

# Circuits, Systems, and Signal Processing

## Correntropy based Multi-objective Multi-channel Speech Enhancement

--Manuscript Draft--

<b>Manuscript Number:</b>	CSSP-D-21-00332R2
<b>Full Title:</b>	Correntropy based Multi-objective Multi-channel Speech Enhancement
<b>Article Type:</b>	Original Research
<b>Keywords:</b>	Speech enhancement, microphone array, correntropy, neural networks, masking
<b>Abstract:</b>	<p>Although deep learning-based methods have greatly advanced the speech enhancement, their performance are intensively degraded under the non-Gaussian noises. To combat the problem, a correntropy based multi-objective multi-channel speech enhancement method is proposed. First, the log-power spectra (LPS) of multichannel noisy speech are feed to the bidirectional long short-term memory (BiLSTM) network with the aim of predicting the intermediate log ideal ratio mask (LIRM) and LPS of clean speech in each channel. Then, the intermediate LPS and LIRM features obtained from each channel are separately integrated into a single-channel LPS and a single-channel LIRM by fusion layers. Next, the two single-channel features are further fused into a single-channel LPS and finally feed to the deep neural network to predict the LPS of clean speech. During training, a multi-loss function is constructed based on correntropy with the aim of reducing the impact of outliers and improving the performance of overall network. Experimental results show that the proposed method achieves significant improvements in suppressing non-Gaussian noises and reverberations, and has good robustness to different noises, signal-noise-ratios (SNRs) and source-array distances.</p>
<b>Response to Reviewers:</b>	<p>The Point-by-point Responses to the Reviewer's Comments</p> <p>We would like to express our appreciation to the reviewers for providing us with valuable comments for improving this manuscript. In the following, we present the point-by-point replies to the reviewers' comments.</p> <p>Reviewer1:</p> <p>Comments No.1: Still, in Abstract, line number 24-29 is not clear, may have to reframe it.</p> <p>Answer: Thanks for the reviewer's advice. We have reframed the sentences in line number 24-29. Please refer to the Abstract of the revised manuscript.</p> <p>Comments No.2: Typo, page-7, line 52, lager--&gt;larger</p> <p>Answer: Sorry for our carelessness. The "lager" has been modified as "larger" with red font in Page 7 in the revised manuscript.</p> <p>Comments No.3: The authors should be consistent with their notations, either IRM or LIRM can be used, If LIRM is used have to change in fig.1 and fig 3 accordingly.</p> <p>Answer: Thanks for the reviewer's valuable advice. We have modified the notations to make them consistent, and Fig. 1 and Fig. 3 have also been changed accordingly.</p> <p>Comments No.4: Still, the computation of LIRM in the intermediate hidden layer is not clear. Are the authors passing both the noisy and clean LPS signal to the intermediate layer to compute the ratio? (clearly mention the same in the respective section)</p> <p>Answer: Sorry for our unclear description. As shown in Fig. 3, LIRM feature is one of the intermediate outputs, which is estimated or computed through the BiLSTM and DNN modules by supervised learning. Both the noisy and clean LPS signals are not used during the computation of LIRM in the intermediate layer. In fact, the ultimate goal of the proposed method is to obtain the LPS feature of the clean speech, thus, according</p>

to Eq. (4), the noisy LPS signal is passed and added with the corresponding LIRM to obtain the indirectly estimated LPS of clean speech in each channel. And the multi-channel indirectly estimated LPS are then fused into one channel and incorporated with directly estimated LPS to form a single-channel LPS feature for further estimating the final LPS of clean speech.

According to the reviewer's advice, we have added the corresponding description in the revised manuscript. Please refer to the first paragraph in Section 2.6.

Comments No.5: In the experiment section, the author showed MRCAE outperforms the proposed approach in low SNR scenarios in the case of untrained noise and dataset conditions. The author should justify it appropriately, and should also mention which attribute of MRCAE helps to perform better in low SNR scenarios. The author can include this in the conclusion section and can provide appropriate motivation for future work.

Answer:

Thanks for the reviewer's comment. Indeed, as shown in Tables 9 and 11, MRCAE outperforms the proposed method in low SNR scenarios. This may because MRCAE is a waveform-based learning method, while the proposed method is a feature-based learning method. Generally, the waveform-based methods process speech in time domain, which causes less damage to frequency bins and could retain more information, including useful clean speech and residual noise. Thus, when the input SNR is low, these retain information make the utterances more smoothly. But when the input SNR is high, the excessive residual noise will significantly degrade the performance, especially for speech intelligibility, as shown in Figs.10 and 11. In contrast, the feature-based learning methods are easy to remove the useful frequency domain information of clean speech during denoising and dereverberation, especially under strong noisy conditions, and lead to speech distortion in time domain.

Meanwhile, new experiments show that the latest work CP-NBDF has the best performance compared with other methods in low SNR scenarios. A possible reason is that CP-NBDF uncouples the inter-frequency dependency from full-band signal, and processes each frequency bin separately to make full use of the narrow-band information. However, MRCAE and the proposed method are all channel-wise methods, which processes the signal in full-band and could not focus on all the information from every frequency bin, and then lead to inferior performance compared to CP-NBDF.

According to the reviewer's advice, the above description about which attribute of MRCAE helps it perform well and the reason why CP-NBDF performs the best in low SNR scenarios have been added in the revised manuscript, respectively. Please refer to the last two paragraphs in Subsection 3.3.4. Besides, they are also included in Conclusion section to provide the motivation for future work.

Reviewer2:

Thanks for your previous comments that are very useful for improving our manuscript.

Reviewer3:

Comments No.1: The authors have now compared with MRCAE based method, which was published in 2018. Although its good to compare with that work, but the reviewer expected comparison with more recent work, at least those that are published after 2020 to have a good comparison.

Following are some, which the authors could have tried:

- (i) Panagiotis Tzirakis et al., "Multi-Channel Speech Enhancement using Graph Neural Networks", ICASSP 2021.
- (ii) Papers under Multi-Channel Speech Enhancement session of Interspeech 2020 and Interspeech 2021 for comparison.

[https://www.isca-speech.org/archive/interspeech\\_2020/index.html#Multi-Channel%20Speech%20Enhancement](https://www.isca-speech.org/archive/interspeech_2020/index.html#Multi-Channel%20Speech%20Enhancement)

Answer:

According to the reviewer's suggestion, a more recent multi-channel method named channel-padding narrow band deep filtering (CP-NBDF) [A1] ([46] in revised manuscript) is employed for comparison, and the results have been added in all the experiments. Please refer to Tables 6-12 and Figs. 5-12 and the corresponding descriptions of Subsections 3.2 and 3.3 in the revised manuscript.

[A1] S. Zhang, X. Li. Microphone array generalization for multichannel narrowband deep speech enhancement, in: Interspeech, pp. 666-670 (2021)

Comments No.2: Although performance measure such as PESQ and STOI are used for speech enhancement, in the recent years the trend is more towards how speech enhancement can help some application. For instance, the current DNS challenge (<https://www.microsoft.com/en-us/research/academic-program/deep-noise-suppression-challenge-icassp-2022/>) focuses on a combined metric that is derived using DNSMOS P.835 (<https://arxiv.org/pdf/2110.01763.pdf>) and ASR word accuracy (as often it is found enhanced speech degrades the speech accuracy). So, I recommend authors to use such things for future works, which will make the works not only interesting, but also useful for the applications.

Answer:

Thanks for the reviewer's comment. DNSMOS P.835 is a perceptual objective metric that serves as a proxy for subjective scores. However, when we tried to connect to DNSMOS API provided by Microsoft, we found that the process speed is very slow and the connection often interrupts. Thus, considering the page and time limitations, the DNSMOS metric is only conducted on untrained noise types experiment. Please refer to Table 10 and the last paragraph of Subsection 3.3.4 in the revised manuscript. Moreover, ASR word accuracy experiment has also been added to further evaluate the proposed method. Please refer to Section 3.4 in the revised manuscript.

# Correntropy based Multi-objective Multi-channel Speech Enhancement

Xingyue Cui<sup>1</sup>, Zhe Chen<sup>1,\*</sup>, Fuliang Yin<sup>1</sup>, Xianfa Xu<sup>1</sup>

*School of Information and Communication Engineering, Dalian University of Technology, Dalian 116023, China*

---

## Abstract

Although deep learning-based methods have greatly advanced the speech enhancement, their performance are intensively degraded under the non-Gaussian noises. To combat the problem, a correntropy based multi-objective multi-channel speech enhancement method is proposed. First, the log-power spectra (LPS) of multichannel noisy speech are feed to the bidirectional long short-term memory (BiLSTM) network with the aim of predicting the intermediate log ideal ratio mask (LIRM) and LPS of clean speech in each channel. **Then, the intermediate LPS and LIRM features obtained from each channel are separately integrated into a single-channel LPS and a single-channel LIRM by fusion layers. Next, the two single-channel features are further fused into a single-channel LPS and finally feed to the deep neural network to predict the LPS of clean speech. During training, a multi-loss function is constructed based on correntropy with the aim of reducing the impact of outliers and improving the performance of overall network.** Experimental results show that the proposed method achieves significant improvements in suppressing non-Gaussian noises and reverberations, and has good robustness to different noises, signal-noise-ratios (SNRs) and source-array distances.

**Keywords:** Speech enhancement, microphone array, correntropy, neural networks, masking

---

## 1. Introduction

When acquiring speech signal in the real-world environment, it is inevitably suffering from noise and reverberation, and such nonideal signal does have a significant impact on the performance of following speech-related tasks. To improve the performance of speech processing tasks, such as speech recognition accuracy [1] and communication quality [2], some speech enhancement (SE) algorithms [3] were proposed to extract the desired speech from noisy-reverberation speech. Generally, speech enhancement methods can be divided into two categories. The first category uses a single microphone (also called monaural), while the second category uses multiple microphones (also called multichannel) to perform speech enhancement.

---

\*Corresponding author: Zhe Chen

Email addresses: xiechoah@mail.dlut.edu.cn (Xingyue Cui), zhechen@dlut.edu.cn (Zhe Chen), flyin@dlut.edu.cn (Fuliang Yin), xuxianfa@mail.dlut.edu.cn (Xianfa Xu)

For monaural based speech enhancement, many approaches have been developed in the past decades. Spectral subtraction [4], Wiener filtering [5], subspace method [6] and statistical model algorithm [7] are all typical noise reduction methods. Weighted prediction error (WPE) [8] and inverse filter [9] are both widely used methods for dereverberation. However, these methods can only achieve either denoising or dereverberation, and the former is less robust to non-stationary noise, while the latter is difficult to fully eliminate the late reverberation. More recently, the supervised learning methods have attracted wide attention since they can adapt to different acoustic conditions. In [10], a single deep neural network (DNN) was employed to simultaneously perform denoising and dereverberation via spectral mapping. In [11], a time-frequency masking called complex ideal ratio mask (cIRM) is proposed by Williamson and Wang for DNN-based enhancement. In the following [12], a two-stage system was proposed, where two DNN-based subsystems are sequentially conducted to denoising and dereverberation, then these pre-trained DNNs were combined into a deeper network for joint training. Subsequently, Li et al. [13] proposed a multi-objective speech enhancement learning framework based on stacked and temporal convolutional neural network, which uses the log-power spectra (LPS), power function compression Mel-frequency cepstral coefficient and ideal ratio mask (IRM) as the target features. The above-mentioned methods are all based on single-channel time-frequency (T-F) information and have achieved good effects in speech enhancement.

Nowadays, with multiple microphones being equipped on modern devices, microphone array techniques have become attractive solutions for speech enhancement. Spatial filtering (or so-called beamforming) is one of the most typical techniques [14]. Its idea is to design a linear filter to enhance or maintain the signal from the target direction while attenuate the interferences from other ones. Another effective method is based on a coherence algorithm that calculates the correlation of two input signals to estimate a filter to weaken the interference components [15], [16]. Besides, post-filtering [17] is commonly used for further noise reduction. Usually, when speech covariance matrices and direction of the arrival can be accurately estimated, the above methods are capable of achieving good speech enhancement performance. However, due to the complex environmental conditions, especially non-stationary noise and room reverberation, the capability of these methods are fundamentally limited. Meanwhile, the number of microphones in the array is also one of the restrictions.

