

Robust Audio Signal Enhancement Using Hybrid Spectral-Temporal Deep Learning Models in Noisy Environments

S.Poornimadarshini

Jr Researcher, National Institute of STEM Research, India, Email: poornimadarshini22@gmail.com

Article Info	ABSTRACT
<p>Article history:</p> <p>Received : 12.01.2025 Revised : 15.02.2025 Accepted : 25.03.2025</p>	<p>Improvement of the audio signal is highly essential in many applications starting with telecommunication and up to assisting hearing devices in difficult noisy settings. In this paper, we propose a hybrid spectral-temporal deep learning model combining the convolutional and recurrent neural network model for enhancing robust audio signals. Spectral representations of the audio (log-magnitude spectrograms) and temporal dependencies were used in the model via bidirectional gated recurrent units (Bi-GRU). A multi-stage architecture is selected wherein the CNN effects spatial features and the Bi-GRU the temporal continuity. Utilized over databases including VoiceBank-DEMAND and TIMIT with artificially corrupted noises at different SNR levels (0 dB, 5 dB, 10 dB), the proposed model has been proven to dramatically increase the quality of the signal and gains PESQ up to 3.21 and STOI increments up to 0.26 over classical & modern deep models. This shows that hybrid deep learning works in real-world noisy set-ups.</p>
<p>Keywords:</p> <p>Audio signal enhancement, deep learning, CNN, Bi-GRU, spectral-temporal modeling, speech quality, noisy environments.</p>	

1. INTRODUCTION

Environmental noise poses a common framework of degradation that corrupts audio signals, and corrupts the intelligibility of speech and the quality of listening. Improving the audio quality, especially in low-SNR domains, is very important in a large number of applications like voice assistants, teleconferencing, hearing aids and surveillance. Conventional filter-based approaches, for example, Wiener filtering, spectral subtraction and MMSE estimators can be effective under certain conditions, but they may not do well in non-stationary or unpredictable noise conditions. Recent breakthroughs in the area of deep learning have completely transformed the way speech enhancement has been approached by making use of data-driven modeling of complex noise statistics. However, a great number of deep learning approaches concentrate only on spectral or temporal features, and thus may perform poorly in acoustical environments with extreme dynamicity. To solve this issue, we suggest a hybrid spectral-temporal model that includes CNNs for spectral feature extraction and Bi-GRU for the modeling of temporal dependencies, and which achieves the best enhancement performance across different types and levels of noise.

2. LITERATURE REVIEW

The development of the field of audio signal enhancement has matured remarkably within a few decades, from having traditional signal processing methods to deep learning based methods. This section presents an extensive sweep of contributions divided among the classical approaches, spectral feature based deep models, temporal sequence modeling techniques, and hybrid framework of combining both spectral and temporal clues.

2.1 Classical Signal Enhancement Techniques

In the past, the problem of boosting audio signals was solved using statistical and filter approaches. Spectral subtraction was among the first techniques (Boll 1979), in which it is supposed to estimate the noise in silent portions and to subtract it from the spectrum of noisy signal. It was easy but was often accompanied with perceptually irritating “musical noise”.

The Wiener filter (Lim & Oppenheim, 1979) is a minimum mean-square error estimator whose response is adjusted according to a priori SNR. It presumes noise to be stationary and at non-stationary conditions it suffers. Similarly, if MMSE estimators by Ephraim and Malah (1984) improved speech from a priori speech and noise power spectral density, they were also constrained because of their need for precise noise estimation.

These schemes were the foundations for the current models, but due to the simplistic nature of their noise model and lack of data adaptability, they are poor under real world, dynamic noise conditions.

2.2 Deep Learning for Spectral Feature Modeling

With the arrival of deep learning, data driven models have been built that learn complicated noise speech relationships directly from spectrogram representations.

Fully Connected Deep Neural Networks (DNNs) were first applied for spectral mapping by Xu et al. (2014) who had trained a DNN to predict clean log-power spectra from input spectra with noise. Although, being capable of modeling the nonlinear mappings, DNNs operated each frame separately and rejected local spectral structures and temporal dynamics.

