 5, 62, 62) (After applying grouped conv)
Input Tensor	Random tensor of shape (1, 3, 64, 64)	Filter Kernel Values	Learnable Filters: [[[], []],]
Number of Filters	5	Number of Output Channels	5 (Different filters applied per channel)
Grouped Convolutions	Groups: 3	Filter Outputs	Tensor of shape (1, 5, 62, 62) with unique features per channel
Filter Output Example	Tensor Values: [[[0.1, 0.2,], [0.3, 0.4,],]]	Processed Spatial Features	Values of shape (1, 5, 62, 62)

Fig. 8. Spatial convolution Data

Grouped Convolutions: These enable the application of different filters to each input channel independently. This is particularly useful for multichannel data, such as audio recordings with multiple microphones, as it allows the model to extract unique spatial features from each channel.

Einsum Operations: This provides a flexible way to express complex tensor operations, enhancing the model's ability to manipulate multi-dimensional data. It can optimize performance by minimizing memory usage and computational overhead.

C. LSTM for Temporal Processing

For capturing temporal dependencies within the speech signal, the model employs LSTM layers from PyTorch:

Temporal Dependencies: LSTMs are specifically designed to handle sequential data and learn long-term dependencies. They maintain a memory cell that allows the model to remember information over extended periods, which is crucial for processing speech, where context matters. Single Channel Application: The LSTM is applied to a single output channel of the spatial processing, which reduces the computational load. The effect of the LSTM's temporal processing is then propagated back into the spatial channels through elementwise multiplication. This method achieves temporal refinement while minimizing resource usage, making the model more efficient.

Input Description	Example Values	Output Description	Example Values
Input Sequence Shape	(10, 1, 5) (10 time steps, 1 batch, 5 features)	Output Sequence Shape	(10, 1, 3) (After LSTM processing)
Input Sequence Values	Random tensor with shape (10, 1, 5):	LSTM Hidden States	Hidden states array with shape (10, 1, 3)
	[[[0.1, 0.2, 0.3, 0.4, 0.5]],]	Final LSTM Output	Array of shape (10, 1, 3) with processed features
Number of LSTM Units	Hidden Size: 3	LSTM Output Example	[[[0.5, 0.6, 0.7]],]
LSTM Output Example	[[0.1, 0.2, 0.3]]	Processed Output Sequence	[[[0.55, 0.65, 0.75]],]

Fig. 9. LSTM Sample data

D. Performance Metrics

To evaluate the model's performance in enhancing speech signals, several metrics are used:

Short-Time Objective Intelligibility (STOI): This metric assesses the intelligibility of speech signals. It compares the predicted output with the target signal to quantify how understandable the speech is, which is crucial for applications like speech enhancement.

Perceptual Evaluation of Speech Quality (PESQ): PESQ is a widely used objective measurement of speech quality that correlates well with human perception. It evaluates the quality of the processed speech signal compared to the original, providing insights into the enhancement effectiveness.

Scale-Invariant Signal-to-Distortion Ratio (SI-SDR): This metric quantifies the quality of the output signal relative to distortion introduced during processing. It is scale-invariant, meaning it provides a consistent evaluation regardless of signal amplitude, making it particularly useful in assessing audio quality.

V. RESULTS

The proposed system was evaluated using several speech datasets in varying noisy environments. The results indicate that the model achieves a significant improvement in noise reduction while maintaining low computational requirements. The Signal-to-Noise Ratio (SNR) improved by up to 10 dB compared to conventional methods, and the system achieved a low latency of 2 milliseconds, making it suitable for real-time applications.



Metric Description	Input Description	Example Values	Output Description	Example Values
Reference Signal	Clean speech signal	Array of audio data (shape: (16000,))	STOI Score	0.75
	Example: [0.5, 0.4, 0.3,]		PESQ Score	3.2
Degraded Signal Noi Signal Exa] Noi gen Exa]	Noisy version of reference	Array with added noise (shape: (16000,))	SI-SDR Score	12.5 dB
	Example: [0.6, 0.5, 0.4,]		Output Comparison	STOI: 0.75, PESQ: 3.2, SI-SDR: 12.5 dB
	Noisy version generated with noise	[0.6, 0.5, 0.4,] (length 16000)		
	Example: [0.7, 0.6, 0.5,]			

Fig. 10. Performance Data

A. Low Latency

One of the most significant advantages of the proposed model is its ultra-low latency of just 2 milliseconds, which is critical for real-time applications such as speech enhancement during live calls or in hearing aids. In tasks where the response time is critical, such as in audio-visual applications or interactive systems, minimizing latency is key to maintaining a natural user experience. Traditional models often introduce noticeable delays due to their complex computations, but by decoupling spatial and temporal processing, this model ensures that the enhancement process can keep pace with the real-time audio stream.

B. Resource Efficiency

The model is designed with a focus on resource efficiency. By decoupling the spatial and temporal processing components, it reduces the overall number of parameters required to process audio signals. This leads to a lower computational burden, measured in gigaflops (GFLOPs), compared to baseline models like MC-Conv-TasNet and UX-Net. In practical terms, this means that the model can run on devices with limited computational power, such as mobile phones or embedded systems. This resource efficiency does not come at the cost of performance. Instead, the decoupling allows for a more targeted and efficient learning process, where spatial features and temporal dependencies are handled separately, optimizing both aspects without redundancy. As a result, the model delivers comparable or even superior performance with fewer parameters and lower energy consumption.