More Recently, deep learning methods have exhibited encouraging performance in multichannel SE tasks. Its main idea is to estimate the time-frequency masks [18] using multi-channel information. For binaural enhancement, Jiang et al. [19] used DNNs to estimate the ideal binary mask (IBM), where monaural feature, binaural features ( interaural time (ITD) and level differences (ILD) are used for network training. In subsequent studies [20, 21], multiple features, such as spatial feature, spectral feature and phase differences between channels, were exploited as the network input to train a DNN or deep auto-encoder (DAE) for speech enhancement. Following in [22], the binaural signals were combined into a monaural complex signal, a complex mask was estimated using the complex DNN and then applied to the monaural complex signal.

In addition, some other methods [23-29] that are not based on the position of speech source have been proposed as well. In [23], a bi-directional long short-term memory (BiLSTM) network was adopted to compute the spectral masks of speech and noise, which then utilized to calculate the corresponding cross-power spectral density (PSD) matrices for getting the beamformer coeffi-

cients. Following, an extended version [24] was developed, where the spatial covariance matrices and beamforming are estimated by complex-valued short-time Fourier transform coefficients, and the magnitude features were used for mask prediction. In recent contributions [25, 26], Chakrabarty proposed a speech enhancement approach based on convolutional recurrent neural network (CRNN), where the estimated T-F masking is either directly applied to the microphone signal or indirect compute the PSD matrices for beamformer. In [27], multiple features of each channel were estimated by using a sharing BiLSTM network, which then fused to obtain the desired single-channel signal. Moreover, in addition to masking- or feature mapping- based methods, some other approaches have also been proposed. In [28], Higuchi et al. derived a frame-by-frame update rule for mask-based minimum variance distortion less response (MVDR) beamformer, which is capable of obtaining the enhanced signals without a long delay. Following, Qi et al. [29] proposed a tensor-to-vector regression approach, which casts the conventional DNN based vector-to-vector regression formulation under a tensor-train network (TTN) framework to address the issue of input size explosion and hidden-layer size expansion.

Although existing methods show strength in suppressing many kinds of noises, they still face with challenges due to non-Gaussian noises, especially impulse noise. The reason lies that most methods are based on mean square error (MSE) criterion which would be fragile to outliers. Thus, another loss function inspired by the statistical measure called correntropy [30, 31, 32] is proposed to improve the robustness of non-Gaussian noises. Correntropy is a nonlinear and local similarity measure that shows the similarity between two random variables in a neighborhood of the joint space controlled by the kernel bandwidth. Compared with MSE, the main advantage of correntropy is its insensitivity to outliers (or impulsive noises), which indicates that it has more potential for robust feature learning. Some related works have been done in recent years. Singh et al. [33, 34] proposed a correntropy based loss function for training network classifier, where the correntropy-loss (C-loss) displays superior robustness against the outlier and can approximate different norms (from L0 to L2) of data. Ref. [35] proposed a robust stacked autoencoder (RSAE) based on maximum correntropy criterion (MCC) to deal with the data containing non-Gaussian noises and outliers. Following, Chen et al. [36] proposed a model based on stacked auto-encoders (SAE) and correntropy-induced loss function (CLF). In this model, the reconstruction loss, the sparsity penalty term and the fine-tuning procedure were all built with CLF, and the results showed an obvious improvement of the model robustness when the data contains outliers and impulsive noise.

The motivation of this work are summarized in Table 1. In order to suppress impulse noise, we extend our previous work [27] to propose a correntropy based multi-objective multi-channel speech enhancement method in this paper. Specifically, a BiLSTM network is first trained to learn the mapping relationship between the multi-channel LPS of noisy speech and their corresponding clean LPS and log-ideal ratio mask (LIRM). Then, the intermediate outputs LPS and LIRM from multiple channels are separately fused into a single-channel LPS and a single-channel LIRM. Finally, the two single-channel features are incorporated into a single-channel LPS and then feed to the fully connected layers for further predicting the LPS of clean speech. Different from conventional speech enhancement network that uses MSE as loss function, the objective function in the proposed method is built based on correntropy. Some experiments under different acoustic conditions verify the superior enhanced performance and the generalization of the proposed method.

**Table 1** The motivation of correntropy based deep learning speech enhancement

Num	Motivation
1	Compared with MSE, correntropy is insensitive to the outliers (or impulse noise).
2	Compared with traditional multi-channel speech enhancement methods, deep learning based methods do not explicitly require the position of speech source but still exhibit encouraging enhanced capability.
3	The performance of the existing deep learning based methods are severely degraded under the non-Gaussian noises, especially impulse noise.

The main contributions of this paper are summarized as follows: 1) A correntropy based multi-objective multi-channel speech enhancement model is proposed, i.e. correntropy is employed to construct a multi-objective loss function to optimize the model; 2) The model structure is designed more reasonable, and displays a superior enhanced performance in untrained acoustic conditions; 3) More severe acoustic conditions have been considered during training to improve the generalization of the proposed model.

The rest of this paper is organized as follows. In Section 2, the formulation of the problem, the target features and related networks involved in this framework are introduced. Then, a detailed description of the correntropy based multi-objective multi-channel network is presented. Following, the experimental evaluations of the proposed method are provided in Section 3. Finally, Section 4 concludes the paper.

## 2. Multi-channel Speech Enhancement based on Correntropy

In the following subsections, the noisy-reverberant signal model for microphone array speech enhancement is first introduced. Then, the related features (LPS and IRM), the architecture of long-short term memory (LSTM) and the definition of correntropy are elaborated in details. Finally, the process of correntropy based multi-objective multi-channel speech enhancement is described, including the overall learning framework and the procedure of network training.

### 2.1. Problem Formulation

Consider a reverberant and noisy environment, the signals received by  $M$  microphones are typically modeled as

$$\mathbf{X}(n, k) = \mathbf{H}(k) \mathbf{S}(n, k) + \mathbf{V}_d(n, k) + \mathbf{V}(n, k) \quad (1)$$

in the frequency domain. Here,  $\mathbf{X}(n, k)$ ,  $\mathbf{S}(n, k)$ ,  $\mathbf{V}_d(n, k)$  and  $\mathbf{V}(n, k)$  are the  $M$ -dimensional vectors that denote the short-time Fourier transform (STFT) of the received signals, clean signals, environmental noises and spatially uncorrelated microphone self-noises at the  $n$ -th time frame and the  $k$ -th frequency bin, i.e.  $\mathbf{X}(n, k) = [X_1(n, k), X_2(n, k), \dots, X_M(n, k)]^T$ , and  $\mathbf{H}(k)$  denotes the STFT of the acoustic transfer function. In this paper, the goal of speech enhancement is to find a function that maps the noisy-reverberant observation  $\mathbf{X}(n, k)$  to the clean speech component that approximates  $\mathbf{S}(n, k)$  as close as possible.

## 2.2. Spectral Features

In recent studies, the speech log power spectrum (LPS) is usually preferred as the learning targets because of the broad dynamic range of spectral magnitude. Given the received  $M$ -dimensional time domain signals, the corresponding frame length and frame shift are set as  $N$  and  $N/2$  samples. Then, a Fourier transform is applied to each overlapping windowed frame, and each frame can be described using a vector  $\mathbf{x}(n)$  as

$$\mathbf{x}(n) = [\mathbf{X}(n, 0), \mathbf{X}(n, 1), \dots, \mathbf{X}(n, N)]^T \quad (2)$$

Since speech is correlated from frame to frame, we incorporate temporal dynamics by joining adjacent frames into a single feature vector. Therefore, the input features for network are extended as

$$\tilde{\mathbf{x}}(n) = [\mathbf{x}(n-l), \dots, \mathbf{x}(n), \mathbf{x}(n+1), \dots, \mathbf{x}(n+l)]^T \quad (3)$$

where  $l$  is the number of adjacent frames on each side, then the total number of involved frames is  $2l+1$ . Besides, similar to the input of network, the desired output of network is also the log power spectra at the  $n$ -th frame, and all the input and output features are synchronously extracted from noisy and clean speeches to keep aligning on each frame.

## 2.3. Ideal ratio mask

Recent studies [18], [37] have shown that ratio mask targets are superior to other ones in terms of the objective intelligibility and quality metrics. Thus, one of the mask targets called ideal ratio mask (IRM) [37] is employed as another output feature in the proposed network. IRM is a kind of soft mask to measure the presence of speech in a T-F unit. Generally, its value ranges from zero to one. In this work, after considering the statistical independence between clean speech and noise, the IRM of each T-F bin is defined as

$$IRM(n, k) = \frac{|S(n, k)|}{|X(n, k)|} \quad (4)$$

where  $S(n, k)$  and  $X(n, k)$  denote the clean- and noisy- spectral magnitudes at the  $k$ -th frequency bin in the  $n$ -th frame, respectively. From Eq. (4), it can be seen that unlike the convention IRM, the values of IRM defined here may greater than one. Thus, in order to avoid unwanted amplification of noise components in the signal due to estimation errors and obtain better numerical stability in backpropagation training, the values are saturated to one.

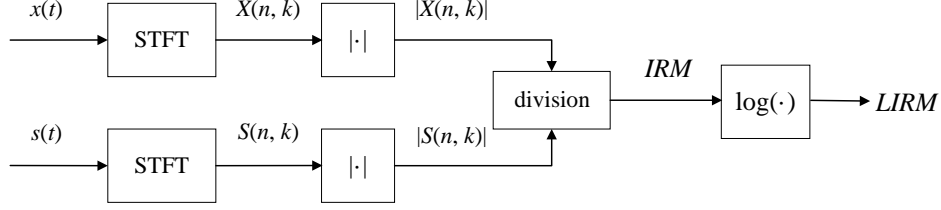
In addition, since IRM is adopted as an intermediate output to merge with another intermediate output LPS, a logarithmic operation is performed on all IRM for convenience, denoted as LIRM.

The whole computation process of LIRM is shown in Fig .1.

## 2.4. Long-short Term Memory Network

Recurrent neural network (RNN) displayed excellent performance in solving time series learning problems by adding the connections between output and hidden layers. However, it has a main obstacle as gradient vanishing problem [38] which limits the capability of learning long-range context dependencies. To address this problem, an improved RNN version is introduced as LSTM





**Fig. 1. Computation process of LIRM.**

[39], which defines memory cell and several gates to regulate the information flow. Fig. 2 illustrates the structure of LSTM memory block. It can be seen that each LSTM block comprises three gates, input gate  $\mathbf{i}_t$ , forget gate  $\mathbf{f}_t$  and output gate  $\mathbf{o}_t$ , and the main implementation can be described as follows

$$\mathbf{i}_t = \kappa(\mathbf{W}_{xi}\mathbf{x}_t + \mathbf{W}_{hi}\mathbf{h}_{t-1} + \mathbf{b}_i) \quad (5)$$

$$\mathbf{f}_t = \kappa(\mathbf{W}_{xf}\mathbf{x}_t + \mathbf{W}_{hf}\mathbf{h}_{t-1} + \mathbf{b}_f) \quad (6)$$

$$\mathbf{c}_t = \mathbf{f}_t \odot \mathbf{c}_{t-1} + \mathbf{i}_t \odot g(\mathbf{W}_{xc}\mathbf{x}_t + \mathbf{W}_{hc}\mathbf{h}_{t-1} + \mathbf{b}_c) \quad (7)$$

$$\mathbf{o}_t = \kappa(\mathbf{W}_{xo}\mathbf{x}_t + \mathbf{W}_{ho}\mathbf{h}_{t-1} + \mathbf{b}_o) \quad (8)$$

$$\mathbf{h}_t = \mathbf{o}_t \odot g(\mathbf{c}_t) \quad (9)$$

where  $\mathbf{x}_t$  and  $\mathbf{h}_t$  are the input and output vectors at the time frame  $t$ ,  $\mathbf{c}_t$  is the memory cell state,  $\mathbf{W}(\mathbf{W}_{xi}, \mathbf{W}_{hi}, \mathbf{W}_{xf}, \mathbf{W}_{hf}, \mathbf{W}_{xc}, \mathbf{W}_{hc}, \mathbf{W}_{xo}, \mathbf{W}_{ho})$  are the weight matrices that need to be learned during training and  $\mathbf{b}(\mathbf{b}_i, \mathbf{b}_f, \mathbf{b}_c, \mathbf{b}_o)$  are the bias term of the corresponding gates.  $\odot$  is the element-wise vector product,  $\kappa(\cdot)$  and  $g(\cdot)$  are the well-known sigmoid function and the hyperbolic tangent (tanh) function, which are usually used as gate activation function and input and output activation function.