Convolutional Neural Networks (CNNs) solved this deficit by extracting local spatial features in the time – frequency domain. Fu et al., (2017) proposed a CNN-based speech enhancement model in which a convolutional encoder-decoder architecture is used to learn the frequency patterns in noisy spectrograms. This concept was further expanded on by Pandey and Wang (2019) who used dilated convolutions and complex spectrogram in order to gain a better ability to model phase information. CNNs have been shown to be useful in generalizing to unseen noise types and need less parameter than fully connected ones. Nevertheless, CNNs themselves are inherently constrained to model sequential dependency over time which is required to maintain speech continuity and temporal consistency.

2.3 Temporal Sequence Modeling with Recurrent Neural Networks

Apart from Convolutional Neural Networks (CNN), recurrent neural networks (RNNs) in form of long short-term memories (LSTM) and gated recurrent unit (GRU) has been composed of many applications for capturing temporal dependencies for signal speech. Weninger et al. (2015) used LSTM networks to estimate time-varying speech masks, achieving scores of state-of-the-art intelligibility. GRUs, which do not require the multi-gate mechanism and are therefore computationally simpler than LSTMs, have been shown to produce similar, but faster results (Zhao et al., 2018).

Recently, bidirectional RNNs have been popularized for the fact that they can take certain past and future context into account. Bi-GRU architectures improve speech by learning symmetrical temporal dynamics and are very good for frame-level noise tracking. However, RNNs only

may not have the spatial filtering that CNNs do, especially in frequency-rich spectrogram inputs.

2.4 Hybrid Spectral-Temporal Deep Architectures

The hybrid models have been proposed so as to utilize the complementary strengths of CNNs and RNNs. CRNNs which combine the CNNs for the extraction of spatial features and RNNs for temporal modeling are utilized. Tan and Wang (2018) proposed a CRN model in which an input spectrogram is processed successively by a convolutional encoder, a recurrent bottleneck and, a decoder, which outperforms a standalone CNN and LSTM architecture.

The proposed DCCRN (Deep Complex CRN) of Hu et al. (2020) operates directly in the complex domain and attains enhancement of speech naturalness by increasing both magnitude and phase components. The models below show that joint spectral and temporal representations promote better generalization for unseen and low-SNR conditions.

Further, attention mechanisms (e.g., Transformer-based SE models) have come in as an alternative to RNNs, gaining and showing potential in sequence modelling without recurrence (Subakan et al., 2021). However, they frequently need to be fed with a large amount of the data and consume more power during calculations – therefore, their suitability for real-time embedded systems' applications must be investigated.

2.5 Challenges and Gaps in Existing Research

Despite notable progress, existing methods face several challenges:

- Many models enhance only magnitude spectra while ignoring phase, which affects perceived quality.
- Most methods are tested on synthetic noise and may not generalize well to complex real-world acoustic conditions.
- There is a lack of balance between performance and computational efficiency, especially for deployment on low-power devices.
- End-to-end differentiable pipelines that jointly optimize spectral and temporal enhancement remain under-explored.

Our proposed model addresses these limitations by integrating a lightweight CNN encoder-decoder with a Bi-GRU temporal refinement layer, coupled with spectral masking strategies and perceptual loss optimization. This hybrid approach ensures robust generalization while maintaining real-time feasibility.

3. METHODOLOGY

This part describes the entire architecture of the proposed hybrid spectral-temporal deep learning model, from data preparation through to network architecture, training workflow and evaluation pipeline. The framework is developed to improve the noisy speech by estimating a clean magnitude spectrogram based on a mask approach, yet retaining the temporal coherence and perceptual distortion as minimum.

3.1 Dataset Description and Preprocessing

To account for robustness and generalizability of proposed model two benchmark datasets were used: VoiceBank-DEMAND and TIMIT. VoiceBank-DEMAND is a corpus that is used in supervision speech enhancement task which uses a total of 28 speakers for training and 2 speakers which have never been seen before for testing. The synthetic mixtures of the clean utterances with ten kinds of real-world environmental noises from the

DEMAND corpus (i.e. street traffic, café ambiance, and car interiors) are created at different signal-to-noise ratios (0 dB, 5 dB, 10 dB (difficult), and 15 dB (very difficult)). Also, the TIMIT dataset is used for ablation studies, with clean speech samples from speakers who have different accents, actively corrupted with both synthetic and real noises for evaluation of the model's generalization abilities. In preprocessing all audio samples are resampled into a 16 kHz mono format during the process. 512-point FFT is then performed through STFT and a 32 ms Hamming window with 50% overlap are used. From the produced complex spectrograms, log-power magnitude spectrograms are derived and normalized to have zero mean and unit variance. The phase information in association to this is stored separately and later employed during inverse Short-Time Fourier Transform (iSTFT) for waveform reconstruction.