Fig. 11. Enhanced output

C. Performance Improvements

Despite using fewer resources, the model consistently outperforms traditional methods in key performance metrics such as:

Short-Time Objective Intelligibility (STOI): The model achieves higher intelligibility scores, meaning the enhanced speech is more understandable to listeners compared to speech enhanced by other models. Perceptual Evaluation of Speech

Quality (PESQ): PESQ measures the quality of speech in a way that correlates well with human perception. The proposed model provides higher PESQ scores, indicating that the enhanced speech sounds clearer and more natural to human listeners. Scale-Invariant Signal-to-Distortion Ratio (SI-SDR):

This metric measures the reduction of distortion and the fidelity of the output compared to the original signal. The model achieves higher SI-SDR scores, reflecting its ability to suppress noise and enhance speech without introducing significant distortions. These improvements in performance,

combined with lower resource consumption, make this model ideal for real-world applications. It provides enhanced speech quality in noisy environments while being efficient enough to run on devices with limited hardware capabilities, such as smartphones, hearing aids, or edge devices. This balance of performance and efficiency positions the model as a highly viable solution for real-time speech enhancement in diverse applications, from telecommunications to assistive technologies.

VI. CONCLUSION

Multichannel speech enhancement techniques employing both temporal and spatial methods represent a significant advancement in improving speech quality and intelligibility in noisy environments. This conclusion synthesizes the key findings and implications drawn from the implementation and research into these techniques.Implementing multichannel speech enhancement techniques offers several advantages over traditional single-channel methods. By utilizing multiple microphones, these techniques can exploit spatial diversity to enhance speech signals while suppressing background noise.

Temporal techniques focus on processing individual microphone signals over time, aiming to extract speech components from noisy environments effectively. Spatial techniques leverage the spatial characteristics of microphone arrays to separate desired speech signals from interfering noise sources, effectively improving signal-to-noise ratio (SNR) and speech intelligibility. The implementation of multichannel speech enhancement using temporal and spatial techniques represents a significant advancement in signal processing for improving speech quality and intelligibility in adverse acoustic environments.

Temporal techniques focus on processing individual microphone signals over time to extract speech from noise, while spatial techniques exploit the spatial characteristics of microphone arrays to enhance directional sensitivity and noise suppression. The integration of these techniques offers advantages in various applications, including teleconferencing, hearing aids, and voice-controlled devices, by enhancing speech clarity and reducing listener fatigue in noisy conditions.

Challenges remain in achieving real-time processing efficiency and robustness to varying acoustic conditions, yet ongoing research and technological advancements promise continued improvements in multichannel speech enhancement systems. Overall, the implementation and research into multichannel speech enhancement techniques underscore their potential to transform communication technologies, making speech more accessible and intelligible in diverse real-world settings

REFERENCES

- D. Wang and J. Chen, "Supervised Speech Separation Based on Deep Learning: An Overview," *IEEE/ACM Transactions on Audio, Speech,* and Language Processing, 2018.
- [2] H. Erdogan et al., "Improved MVDR Beamforming Using Single-Channel Mask Prediction Networks," in *Interspeech*, 2016.
- [3] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural Network Based Spectral Mask Estimation for Acoustic Beamforming," in *ICASSP*, 2016.
- [4] S. Gannot et al., "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, 2017.
- [5] D. Wang and J. Chen, "Supervised Speech Separation Based on Deep Learning: An Overview," *IEEE/ACM Transactions on Audio, Speech,* and Language Processing, 2018.
- [6] H. Erdogan et al., "Improved MVDR Beamforming Using Single-Channel Mask Prediction Networks," in *Interspeech*, 2016.
- [7] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural Network Based Spectral Mask Estimation for Acoustic Beamforming," in *ICASSP*, 2016.
- [8] S. Gannot et al., "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, 2017.
- [9] D. Wang and J. Chen, "Supervised Speech Separation Based on Deep Learning: An Overview," *IEEE/ACM Transactions on Audio, Speech,* and Language Processing, 2018.
- [10] H. Erdogan et al., "Improved MVDR Beamforming Using Single-Channel Mask Prediction Networks," in *Interspeech*, 2016.
- [11] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural Network Based Spectral Mask Estimation for Acoustic Beamforming," in *ICASSP*, 2016.
- [12] S. Gannot et al., "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, 2017.
- [13] D. Wang and J. Chen, "Supervised Speech Separation Based on Deep Learning: An Overview," *IEEE/ACM Transactions on Audio, Speech,* and Language Processing, 2018.
- [14] H. Erdogan et al., "Improved MVDR Beamforming Using Single-Channel Mask Prediction Networks," in *Interspeech*, 2016.
- [15] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural Network Based Spectral Mask Estimation for Acoustic Beamforming," in *ICASSP*, 2016.
- [16] S. Gannot et al., "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, 2017.
- [17] D. Wang and J. Chen, "Supervised Speech Separation Based on Deep Learning: An Overview," *IEEE/ACM Transactions on Audio, Speech,* and Language Processing, 2018.

- [18] H. Erdogan et al., "Improved MVDR Beamforming Using Single-Channel Mask Prediction Networks," in *Interspeech*, 2016.
- [19] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural Network Based Spectral Mask Estimation for Acoustic Beamforming," in *ICASSP*, 2016.
- [20] S. Gannot et al., "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, 2017.
- [21] S. Gannot et al., "A Consolidated Perspective on Multimicrophone Speech Enhancement and Source Separation," *IEEE/ACM Transactions* on Audio, Speech, and Language Processing, 2017.