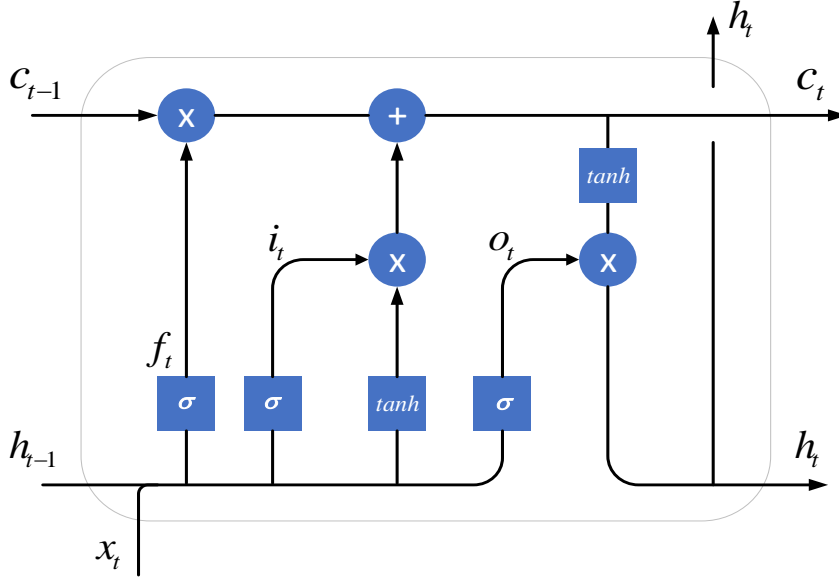
In this work, the bidirectional long short-term memory (BiLSTM) network [40] is employed. A BiLSTM network has recurrent connections in both forward and backward directions, and this structure provides the output layer with complete past and future context in the input sequence. Thus, the BiLSTM can make full use of the temporal information and more suitable for speech enhancement.

## 2.5. Correntropy and maximum correntropy criterion

For neural network training, one of the key issues is to select an appropriate loss function to measure the error between the output and the target so that the performance of whole network is the most optimal. As a classical second-order statistic, mean square error (MSE) is preferred as a loss function in many supervised learning networks. However, it is sensitive to non-Gaussian noises and outliers so that the capability of feature learning would be fragile when the input is high noisy-reverberation data. Hence, in this paper, correntropy [30] is employed as the loss function to optimize network performance.

Given two random variables  $A$  and  $B$ , the cross correntropy (CC) (usually simply correntropy) is defined as [31]

$$V_\sigma(A, B) = E[\kappa_\sigma(A - B)] \quad (10)$$



**Fig. 2.** An illustration of the LSTM block.

where  $E$  is the expectation operator,  $\kappa_\sigma(\cdot)$  is the Gaussian kernel with  $\sigma$  is the kernel size:

$$\kappa_\sigma(a - a_i) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{\|a - a_i\|^2}{2\sigma^2}\right) \quad (11)$$

Usually, due to the fact that the joint probability density function (PDF) is difficult to obtain and only a finite number of data  $\{(a_i, b_i)\}_{i=1}^N$  are available, the sample estimation of the cross correntropy can be obtained by

$$\hat{V}_{N,\sigma}(A, B) = \frac{1}{N} \sum_{i=1}^N \kappa_\sigma(a_i - b_i) \quad (12)$$

Correntropy is a measure of the similarity between two random variables in a small neighborhood determined by the kernel size. Compared to MSE, these two measures have both similarities and significant differences. Generally, correntropy and MSE can be used as a certain optimization criterion as well as a measurement of similarity between random variables  $A$  and  $B$ . However, MSE is a global statistical function, which is more suitable for Gaussian distribution. When the samples far away from the distribution center, MSE will significantly amplify the errors, thus, it cannot be optimal in the case of asymmetric error distribution, non-zero center and outliers. In contrast, correntropy is a local criterion of similarity, and has the capability to suppress the outliers because of the nonlinear kernel function, such as Gaussian kernel function. According to the characteristic of Gaussian function, the **larger** the outliers, the better the suppression. Thus, it is very suitable for cases when the error samples are non-Gaussian with large outliers [31, 34].

As a result, correntropy of the error is employed as a cost function in our network training, the goal of which is to maximize the similarity between the predicted output and the true values in the sense of correntropy. This is also called maximum correntropy criterion (MCC). In addition,

the MCC is based on the property that maximizing the correntropy is equivalent to minimizing another metric named correntropy induced metric (CIM), i.e. the smaller the CIM, the larger the MCC, and the higher similarity between two features. In [31], the CIM is defined as

$$\text{CIM}(X, Y) = (\kappa_{\sigma}(0) - V_{\sigma}(A, B))^{1/2} \quad (13)$$

which can be seen as a linear operation of correntropy and a constant. Therefore, in this paper, to facilitate using the MCC for network training, CIM is adopted as the loss function.

## 2.6. Multi-objective Multi-channel Speech Enhancement Scheme based on correntropy

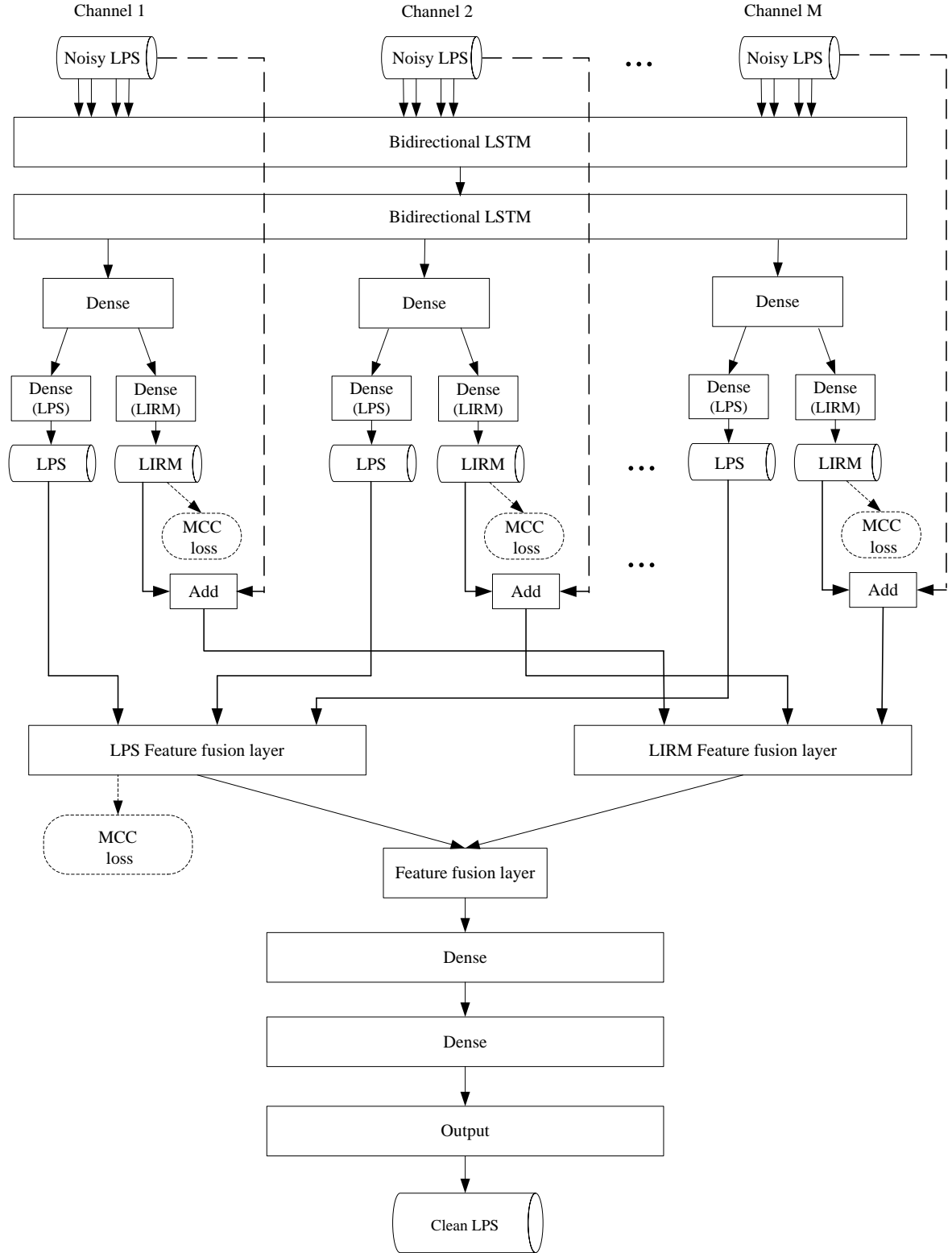
Fig. 3 presents the architecture of the proposed microphone array speech enhancement network, which mainly consists of two stages: (1) multi-objective based multi-channel feature learning; (2) fusion and output. Overall, by using an array with  $M(M > 2)$  microphones, a multi-channel signal can be collected from reverberant-noisy speech. Then, LPS features of these signals are extracted and used as the input of the network to simultaneously predicted the intermediate outputs (LPS and LIRM) of multiple channels. Subsequently, these multi-channel intermediate outputs are fused into a single-feature LPS and finally acquired the output of the proposed network. More specifically, in the first stage, a BiLSTM module composed of two layers BiLSTM is employed for the sake of learning the temporal spatial characteristics between the multi-channel LPS input and their corresponding clean LPS and LIRM. Following is the fully connected module, which  $M$  Dense layers are first performed as  $M$  output channels, and then each output channel predicts intermediate LPS and LIRM output by using two independent Dense layers. Here, LIRM is estimated through the BiLSTM and DNN modules by supervised learning. In the second stage, fusion layers are separately operated to cope with the directly predicted LPS and the indirect LPS, where the indirect LPS is computed by the estimated LIRM and the corresponding noisy LPS, and these two features are then fused to form the final single-channel LPS feature. Here, the average sum operation is used for all the fusion layers

$$\hat{Y}_d^f(n, k) = \frac{1}{M} \sum_{m=1}^M \hat{Y}_m(n, k) \quad (14)$$

$$\hat{Y}_i^f(n, k) = \frac{1}{M} \sum_{m=1}^M [Y_m(n, k) + I_m(n, k)] \quad (15)$$

where,  $\hat{Y}_m(n, k)$ ,  $Y_m(n, k)$  and  $I_m(n, k)$  are the directly predicted clean LPS, the corresponding noisy LPS input and mask from channel  $m$ ,  $Y_m(n, k) + I_m(n, k)$  is the masking-based indirect LPS feature. Ultimately, the ensemble results  $\hat{Y}_d^f(n, k)$  and  $\hat{Y}_i^f(n, k)$  are merged into a single channel output and feed to the fully connected layers to get the desired clean LPS. It should be pointed out that because the fusion LPS is used as secondary prediction target, otherwise, all the  $2M$  channel features can be directly fused into a single one.

As for the network configuration, 1024 hidden units are adopted in every BiLSTM layer. For fully connected module in each channel, 512 hidden units are first used and then 129 units are introduced with the aim of predicting the intermediate LPS and LIRM. Meanwhile, after the fusion



**Fig. 3.** Block diagram of correntropy based multi-objective multi-channel speech enhancement scheme.

layers, 512 hidden units are used for two DNN layers, and 129 units are employed in output layer to estimate the clean LPS. Moreover, tanh and hard-sigmoid are used as the activation function and recurrent activation function for recurrent layers, respectively. All the DNNs are trained with the rectified linear unit (ReLU) and the linear unit is utilized for the estimation of LIRM and LPS.