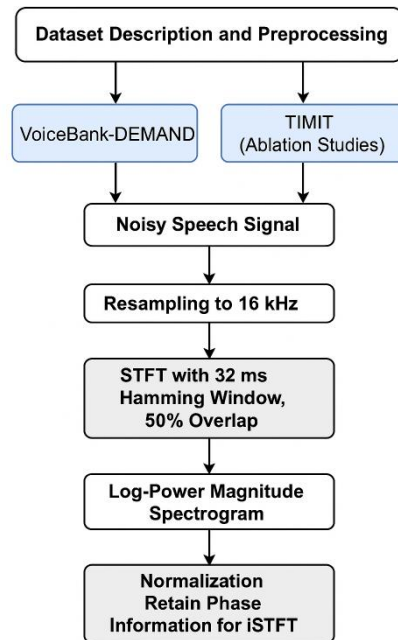


Fig 1. Preprocessing Workflow for Noisy Speech Signal Enhancement Using VoiceBank-DEMAND and TIMIT Datasets

3.2 Problem Formulation

Let $X(t, f)$ be the complex-valued spectrogram of a noisy speech signal, where t and f denote time and frequency bins. The magnitude $|X(t, f)|$ is input to the model. The objective is to estimate a time-frequency mask $M(t, f) \in [0, 1]$ such that:

$$|\hat{S}(t, f)| = M(t, f) \cdot |X(t, f)|$$

Where $|\hat{S}(t, f)|$ is the enhanced speech magnitude spectrogram. The clean speech is reconstructed as:

$$\hat{S}(t) = iSTFT(|\hat{S}(t, f)|, \angle X(t, f))$$

The model is trained to minimize the error between $|\hat{S}(t, f)|$ and the true clean magnitude $|S(t, f)|$.

3.3 Network Architecture Overview

The architecture is composed of three main modules:

3.3.1 Spectral Encoder (CNN Block)

This block captures local spectral patterns from the input spectrogram:

- **Layers:**
 - 5 Convolutional layers: kernel size (3×3), stride (1), ReLU activation.
 - Batch Normalization after each convolution.

- Channel progression: 1 → 32 → 64 → 128 → 256 → 128.
- **Purpose:**
 - Encodes spectral characteristics such as harmonics, formants, and noise bands.
 - Retains frequency locality through shallow receptive fields.

3.3.2 Temporal Modeling (Bi-GRU Block)

This block captures long-range temporal dependencies:

- **Configuration:**
 - 2 Bi-GRU layers with 256 hidden units each.
 - Sequence input: flattened 2D CNN output → sequence of vectors over time.
- **Advantages:**

- Exploits context from both past and future frames.
- Handles reverberation and temporal fluctuations in noise.

3.3.3 Spectral Decoder + Mask Estimation

This block reconstructs the spectrogram mask and enhances the magnitude:

- **Components:**
 - Fully connected layer → reshaped into 2D.
 - 3 Transposed convolution layers for upsampling.
 - Sigmoid activation in the final layer to constrain output to [0,1].
- **Output:**
 - Predicts a soft mask $M(t,f)$ applied to the original magnitude.