Since the network contains multiple output layers with multiple targets to be predicted, the following multi-objective MCC (minimum CIM) is constructed as the loss function

$$L = \frac{1}{T} \left[ \sum_{m=1}^M \text{CIM}(\hat{I}_m, I_m) + \text{CIM}(\hat{S}_{LPS}, S) + \text{CIM}(\hat{S}, S) \right] \quad (16)$$

where  $\hat{I}_m$  and  $I_m$  indicate the predicted LIRM and its corresponding clean LIRM in the  $m$ -th channel; Similarly,  $\hat{S}_{LPS}$ ,  $\hat{S}$  and  $S$  indicate the intermediate fused LPS, the final predicted LPS and their corresponding clean LPS; and  $T$  is the mini-batch size. Here, all the time and frequency indexes are omitted for more intuitive representation. Besides, Adam is utilized as the optimizer to train the network.

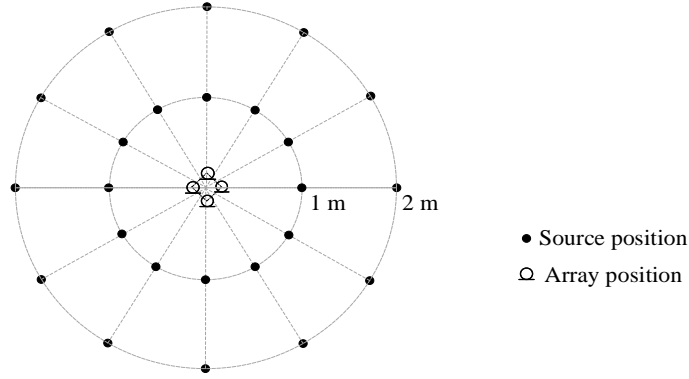
As described above, the exploitation of multiple objectives and multiple channels can fully utilize the complementarity of speech features and increase the diversity of training data [41]. Moreover, the fusion layer can alleviate the problem of overestimation or underestimation of the enhanced spectrum. Besides, correntropy based loss function is insensitive to outliers and more robust to non-Gaussian noises. In summary, the proposed scheme has a better potential for denoising and dereverberation.

### 3. Experimental Results and Discussions

In this section, several comparative experiments are presented to evaluate the performance of the proposed method. First, the data generation and experimental setup are introduced. Then, the generalization capability of the proposed scheme is evaluated under different acoustic conditions (RIRs, SNRs and noise types). Finally, the source-array distance is also discussed to further verify the robustness of the algorithm.

#### 3.1. Dataset and Experimental Setup

As shown in Fig. 4, a square array with  $M = 4$  microphones is considered in this paper, and the distance from each microphone to the center of the square is 10 cm. For both training and testing, since the proposed method is aimed to be independent of the location of the speech source, the complete angular range of the array center is discretized with a step 5 degree, and 72 different angular positions of the speech source are obtained. Specifically, to generate the training data, five rooms with different acoustic conditions are employed, as shown in Table 2. For each array position and source-array position in each room, 8 speech sentences from TIMIT training set [42], are used for each angular position, i.e. convolved with the room impulse responses (RIRs) corresponding to the specific setup. Here, the RIRs that simulate different acoustics conditions are generated by the RIR generator [43]. Therefore, a total of  $8 \times 72 = 576$  clean speech sentences are available in each array position during the training phase.



**Fig. 4.** Geometric setup for data generation.

As for diffuse noise used in training phase, the NOISEX-92 database [44] which contains 15 types of noises are considered. By randomly choosing one of 15 types, the noisy speech are generated with signal-noise-ratio (SNR) in -5dB, 0dB, 5dB and 10dB. Meanwhile, spatially uncorrelated white Gaussian noise with 20dB SNR is also added as microphone self-noise.

**Table 2** Acoustic Conditions in Five Different Rooms for Training

Parameters	Configuration on training data
Speech	TIMIT training set
Noise	NOISEX-92
Room	R1:(5×4×3)m <sup>3</sup> , R2:(6×5×3)m <sup>3</sup> , R3:(10×5×3)m <sup>3</sup> , R4:(10×8×3)m <sup>3</sup> , R5:(11×8×3)m <sup>3</sup>
Array positions	5 different positions in each room
Source-array distance	1m, 2m
RT60	R1: 0.3s, R2: 0.4s, R3: 0.6s, R4: 0.8s, R5: 0.9s
SNR	Diffuse: -5dB to 10dB, Spatially white Gaussian: 20dB

For testing, different from training configuration, two mismatched acoustic conditions with arbitrary array positions as well as different source array distances are taken into account, as shown in Table 3. Moreover, for each array setup in each room, five sentences from TIMIT testset are used for each angular position of the source, thus, a total of 360 different sentences are available in each array position during the testing phase.

In the proposed network, both BiLSTM and Dense networks are initialized with random weights, the hidden units are 1024 and 512, respectively. During training, the mini-batch size is 1024, the learning rate is 0.001 in the first 10000 epoch and then decreased to 0.0001. Besides, for each channel, all the clean and noisy waveforms are resampled to 8KHz, the corresponding frame length and frame shift are set as 256 samples (i.e. 32 ms) and 128 samples, respectively. Then, the log-power spectra of each frame are extracted as input feature, and its vector dimension is 129. All the implementations are done in Keras [45].

**Table 3** Acoustic Conditions in Five Different Rooms for Test

Parameters	Configuration on test data
Speech	TIMIT test set
Noise	Aurora2
Room	R1:(8×5×3)m <sup>3</sup> , R2:(9×7×3)m <sup>3</sup>
Array positions	3 arbitrary positions in each room
Source-array distance	1.2m, 1.7m
RT60	R1: 0.45s, R2: 0.7s
SNR	Diffuse: -2dB to 7dB, Spatially white Gaussian: 20dB

### 3.2. Reference Methods and Evaluation Metrics

For comparison purposes, we construct mapping-based, masking-based and waveform-based methods as baseline. For the mapping-based method, the log power spectra of clean speech are directly estimated from those of noisy reverberant speech. For the masking-based method, time-frequency masks (IRM) are predicted and then combined with noisy reverberant speech to obtained clean speech. For the waveform-based method, time-domain waveforms of clean speech are directly estimated from corresponding noisy speech, without feature extraction. During training in mapping- and masking-based networks, MSE and correntropy are separately adopted as loss function to evaluate its effect on network performance. For convenience, in the following experiments, MSE-based methods are referred as **M-LPS** and **M-IRM** while correntropy-based methods are referred as **C-LPS** and **C-IRM**, respectively. Moreover, to make a fair comparison, the structure of above mapping- and masking-based methods are basically same as that of the proposed method, where two-layer BiLSTM and two-layer Dense are employed, followed by a fusion block which combines the intermediate multi-channel output into single one, and lastly a Dense module is incorporated to acquire the final output. Meanwhile, all the parameters of above comparison methods are set same as the proposed network. **In addition, a latest narrow-band multi-channel method named as channel-padding narrow band deep filtering (CP-NBDF) [46] is also employed as baseline. It is a masking-based method, where each frequency of one signal is processed independently. The model is composed of two-layer BiLSTM and one-layer Dense, and the input is multi-channel STFT coefficients associated with a single frequency bin, and the output is the corresponding IRM. With the aim of fair comparison, for CP-NBDF, the power spectra of  $M = 4$  channels are first organized and the last  $M' - M$  channels are padded with zeros ( $M'$  is the upper bound for channels). The network setup and the input sequence length are set as Ref. [46]. As for the waveform-based method, a multi-channel input and output network, called multi-resolution convolutional auto-encoder (MRCAE), is implemented according to [47], where the encoder and decoder are composed of two convolutional layers with sets of filters and two transposed convolution layers with sets of filters, respectively, and another transposed convolution layer is used for the final output. For MRCAE, the number of the filters for each convolutional layer is set as Ref. [47].**

The enhanced speech signals from each approach are evaluated with **four** objective metrics, namely perceptual evaluation of speech quality (PESQ) [48], short-time objective intelligibility

(STOI) score [49], extended STOI (ESTOI) [50] and **deep noise suppression mean opinion score (DNSMOS) [51]**. PESQ is computed by comparing the enhanced speech with the corresponding clean speech, producing scores in range  $[-0.5, 4.5]$  where a higher score indicates a better quality. STOI and ESTOI measure speech intelligibility by computing the correlation of short-time temporal envelopes between clean and enhanced speech, resulting in scores in the range of  $[0, 1]$  where a higher score indicates a better intelligibility. **DNSMOS is a perceptual speech quality metric that serves as a proxy for subjective scores, ranging from 0 to 5, where a higher score indicates a better hearing for listeners.**

### 3.3. Experiments with different acoustic conditions

In this section, a experiment is first designed to select a suitable value of kernel size for correntropy. Then, the generalization capability of the proposed method is explored by introducing two different acoustic environments in the following experiments. Besides, all the parameters configurations and comparison methods are depicted in Subsection 3.1 and 3.2, and the array positions showed in Table 3 are randomly selected in each room.

#### 3.3.1. Selection of kernel size $\sigma$

Since the correntropy is based on a Gaussian kernel function, the choice of the kernel size should always be relative to the dynamic range of the variable or signal. In other words, a suitable value of  $\sigma$  is crucial to the performance of the network. However, speech enhancement is a regression task, the error between true values and predicted values are not fixed in a range like the classification tasks, thus, when using the correntropy as loss function, the kernel size  $\sigma$  has to be chosen through the experimental results.

Generally speaking, for denoising and dereverberation, we are more concerned about whether the proposed method has better adaptability under different acoustic environments (Rooms). Therefore, this paper gives priority to the enhancement results of the C-LPS+IRM under different RIR conditions when choosing the kernel size. Table 4 lists the enhancement performance of the proposed C-LPS+IRM with different kernel size  $\sigma$  under different RIRs. From Table 4, when the value of  $\sigma$  is 1.0, 1.3, 1.4 and 1.5, the performance of the model tends to be the best. However, by looking in Table 5, although the PESQ results are similar when  $a = 1.0, 1.3, 1.4$  and  $1.5$ , more attention is paid to low SNR conditions for the overall intelligence of speech. Thus,  $\sigma = 1.0$  is the most suitable kernel size and is fixed in the following experiments.

**Table 4** Average PESQ with different kernel size  $\sigma$ .

$\sigma$	PESQ	$\sigma$	PESQ	$\sigma$	PESQ
0.5	2.212	1.2	2.215	1.7	2.215
0.8	2.215	1.3	2.218	1.8	2.217
0.9	2.212	1.4	2.218	1.9	2.215
1.0	2.219	1.5	2.218	2.0	2.216
1.1	2.215	1.6	2.217	2.5	2.215



**Table 5** PESQ results at different input SNRs with different kernel size  $\sigma$ .

SNR(dB)	-5	0	5	10
$\sigma=1.0$	1.937	2.166	2.334	2.439
$\sigma=1.3$	1.931	2.166	2.334	2.440
$\sigma=1.4$	1.932	2.169	2.330	2.439
$\sigma=1.5$	1.934	2.166	2.331	2.440

### 3.3.2. Generalization to different RIRs

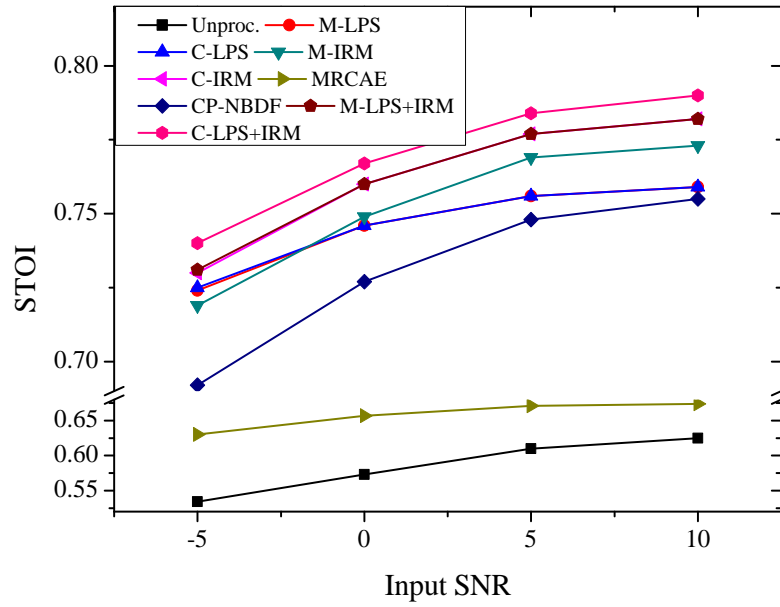
1) Generalization to impulse noise: To evaluate the robustness of correntropy based model against impulse noise, a typical impulse noise machinegun is tested in this experiment. As shown in Table 3, different RIRs and testing speech are employed, and the array positions in each room are randomly selected as well. Table 6 and Fig. 5 display the PESQ and STOI results under different noisy-reverberant conditions by using different methods. Obviously, C-LPS+IRM has a great advantage in suppressing machinegun noise, especially in the case of low input SNRs. However, in turn, CP-NBDF has the worst performance. Specifically, for PESQ metric, the average enhanced performance of C-LPS+IRM outperforms that of the second best C-IRM with 0.02 improvement. In addition, for low input SNRs (-5dB and 0dB), the gaps between these two models achieve to 0.038 and 0.025, respectively. Furthermore, for STOI metric, C-LPS+IRM still exhibits superior enhanced performance in all noisy-reverberant conditions. M-LPS+IRM and C-IRM jointly rank the second best with a comparable performance.