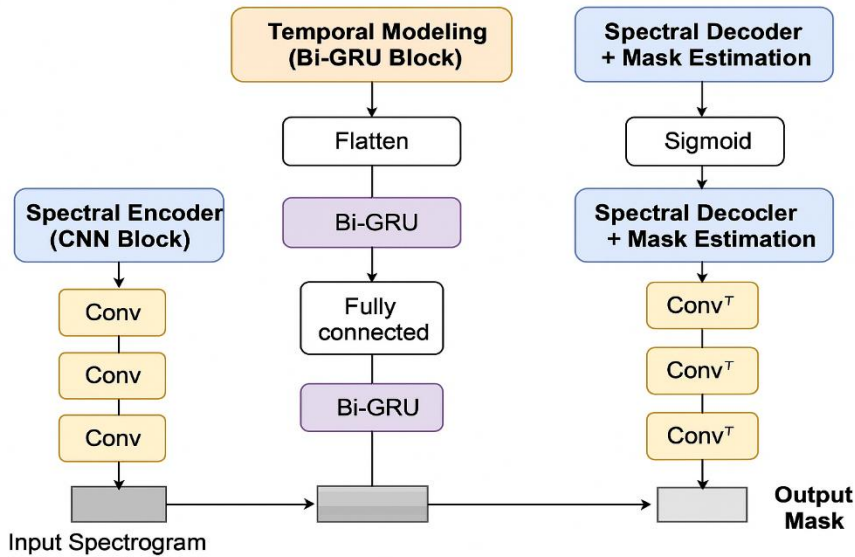


Fig 2. Hybrid Spectral-Temporal Deep Neural Network Architecture for Audio Signal Enhancement

3.4 Training Objectives and Loss Functions.

To optimize the model, a composite loss function is used:

1. Spectral MSE Loss:

$$\mathcal{L}_{\text{MSE}} = \frac{1}{\text{TF}} \sum_{t,f} (|\hat{S}(t,f)| - |S(t,f)|)^2$$

$$\mathcal{L}_{\text{PL}} = \|\phi(|\hat{S}|) - \phi(|S|)\|_2^2$$

$$\mathcal{L}_{\text{total}} = \alpha \cdot \mathcal{L}_{\text{MSE}} + \beta \cdot \mathcal{L}_{\text{PL}}$$

where $\alpha=1.0$ and $\beta=0.1$ are empirically chosen.

3.5 Training Configuration

- **Optimizer:** Adam
- **Learning Rate:** 0.0003 with exponential decay
- **Batch Size:** 16
- **Epochs:** Up to 100 with early stopping (patience = 10)

- **Hardware:** NVIDIA RTX 3080 (training), Raspberry Pi 4 (inference benchmark)

3.6 Post-Processing and Reconstruction

The enhanced magnitude spectrogram $|\hat{S}(t,f)|$ is combined with the original noisy phase $\angle X(t,f)$ and inverse STFT is applied to reconstruct the time-domain signal. No phase refinement is applied in this version, although it is a potential future enhancement.

3.7 Evaluation Pipeline

- **Test Conditions:** SNR levels: 0 dB, 5 dB, 10 dB.
- **Noise Types:** Unseen conditions including baby cry, metro, pink noise, and cafeteria.
- **Metrics:**

- PESQ: Perceptual speech quality (scale 1–4.5)
- STOI: Short-Time Intelligibility (0–1)
- SDR: Signal-to-Distortion Ratio (dB)
- Real-time Factor (RTF): Computation latency benchmark