In addition to PESQ and STOI, another metric called extended STOI (ESTOI) is also employed to further evaluate the proposed method. Different from STOI, ESTOI works for a larger range of input signals, and does not assume mutual independence between frequency bands. Fig. 6 presents the ESTOI results of machinegun noise among seven models. It can be observed that although the results of ESTOI are decreased compared with those of convention STOI, the trends of model performance are basically consistent with STOI. C-LPS+IRM ranks first with a larger advantage compared with the second best C-IRM and M-LPS+IRM. M-IRM has the third rank, and MRCAE has the worst performance. From above, it can be stated that compared with MSE based models, the proposed correntropy based model is more capable of dealing with impulse noise.

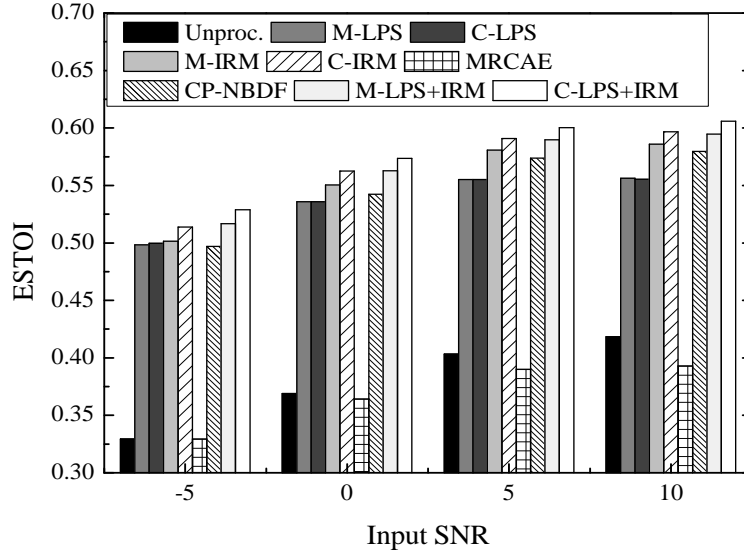
2) Overall performance evaluation: To analyze the proposed model more comprehensively, the whole 15 noise types from NOISEX-92 are tested among different models. Same as the above experiment, the configurations of acoustic conditions are set as Table 3. Table 7 and Fig. 7 list the average PESQ and STOI results among **eight** models under different RIRs with trained noise types at different SNRs. It can be observed that first, all the seven models have the capability to obtain clean speech from noisy-reverberation speech. Second, for the mapping- and masking-based models, the enhanced performance of correntropy based models (C-LPS+IRM, C-IRM and C-LPS) are generally superior to those of MSE based models (M-LPS+IRM, M-IRM and M-LPS) in terms of both PESQ and STOI improvements. Third, although the IRM-based models generally perform better than LPS-based models, a further performance gain can be obtained by using the two features in combination, which implies that the LPS and IRM features are complementary to some extent. More specifically, PESQ results in Table 7 illustrate that **in most cases**, C-LPS+IRM

**Table 6** PESQ results of machinegun noise for different models.

SNR(dB)	PESQ			
	-5	0	5	10
Unprocessed	1.423	1.725	1.924	2.032
M-LPS	2.176	2.248	2.285	2.304
C-LPS	2.180	2.247	2.292	2.308
M-IRM	2.158	2.321	2.431	2.469
C-IRM	2.212	2.354	2.447	2.491
MRCAE	2.110	2.156	2.182	2.199
CP-NBDF	2.057	2.153	2.183	2.196
M-LPS+IRM	2.215	2.352	2.440	2.476
C-LPS+IRM	<b>2.253</b>	<b>2.379</b>	<b>2.462</b>	<b>2.501</b>



**Fig. 5.** STOI results of machinegun noise for different models.



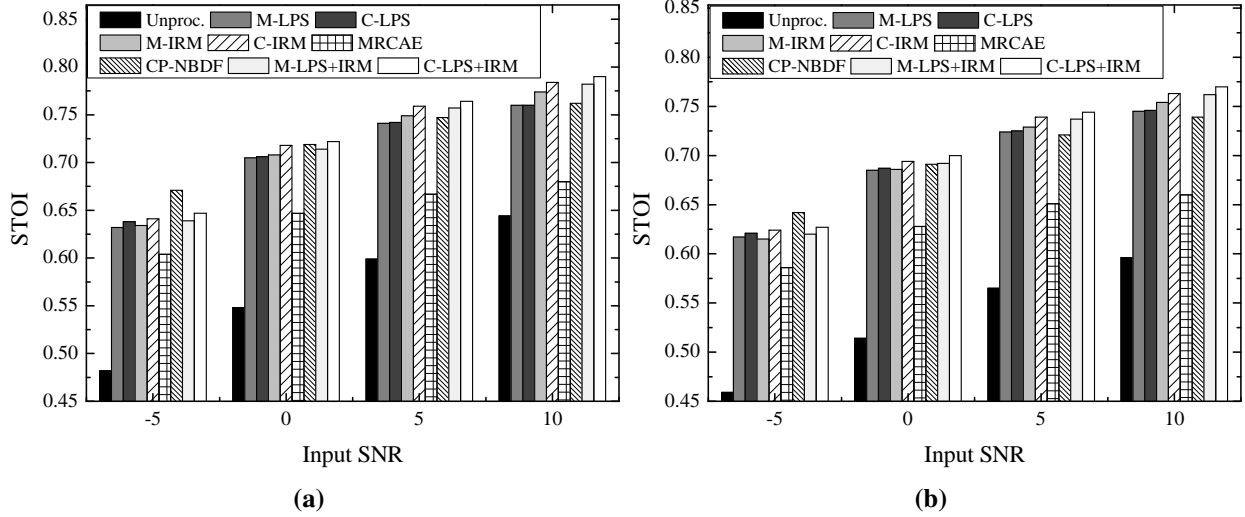
**Fig. 6.** ESTOI results of machinegun noise for different models.

significantly outperforms other models with an average PESQ increment from 1.753 to 2.254 in Room 1, 1.682 to 2.183 in Room 2. But when the input SNR is low (-5dB), CP-NBDF exhibits the best performance, and C-LPS+IRM ranks the second with a gap of 0.046. This may because CP-NBDF processes each frequency independently, which makes less frequency distortion at low input SNR. Following are C-IRM and M-LPS+IRM with comparable performance, the improvements of PESQ are 0.485 and 0.48 in Room 1, 0.481 and 0.483 in Room 2, respectively. However, it can be noted that the performance of C-LPS is inferior to M-IRM, which indicates that a proper loss function does improve the performance of network, but still has a certain limitation. That is to say, many factors, such as feature selection, have effects on network performance. Moreover, STOI results in Fig. 7 show a similar trend to those of PESQ. C-LPS+IRM still ranks the first in most SNR conditions, with the average improvement of 0.163 and 0.178 in Room 1 and Room 2, respectively. C-IRM and M-LPS+IRM have second and third rank with a gap of 0.5 and 0.7 percent, respectively.

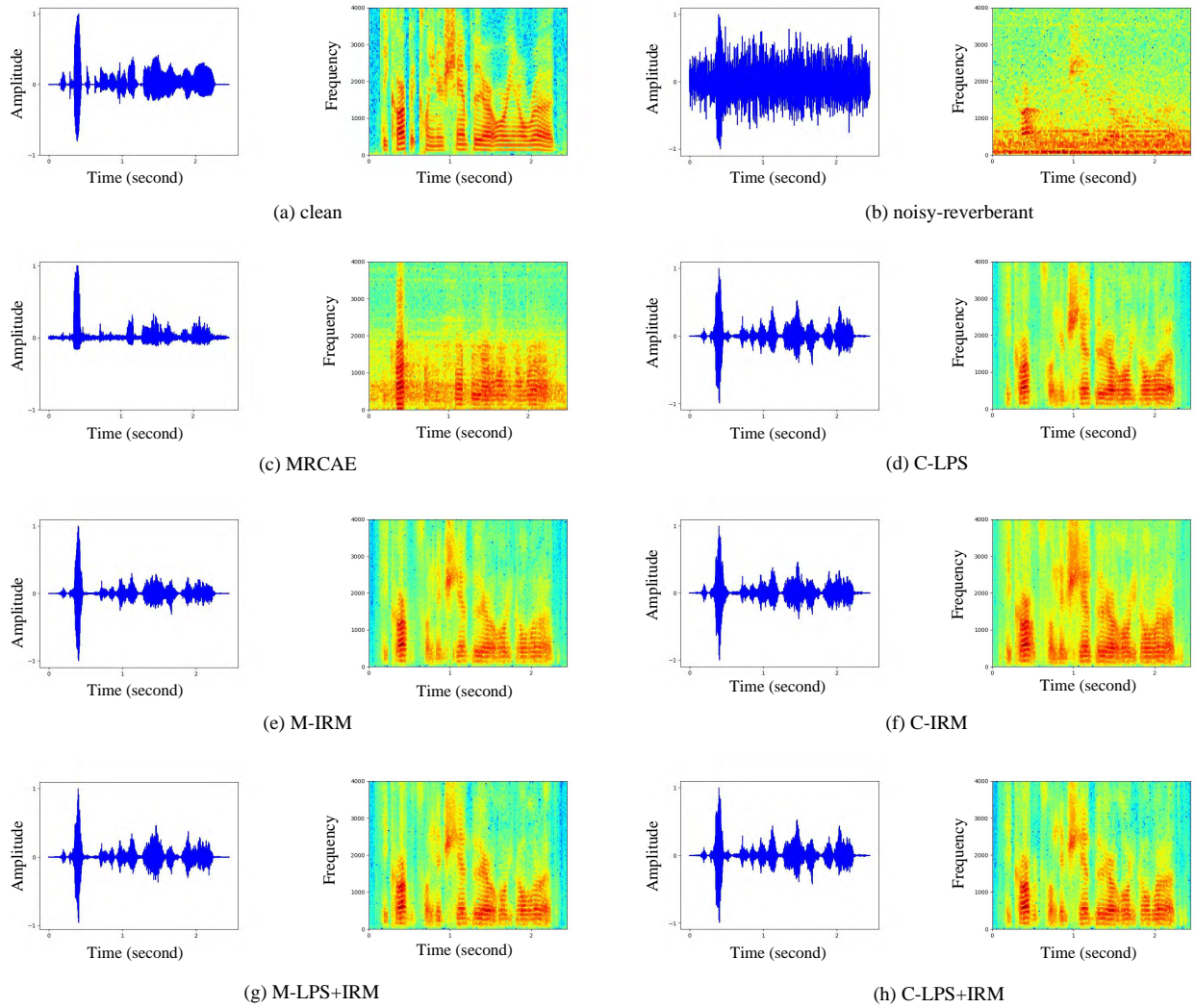
To further illustrate the proposed method, an enhancement example is presented in Fig. 8. The waveforms and spectrograms of clean, noisy-reverberant and the corresponding enhanced speech obtained by MRCAE, C-LPS, M-IRM, C-IRM, M-LPS+IRM and C-IRM+IRM are shown in Fig. 8 (a), (b), (c), (d), (e), (f), (g) and (h), respectively. By comparing waveforms in Fig. 8(a) and Figs. 8(c)-(h), it is clear that the enhanced utterance from MRCAE is distorted notably, IRM-based methods lead to amplitude deviations, and LPS method loses speech segments at some frames. Moreover, by looking at the spectrograms, the utterance restored by C-LPS+IRM has a clearer structure, M-LPS+IRM is slightly inferior but still competitive. Following, C-IRM and C-LPS preserve excessive noise in the high-frequency components, and M-LPS exhibits some distortion in the low-frequency components. Besides, MRCAE fails to remove the background noise and reverberation.

**Table 7** Average PESQ results under untrained RIRs at different SNRs.

SNR(dB)	Room1				Room2			
	-5	0	5	10	-5	0	5	10
Unprocessed	1.454	1.678	1.867	2.014	1.427	1.614	1.783	1.906
M-LPS	1.901	2.104	2.232	2.304	1.888	2.061	2.170	2.242
C-LPS	1.912	2.109	2.243	2.314	1.894	2.066	2.180	2.252
M-IRM	1.904	2.158	2.341	2.467	1.871	2.092	2.258	2.355
C-IRM	1.933	2.183	2.356	2.482	1.898	2.118	2.270	2.368
MRCAE	1.951	2.073	2.154	2.211	1.896	2.006	2.075	2.114
<b>CP-NBDF</b>	<b>2.018</b>	2.141	2.194	2.239	<b>1.948</b>	2.031	2.089	2.112
M-LPS+IRM	1.936	2.174	2.351	2.472	1.904	2.114	2.278	2.367
C-LPS+IRM	1.952	<b>2.197</b>	<b>2.374</b>	<b>2.493</b>	1.922	<b>2.133</b>	<b>2.294</b>	<b>2.386</b>



**Fig. 7.** Average STOI results in (a) Room 1 and (b) Room 2 at different SNRs.



**Fig. 8.** Waveforms and spectrograms of an example utterance: (a) clean speech; (b) noisy-reverberant speech (m109 noise, SNR=-5dB, RT60=0.7s); (c) enhanced speech by MRCAE; (d) enhanced speech by C-LPS; (e) enhanced speech by M-IRM; (f) enhanced speech by C-IRM; (g) enhanced speech by M-LPS+IRM; (h) enhanced speech by C-LPS+IRM.

### 3.3.3. Generalization to untrained SNRs

For the sake of illustrating the proposed C-LPS+IRM model is not sensitive to untrained SNRs, in this experiment, we choose -2dB and 7dB that are not involved in training as input SNRs. Table 8 and Fig. 9 present the average PESQ and STOI scores among six models with trained noise types at untrained SNRs. On the whole, it can be found that the model still attain excellent enhancement capability even facing with untrained SNRs, and are basically consistent with the results of trained SNRs shown in Table 7 and Fig. 7. For PESQ, C-LPS+IRM provides improvements of 0.518 and 0.487 over unprocessed speech in Room 1 and 2, ranks the first. Next, M-LPS+IRM has the second best with the 0.504 and 0.512 improvement in Room 1 and Room 2. For STOI, overall, the proposed C-LPS+IRM still outperforms other models with nearly 2% improvements compared with noisy speech. **But by looking more closely, CP-NBDF is slightly better under low SNR condition.** Besides, the performance of M-LPS+IRM model are inferior to those of C-IRM model, which is slightly different from PESQ results. But taken together, they have a comparable performance, jointly rank the second best.

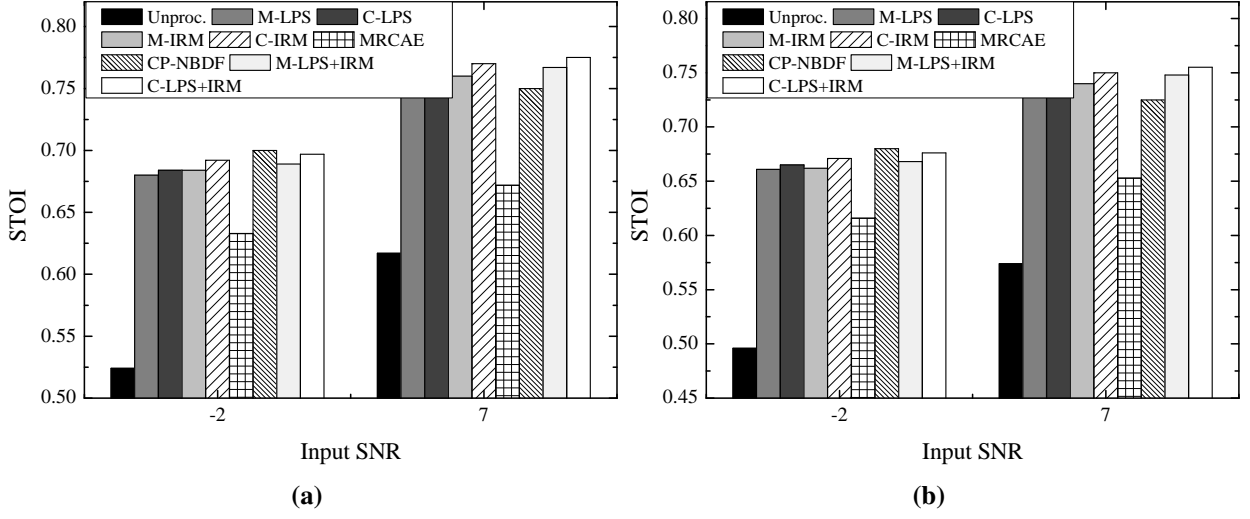
**Table 8** Average PESQ results at different untrained SNRs.

SNR(dB)	Room1		Room2	
	-2	7	-2	7
Unprocessed	1.588	1.934	1.531	1.835
M-LPS	2.035	2.261	2.002	2.201
C-LPS	2.040	2.272	2.013	2.211
M-IRM	2.072	2.394	2.022	2.298
C-IRM	2.087	2.409	2.043	2.312
MRCAE	2.042	2.176	1.973	2.090
<b>CP-NBDF</b>	2.102	2.211	2.010	2.097
M-LPS+IRM	2.092	2.398	2.050	2.315
C-LPS+IRM	<b>2.106</b>	<b>2.421</b>	<b>2.066</b>	<b>2.332</b>

### 3.3.4. Generalization to untrained noise types

In addition to the previous experiment that verified the adaptability of the proposed model to untrained SNRs, the robustness to unseen noise types is also one of the evaluation criteria. Thus, in this experiment, Aurora2 dataset [52] is incorporated as unseen test noises to further assess the robustness of seven models to various non-stationary noises. Tables 9 and Fig. 10 display the PESQ and STOI results of unprocessed and processed speech under different noisy-reverberant conditions by using different methods. Compared with results in Table 7 and Fig. 7, it is clear that when processing untrained noise types, the enhanced performance of seven models are all decreased, which implies that their robustness are declined more or less. However, the trends of their performance are similar to those of matching noises conditions but still have some difference. Specifically, the PESQ results in Table 9 state that CP-NBDF outperforms other models at low input SNRs (-5dB and 0dB), while C-LPS+IRM yields a superior performance at high input SNRs (5dB and 10dB). **This may because CP-NBDF is a frequency-wise method, while others are**





**Fig. 9.** Average STOI results in (a) Room 1 and (b) Room 2 at different untrained SNRs.

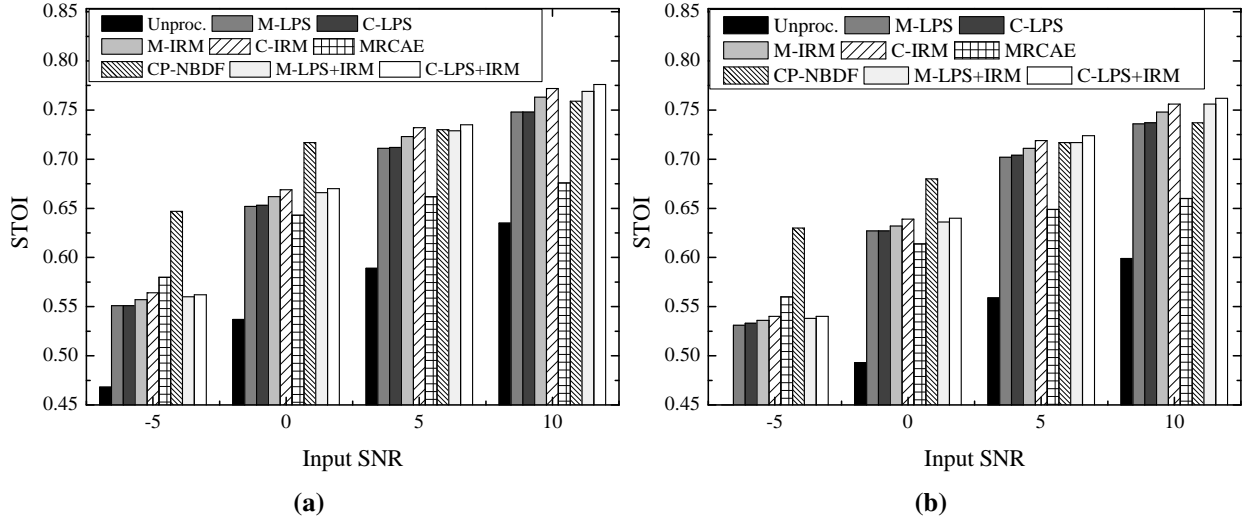
channel-wise methods. Specifically, frequency-wise method uncouples the inter-frequency dependency from full-band signal, and processes each frequency bin separately to make a full use of the narrow-band spatial and temporal information. However, in contrast, channel-wise methods process signal in full-band, and could not focus on all the information from every frequency bin, especially under strong noisy conditions, which in turn remove the useful frequency domain information of clean speech and lead to speech distortion in time domain. Meanwhile, the average STOI results in Fig. 10 also reveal that at low input SNRs (-5dB and 0dB), the performance of C-LPS+IRM are inferior to those of CP-NBDF. For other cases, C-LPS+IRM usually produces better STOI results than others. Thus, taken together, the proposed method has certain advantages in coping with various noise types and is suitable for common environments.

To further assess the proposed method, subjective measurement is also carried out by using DNSMOS test. DNSMOS is a proxy for subjective scores and contains three aspects: signal distortion (CSIG), intrusiveness of background noise (CBAK) and the overall quality (COVL). In this test, 200 utterances are randomly selected from Room 1 and Room 2, and the average results among different models are presented in Table 10. It is clear that the speech processed by MRCAE have the least distortion and the best overall quality, while other methods damage the origin speech more or less. A possible reason is that MRCAE is a waveform-based method, which causes less damage to frequency bins and could retain more information, including useful clean speech and residual noise. Thus, the enhanced speech are less distortion and more smoothly for listeners. However, for background noise intrusiveness, the proposed C-LPS+IRM performs the best, while the capability of MRCAE is degraded significantly. This is consistent with the above analysis, where MRCAE retains more noise component and C-LPS+IRM removes more noise component as well as the useful clean component, resulting in more speech distortion, and ranks the second place in CSIG and CVOL metrics. Besides, CP-NBDF has the worst performance in DNSMOS metric. After listening the speech processed by CP-NBDF, we found that although the speech are

smooth, some other frequencies of sound are introduced into the background, making it hard for listeners to distinguish the specific contents.