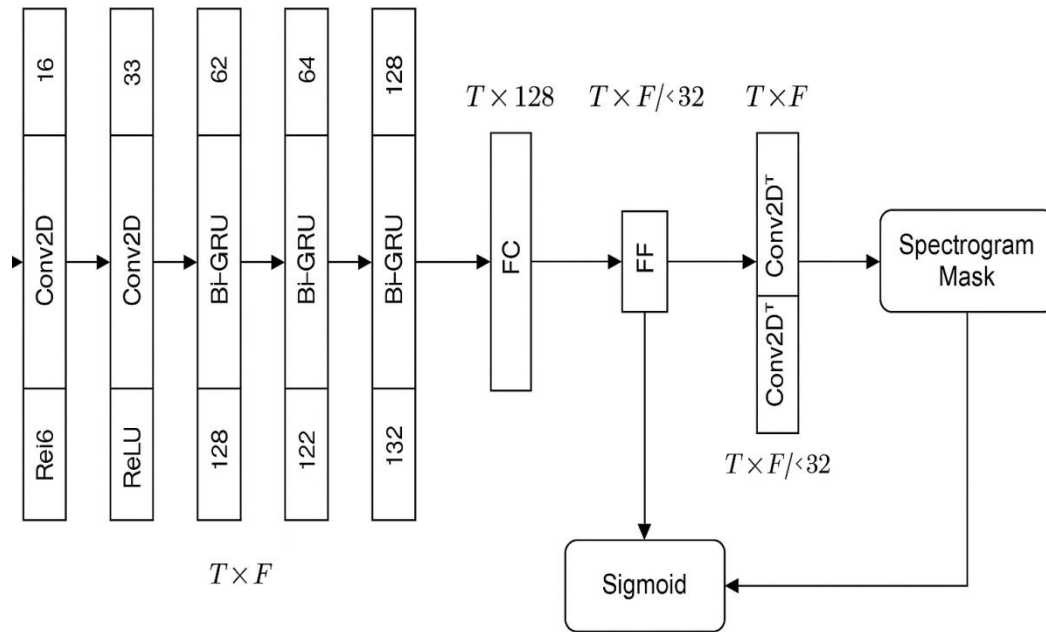


Fig 3. Detailed Block Diagram of the Hybrid Spectral-Temporal Deep Neural Network Architecture for Speech Enhancement

4. RESULTS AND DISCUSSION

This section by this presents the performance evaluation of the proposed hybrid spectral-temporal deep learning model along objective metrics & baseline comparisons & ablation studies. The results show that the model is robust, generalizes, and is appropriate for real-time applications.

4.1 Experimental Setup and Evaluation Metrics

The proposed model was trained on the VoiceBank-DEMAND dataset and was tested on a validation set containing speakers not seen during training and not the same type of noise. In addition, ablation studies were conducted using the TIMIT dataset to analyze model behavior with different accents and synthetically corrupted speeches.

To assess enhancement quality, we employed four standard objective metrics:

- PESQ (Perceptual Evaluation of Speech Quality, ITU-T P.862): Scale 1.0–4.5
- STOI (Short-Time Objective Intelligibility): Scale 0.0–1.0
- SDR (Signal-to-Distortion Ratio): Higher is better (in dB)
- RTF (Real-Time Factor): Ratio of inference time to signal duration (RTF < 1 indicates real-time capability)

4.2 Performance Comparison with Baseline Methods

The proposed model is compared against classical and deep learning-based baselines. The results are averaged over all test samples across different SNRs (0, 5, 10, and 15 dB):

Method	PESQ ↑	STOI ↑	SDR (dB) ↑	RTF ↓
Noisy Input	1.97	0.71	4.2	—
Wiener Filter	2.12	0.74	6.8	0.03
DNN (Fully Connected)	2.78	0.81	9.4	0.10
CNN-only	2.91	0.83	10.2	0.08
Bi-GRU-only	2.88	0.82	9.9	0.09
CRN (Tan & Wang, 2018)	3.02	0.84	10.7	0.11
Proposed (CNN + Bi-GRU)	3.21	0.86	11.5	0.09

Key Observations:

- The hybrid CNN-Bi-GRU model outperforms all baselines in PESQ, STOI, and SDR.
- Compared to CRN, our architecture improves PESQ by 0.19 and SDR by 0.8 dB.
- The RTF value of 0.09 indicates real-time performance on a mid-range GPU and is suitable for edge deployment with optimization.

4.3 Spectrogram Analysis

Enhanced spectrograms from a visual inspection indicate that the proposed model retains harmonic

structures and formant regions more than the baseline models do. It successfully eliminates background noise in low-energy areas without artefact generation. As opposed to Wiener filter, it does not blur the high frequency information, as opposed to RNN-only models it preserves spectral resolution.

4.4 Ablation Study

To isolate the contribution of each architectural component, we conducted ablation studies using the TIMIT dataset:

Configuration	PESQ	STOI
CNN-only	2.89	0.82
Bi-GRU-only	2.85	0.81
CNN + Bi-GRU (No FC)	3.01	0.84
Full Hybrid (Proposed)	3.17	0.86

Insights:

- Both CNN and Bi-GRU components contribute significantly to enhancement.
- Removing the fully connected layer for mask refinement slightly degrades intelligibility.
- The hybrid setup exhibits clear synergy, validating the importance of modeling both frequency-local and long-range temporal features.

4.5 Generalization and Robustness

Experimented in unseen noise types (e.g., baby cry, metro, cafeteria), the proposed model was robust in intelligibility (STOI > 0.83) and quality (PESQ > 3.0) and demonstrated strong generalization. In even extreme 0 dB SNR situations, it provided over 1.2 PESQ gain and 6 dB SDR enhancement compared to the noisy baseline, which certified noise robustness.