**Table 9** Average PESQ results with untrained noise types at different SNRs.

SNR(dB)	Room1				Room2			
	-5	0	5	10	-5	0	5	10
Unprocessed	1.452	1.668	1.853	2.008	1.441	1.603	1.792	1.896
M-LPS	1.635	1.895	2.089	2.224	1.616	1.842	2.056	2.179
C-LPS	1.634	1.895	2.095	2.233	1.621	1.852	2.066	2.182
M-IRM	1.658	1.948	2.194	2.371	1.627	1.886	2.135	2.286
C-IRM	1.679	1.966	2.205	2.378	1.641	1.907	2.149	2.304
MRCAE	1.898	2.065	2.134	2.3198	1.850	1.957	2.074	2.100
<b>CP-NBDF</b>	<b>1.963</b>	<b>2.098</b>	2.172	2.218	<b>1.899</b>	<b>2.001</b>	2.075	2.093
M-LPS+IRM	1.664	1.954	2.194	2.371	1.639	1.895	2.148	2.298
C-LPS+IRM	1.663	1.964	<b>2.214</b>	<b>2.386</b>	1.643	1.904	<b>2.163</b>	<b>2.310</b>



**Fig. 10.** Average STOI results in (a) Room 1 and (b) Room 2 with untrained noise types at different SNRs.

### 3.3.5. Generalization to untrained dataset

In order to comprehensively explore the generalization of the proposed method, an untrained dataset named Deep Noise Suppression (DNS) [53] is employed in this experiment to further evaluate the enhanced performance among seven models. The configurations of test condition are the same as Table 2 except that the noises are randomly selected from DNS wideband noise dataset. Besides, each of the clean speech is mixed with different noise types to fully validate the



**Table 10** Average DNSMOS results with untrained noise types at different SNRs.

SNR(dB)	CSIG				CBAK				COVL			
	-5	0	5	10	-5	0	5	10	-5	0	5	10
Unprocessed	3.056	3.215	3.375	3.545	1.886	2.166	2.564	2.769	2.350	2.519	2.727	2.867
M-LPS	2.795	3.062	3.226	3.396	2.816	3.054	3.238	3.386	2.267	2.380	2.527	2.684
C-LPS	2.828	3.083	3.230	3.407	2.829	3.044	3.234	3.390	2.284	2.378	2.530	2.700
M-IRM	2.746	2.961	3.194	3.370	2.807	3.049	3.222	3.431	2.294	2.372	2.539	2.721
C-IRM	2.797	3.050	3.208	3.449	2.762	3.018	3.251	3.444	2.195	2.399	2.568	2.790
MRCAE	<b>3.352</b>	<b>3.471</b>	<b>3.511</b>	3.453	2.694	2.847	2.987	3.069	<b>2.649</b>	<b>2.780</b>	<b>2.880</b>	<b>2.835</b>
CP-NBDF	2.434	2.619	2.709	2.861	2.272	2.329	2.501	2.500	1.824	1.884	1.971	2.100
M-LPS+IRM	2.795	3.051	3.219	3.403	2.773	3.041	3.278	3.438	2.280	2.375	2.581	2.740
C-LPS+IRM	2.844	3.085	3.282	<b>3.467</b>	<b>2.853</b>	<b>3.077</b>	<b>3.259</b>	<b>3.462</b>	2.285	2.399	2.584	2.805

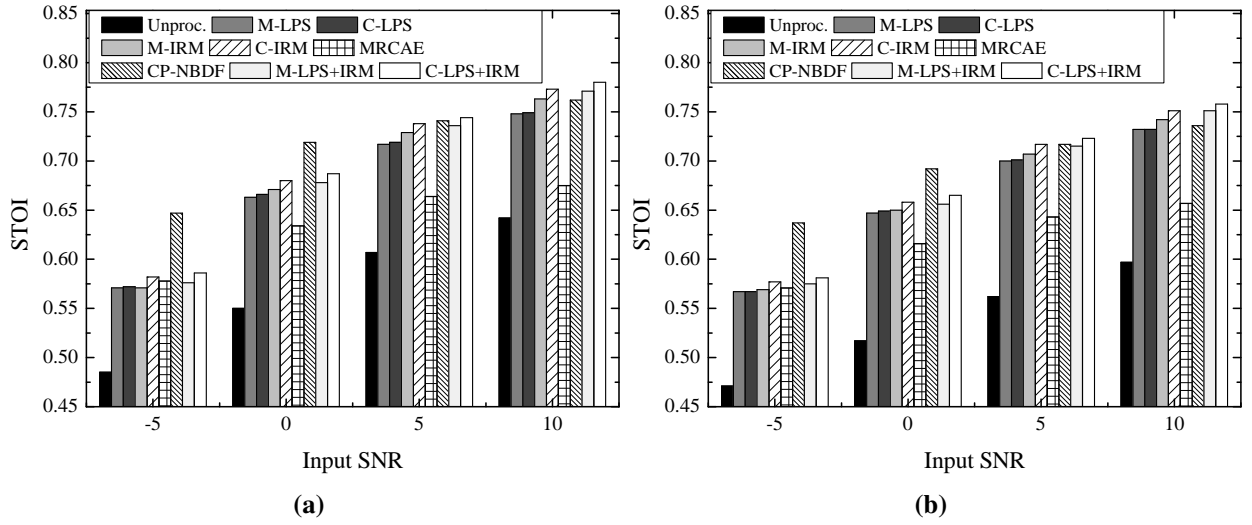
effectiveness of the proposed method. Table 11 and Fig. 11 show the average PESQ and STOI results in processing DNS dataset with different models. From Table 11, it can be seen that the PESQ results of dealing with DNS dataset are similar to those of untrained noise types in Table 9. The frequency-wise method CP-NBDF has a better performance at low input SNRs, followed by MRCAE and C-LPS+IRM. However, when the input SNRs are 5dB and 10dB, C-LPS+IRM ranks the first place, and following are CP-NBDF and MRCAE. Moreover, for other mapping- and masking- based methods, C-LPS+IRM still outperforms others under different acoustic conditions, which confirms its generalization capability. Meanwhile, STOI results also exhibit encouraging performance of the proposed model. C-LPS+IRM stands out among all models at high SNRs (5dB and 10dB) and has the second best at low SNRs (-5dB and 0dB). C-IRM has the third rank, and MRCAE has the fourth rank only at -5dB input SNR.

**Table 11** Average PESQ results with DNS dataset at different SNRs.

SNR(dB)	Room1				Room2			
	-5	0	5	10	-5	0	5	10
Unprocessed	1.401	1.597	1.768	1.893	1.394	1.546	1.708	1.805
M-LPS	1.590	1.852	2.063	2.215	1.596	1.838	2.027	2.216
C-LPS	1.595	1.872	2.076	2.213	1.601	1.839	2.042	2.178
M-IRM	1.572	1.901	2.161	2.358	1.569	1.857	2.100	2.269
C-IRM	1.607	1.923	2.187	2.373	1.602	1.884	2.120	2.295
MRCAE	1.850	2.004	2.112	2.166	1.838	1.966	2.041	2.090
CP-NBDF	<b>1.967</b>	<b>2.101</b>	2.168	2.209	<b>1.905</b>	<b>1.996</b>	2.062	2.081
M-LPS+IRM	1.605	1.943	2.171	2.351	1.596	1.871	2.111	2.288
C-LPS+IRM	1.616	1.945	<b>2.195</b>	<b>2.381</b>	1.616	1.903	<b>2.133</b>	<b>2.301</b>

### 3.3.6. Generalization to untrained source-array distances

In this experiment, two source-array distances (1.2m and 1.7m) that are not used in training are selected to further explore the generalization capability of the proposed model. All the other



**Fig. 11.** Average STOI results in (a) Room 1 and (b) Room 2 with DNS dataset at different SNRs.

conditional configurations, including the noise types and input SNRs, are set as same as training stage. Table 12 shows the average STOI and PESQ results with different source-array distances among seven models. In general, no matter in short or long distance, C-LPS+IRM always exhibits substantial improvements, ranking the first; IRM and MS-LPS+IRM are both inferior to C-LPS+IRM, having the second best and the third rank. By observing the results in details, when the distance is 1.2m, the average PESQ increment of C-LPS+IRM is 0.014 compared to the second best C-IRM; However, when the distance is 1.7m, the increment of C-LPS+IRM is considerably increased to 0.027, which is almost twice compared to that of 1.2m distance. Meanwhile, STOI results shows that for short source-array distance, C-LPS+IRM has an average 0.148 improvements, which just slightly better than C-IRM with 0.146; But for long source-array distance, the STOI improvement of C-LPS+IRM is increased to 0.189, and yields 0.7 percent superior than that of C-IRM. In addition, it can be found that for STOI metric, the advantage of proposed C-LPS+IRM are more obvious in Room 2 (strong noisy-reverberation condition). The improvements achieve to 0.148 and 0.206 which comparable or even better than 0.148 and 0.171 in Room 1.

From above analyses, it can be concluded that the proposed C-LPS+IRM models significantly improves the dereverberation and denoising performance over unprocessed signals in untrained acoustic conditions, and this indicates that it has great prospect for speech enhancement under severe environments.

### 3.4. ASR Experiments

Previous experiments have demonstrated the capability of the correntropy based multi-objective multichannel (C-LPS+IRM) method in speech enhancement task. To further explore the potential of the proposed method in practical applications, it is used as a front end of automatic speech recognition (ASR) under noisy-reverberant conditions. Here, it should be noted that, similar to PESQ and STOI metrics, ASR is employed as an additional objective evaluation metric, rather

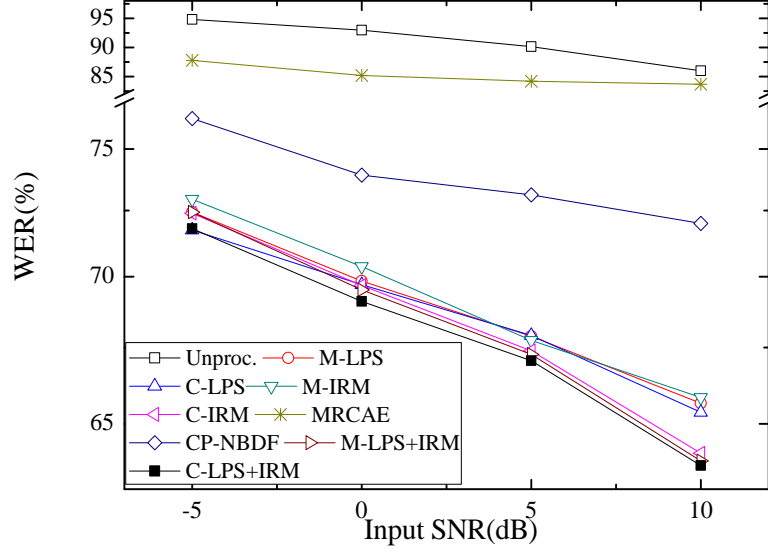
**Table 12** Average STOI and PESQ results with different source-array distances.

Distance(m)	PESQ				STOI			
	Room1		Room2		Room1		Room2	
	1.2	1.7	1.2	1.7	1.2	1.7	1.2	1.7
Unprocessed	1.755	1.744	1.719	1.655	0.586	0.546	0.590	0.480
M-LPS	2.133	2.121	2.123	2.058	0.711	0.702	0.710	0.676
C-LPS	2.137	2.133	2.136	2.063	0.713	0.704	0.712	0.679
M-IRM	2.216	2.185	2.215	2.087	0.722	0.701	0.728	0.668
C-IRM	2.236	2.200	2.227	2.109	0.731	0.711	0.736	0.679
MRCAE	2.075	2.010	2.117	1.934	0.667	0.618	0.643	0.597
CP-NBDF	2.155	2.125	2.087	2.009	<b>0.737</b>	0.716	0.732	0.672
M-LPS+IRM	2.229	2.206	2.221	2.115	0.727	0.710	0.729	0.679
C-LPS+IRM	<b>2.249</b>	<b>2.229</b>	<b>2.241</b>	<b>2.133</b>	0.734	<b>0.717</b>	<b>0.738</b>	<b>0.686</b>

than pursuing a superior ASR performance. In this experiment, we employed a trained ASR system provided by Google (Google Speech Recognition) to check the recognition performance, and it is measured by word error rate (WER). The acoustic environment is the same as the overall performance evaluation test in Subsection 3.3.2, and the averaged WER results of different models in Room 2 are displayed in Fig. 12. It can be observed that all the enhanced speech could obtain lower WER compared to the noisy ones, and the average reduction is about 23%. Specifically, the proposed C-LPS+IRM outperforms other models under all SNR conditions, M-LPS+IRM and C-IRM have the second and third rank with a small gap. Meanwhile, M-LPS, C-LPS and M-IRM show competitive performance at low input SNRs, and CP-NBDF and MRCAE have the worst ASR results. Besides, by looking in details, although all the models can reduce WER, the results are not remarkable, which indicates that their enhanced capabilities are limited. This might be due to the strong reverberation in Room 2 (RT60=0.7s), making it more difficult to extract clean speech from noisy environments, and further result in speech distortion and accordingly deteriorate the performance of ASR. Nonetheless, C-LPS+IRM still stands out among all the models and can be simply used as a pre-processing of ASR system. Besides, it also shows that the proposed C-LPS+IRM has great potential in practical applications.