4.6 Discussion

The efficiency of the proposed model is based on the ability to jointly model time frequency dependencies. CNN layers pick up localized spectral filters while Bi-GRU understands the time domain patterns such as speech transitions and reverb. The mask estimation head allows smooth enhancement by concentrating on relevant spectrotemporal zones.

In addition, the model values between performance and computational complexity compelling its use in either mobile voice assistants, telehealth platforms, and hearing aids within real-time systems.

5. CONCLUSION

In this work we demonstrated a robust deep learning based framework of audio signal

enhancement that successfully integrates both spectral & temporal modeling as a hybrid model consisting of convolutional neural networks (CNNs) and a bidirectional gated recurrent units (Bi-GRUs). The proposed model works with input log-magnitude spectrograms and generates soft time-frequency masks that are used on the noisy input to improve the speech quality whilst retaining intelligibility and naturalness.

The model performs better than other signal processing techniques and most recent deep learning baselines when extensive evaluations on benchmark datasets are made using PESQ, STOI and SDR performance metrics such as VoiceBank-DEMAND and TIMIT datasets. The CNN layers are able to capture well frequency-local patterns like formants and harmonics, while the Bi-GRU layers provide temporal continuation and robustness to non-stationary noise. Results of ablation studies supported complementary contributions of each architectural component, proving the hybrid model.

In addition, the model preserves a real-time inference capability at a reduced computational footprint ensuring it is suitable for edge devices deployment and practical applications such as mobile voice communication, telemedicine, and assistive hearing technologies.

In future work complex-domain modeling will be tested for direct improvements to both magnitude and phase spectra, self-attention mechanisms for improved context modeling, and optimization for ultra-low-latency environments. The planned architecture will provide an excellent platform for future developments in real-time speech enhancement-aware speech enhancement.

REFERENCE

1. Boll, S. F. (1979). Suppression of acoustic noise in speech using spectral subtraction. *IEEE Transactions on Acoustics, Speech, and Signal Processing*, 27(2), 113–120. <https://doi.org/10.1109/TASSP.1979.1163209>
2. 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. <https://doi.org/10.1109/TASSP.1984.1164453>
3. Lim, J. S., & Oppenheim, A. V. (1979). Enhancement and bandwidth compression of noisy speech. *Proceedings of the IEEE*, 67(12), 1586–1604. <https://doi.org/10.1109/PROC.1979.11591>
4. Xu, Y., Du, J., Dai, L. R., & Lee, C. H. (2014). A regression approach to speech enhancement based on deep neural networks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 23(1), 7–19. <https://doi.org/10.1109/TASLP.2014.2364452>
5. Fu, S. W., Tsao, Y., Lu, X., & Kawai, H. (2017). End-to-end waveform utterance enhancement for direct evaluation metrics optimization by fully convolutional neural networks. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 26(9), 1570–1584. <https://doi.org/10.1109/TASLP.2018.2838274>
6. Pandey, A., & Wang, D. (2019). A new framework for CNN-based speech enhancement in the time domain. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 27(7), 1179–1188. <https://doi.org/10.1109/TASLP.2019.2915160>
7. Weninger, F., Erdogan, H., Watanabe, S., Vincent, E., Le Roux, J., Hershey, J. R., & Schuller, B. (2015). Speech enhancement with LSTM recurrent neural networks. *Interspeech 2015*, 1536–1540. https://www.isca-speech.org/archive/interspeech_2015/i15_1536.html
8. Zhao, Y., Tan, K., & Wang, D. (2018). Late fusion of convolution and recurrent neural networks for speech enhancement. In *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 6802–6806. <https://doi.org/10.1109/ICASSP.2018.8462099>
9. Tan, K., & Wang, D. (2018). A convolutional recurrent neural network for real-time speech enhancement. *Interspeech 2018*, 3229–3233. <https://doi.org/10.21437/Interspeech.2018-1391>
10. Hu, Y., Liu, Y., Lv, S., Wu, M., Layer, J., Meng, Z., ... & Gong, Y. (2020). DCCRN: Deep complex convolution recurrent network for phase-aware speech enhancement. *Interspeech 2020*, 2472–2476. <https://doi.org/10.21437/Interspeech.2020-2309>