#### 4. Conclusion

In order to address both the non-Gaussian noise and reverberation problems, we propose a multi-objective based multi-channel speech enhancement network by using correntropy as the loss function. First, the log-power spectra (LPS) of multiple channels noisy speech are employed as the input of BiLSTM network. Then, multiple correntropy-based loss functions are used to simultaneously estimate the intermediate LPS and log-ideal ratio mask (LIRM) of each channel. Subsequently, two fusion layers are adopted to separately fuse the intermediate LPS and LIRM into two single-channel features, which are then integrated into a single-channel LPS. Ultimately, a deep neural network is brought in to further approximate the nonlinear mapping from the final fused LPS to the clean speech LPS. Moreover, during training, the correntropy based loss function



**Fig. 12.** WER results of Google ASR with different methods. (The WER for clean speech is 30.42%)

is employed due to its insensitivity to the outliers (or impulsive noise) for the sake of improving the network performance. Experimental evaluations show that the proposed method has the superiority in suppressing non-Gaussian noise and exhibits good generalization capability to different noises, SNRs, datasets and source-array distances. However, one shortcoming of this method is the limited performance under low SNR conditions. Thus, in the future, we will design the network to simultaneously learn narrow-band and full-band information, or combined the time and frequency information to further improve the learning capability of the method.

## Acknowledgments

This work was supported by National Natural Science Foundation of China (Nos. 61771091, 61871066), National High Technology Research and Development Program (863 Program) of China (No. 2015AA016306), Natural Science Foundation of Liaoning Province of China (No. 20170540159), and Fundamental Research Funds for the Central Universities of China (Nos. DUT17LAB04).

## Data Availability Statement

The datasets generated during the current study are available from the corresponding author on reasonable request.

## Conflict of Interest Statement

We declare that we do not have any commercial or associative interest that represents a conflict of interest in connection with the work submitted.

## References

- [1] J. Li, L. Deng, R. Haeb-Umbach, Y. Gong, *Robust Automatic Speech Recognition: A bridge to practical applications* (Academic Press, New York, 2015)
- [2] J. Benesty, S. Makino, J. Chen, *Speech Enhancement* (Springer, Berlin, 2005)
- [3] P. C. Loizou, *Speech Enhancement: Theory and Practice* (CRC Press, Florida, 2013)
- [4] S. Boll, Suppression of acoustic noise in speech using spectral subtraction. *IEEE Trans. Acoust. Speech Signal Process.* 27(2), 113-120 (1979)
- [5] J. S. Lim, A. V. Oppenheim, Enhancement and bandwidth compression of noisy speech. *Proc. of IEEE* 67(12), 1586-1604 (1979)
- [6] Y. Ephraim, H. L. Van Trees, A signal subspace approach for speech enhancement. *IEEE/ACM Trans. Audio Speech Lang. Process.* 3(4), 251-266 (1995)
- [7] I. Cohen, S. Gannot, *Springer Handbook of Speech Processing* (Springer, Berlin, 2008)
- [8] T. Yoshioka, T. Nakatani, Generalization of multi-channel linear prediction methods for blind MIMO impulse response shortening. *IEEE Trans. Acoust. Speech Signal Process.* 20(10), 2707-2720 (2012)
- [9] S. T. Neely, J. B. Allen, Invertibility of a room impulse response. *J. Acoust. Soc. Amer.* 66, 165-169 (1979)
- [10] K. Han, Y. Wang, D. L. Wang, W. S. Woods, I. Merks, T. Zhang, Learning spectral mapping for speech dereverberation and denoising. *IEEE/ACM Trans. Audio Speech Lang. Process.* 23(6), 982-992 (2015)
- [11] D. S. Williamson, D. L. Wang, Time-frequency masking in the complex domain for speech dereverberation and denoising. *IEEE/ACM Trans. Audio Speech Lang. Process.* 25(7), 1492-1501 (2017)
- [12] Y. Zhao, Z. Wang, D. Wang, Two-stage deep learning for noisy-reverberant speech Enhancement. *IEEE/ACM Trans. Audio Speech Lang. Process.* 27(1), 53-62 (2019)
- [13] R. Li, X. Sun, T. Li, F. Zhao, A multi-objective learning speech enhancement algorithm based on IRM post-processing with joint estimation of SCNN and TCNN. *Digit. Signal Process.* 101, 1-11 (2020)
- [14] S. Gannot, E. Vincent, S. Markovich-Golan, A. Ozerov, A consolidated perspective on multi-microphone speech enhancement and source separation. *IEEE/ACM Trans. Audio Speech Lang. Process.* 25(4), 692-730 (2017)
- [15] N. Yousefian, P. C. Loizou, A dual-microphone speech enhancement algorithm based on the coherence function. *IEEE/ACM Trans. Audio Speech Lang. Process.* 20(2), 599-609 (2012)
- [16] T. Shan, T. Kailath, Adaptive beamforming for coherent signals and interference. *IEEE Trans. Acoust. Speech Signal Process.* 33(3), 527-536 (1985)
- [17] I. Tashev, A. Acero, Microphone array post-processor using instantaneous direction of arrival, in *International Workshop on Acoustic, Echo and Noise Control (IWAENC)*, Paris, France, (2006)
- [18] Y. Wang, A. Narayanan, D. Wang, On training targets for supervised speech separation. *IEEE Trans. Acoust. Speech Signal Process.* 22(12), 1849-1858 (2014)
- [19] Y. Jiang, D. Wang, R. Liu, Z. Feng, Binaural classification for reverberant speech segregation using deep neural networks. *IEEE Trans. Acoust. Speech Signal Process.* 22(12), 2112-2121 (2014)
- [20] X. Zhang and D. Wang, Deep learning based binaural speech separation in reverberant environments. *IEEE Trans. Acoust. Speech Signal Process.* 25(5), 1075-1084 (2017)
- [21] S. Araki, T. Hayashi, M. Delcroix, M. Fujimoto, K. Takeda, T. Nakatani, Exploring multi-channel features for denoising-autoencoder-based speech enhancement, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brisbane, Australia (2015) pp. 116-120
- [22] X. Sun, R. Xia, J. Li, Y. Yan, A deep learning based binaural speech enhancement approach with spatial cues preservation, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brighton, United Kingdom (2019) pp. 5766-5770

- [23] J. Heymann, L. Drude, R. Haeb-Umbach, Neural network based spectral mask estimation for acoustic beamforming, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Shanghai, China (2016) pp. 196-200
- [24] X. Xiao, S. Zhao, D. L. Jones, E. S. Chng, H. Li, On time-frequency mask estimation for MVDR beamforming with application in robust speech recognition, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, New Orleans, USA (2017) pp. 3246-3250
- [25] S. Chakrabarty, D. Wang, E. A. P. Habets, Time-frequency masking based online speech enhancement with multi-channel data using convolutional neural networks, in *International Workshop on Acoustic Signal Enhancement (IWAENC)*, Tokyo, Japan (2018) pp. 476-480
- [26] S. Chakrabarty, E. A. P. Habets, Time-frequency masking based online multi-channel speech enhancement with convolutional recurrent neural networks. *IEEE J. Sel. Topics Signal Process.* 13(4), 787-799 (2019)
- [27] X. Cui, Z. Chen, F. Yin, Multi-objective based multi-channel speech enhancement with BiLSTM network. *Appl. Acoust.* 177, (2021)
- [28] T. Higuchi, K. Kinoshita, N. Ito, S. Karita, T. Nakatani, Frame-by-frame closed-form update for mask-based adaptive MVDR beamforming, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Calgary, Canada (2018) pp. 531-535
- [29] J. Qi, H. Hu, Y. Wang, C. H. Yang, S. Marco Siniscalchi, C. Lee, Tensor-to-vector regression for multi-channel speech enhancement based on tensor-train network, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Barcelona, Spain (2020) pp. 7504-7508
- [30] I. Santamaria, P. P. Pokharel, J. C. Principe, Generalized correlation function: Definition, properties, and application to blind equalization. *IEEE Trans. Signal Process.* 54(6), 2187-2197 (2006)
- [31] W. Liu, P. P. Pokharel, J. C. Principe, Correntropy: Properties and applications in non-Gaussian signal processing. *IEEE Trans. Signal Process.* 55(11), 5286-5298 (2007)
- [32] P. P. Pokharel, W. Liu, J. C. Principe, A low complexity robust detector in impulsive noise. *Signal Process.* 89(10), 1902C1909 (2009)
- [33] A. Singh, J. C. Principe, A loss function for classification based on a robust similarity metric, in *International Joint Conference on Neural Networks (IJCNN)*, Barcelona, Spain (2010) pp. 1-6
- [34] A. Singh, R. Pokharel, J. C. Principe, The c-loss function for pattern classification. *Pattern Recognit.* 47(1), 441-453 (2014)
- [35] Y. Qi, Y. Wang, X. Zheng, Z. Wu, Robust feature learning by stacked autoencoder with maximum correntropy criterion, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Florence, Italy (2014) pp. 6716-6720
- [36] L. Chen, H. Qu, J. Zhao, B. Chen, J. C. Principe, Efficient and robust deep learning with correntropy-induced loss function. *Neural Computing Appl.* 27(4), 1019-1031 (2016)
- [37] D. Wang, J. Chen, Supervised speech separation based on deep Learning: An overview. *IEEE Trans. Acoust. Speech Signal Process.* 26(10), 1702-1726 (2018)
- [38] K. Greff, R. K. Srivastava, J. Koutník, B. R. Steunebrink, J. Schmidhuber, LSTM: A search space odyssey. *IEEE Trans. Neural Netw. Learn. Syst.* 28(10), 2222-2232 (2017)
- [39] S. Hochreiter, J. Schmidhuber, Long short-term memory. *Neural Comput.* 9(8), 1735-1780 (1997)
- [40] A. Graves, J. Schmidhuber, Framewise phoneme classification with bidirectional LSTM and other neural network architectures. *Neural Netw.* 18(5-6), 602-610 (2005)
- [41] L. Sun, J. Du, L. Dai, C. Lee, Multiple-target deep learning for LSTM-RNN based speech enhancement, in *Hands-Free Speech Communications and Microphone Arrays, (HSCMA)*, San Francisco, CA (2017) pp. 136C140
- [42] J. S. Garofolo, L. Lamel, W. M. Fisher, J. G. Fiscus, D. S. Pallett, N. L. Dahlgren, Darpa TIMIT acoustic-phonetic continuous speech corpus. 1993. [Online]. Available: <https://github.com/philipperemy/timit>
- [43] E. A. P. Habets, Room impulse response (RIR) generator. 2016. [Online]. Available: <https://github.com/ehabets/RIR-Generator>
- [44] A. Varga, H. J. M. Steeneken, Assessment for automatic speech recognition: II. NOISEX-92: A database and an experiment to study the effect of additive noise on speech recognition systems. *Speech Commun.* 12(3), 247-251 (1993)

- [45] F. Chollet *et al.*, Keras. 2015. [Online]. Available: <https://github.com/fchollet/keras>
- [46] S. Zhang, X. Li. Microphone array generalization for multichannel narrowband deep speech enhancement, in *Interspeech*, Brno, Czech (2021) pp. 666-670
- [47] E. M. Grais, D. Ward, M. D. Plumbley, Raw multi-channel audio source separation using multi-resolution convolutional auto-encoders, in *European Signal Processing Conference (EUSIPCO)*, Rome, Italy (2018) pp. 1577-1581
- [48] ITU-T, Recommendation P.862: Perceptual evaluation of speech quality (PESQ): An objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs. Technical Report, 2001.
- [49] C. H. Taal, R. C. Hendriks, R. Heusdens, J. Jensen, A short-time objective intelligibility measure for time-frequency weighted noisy speech, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Dallas, USA (2010) pp. 4214-4217
- [50] J. Jensen, C. H. Taal, An algorithm for predicting the intelligibility of speech masked by modulated noise maskers. *IEEE Trans. Acoust Speech Signal Process.* 24(11), 2009-2022 (2016)
- [51] C. K. A. Reddy, V. Gopal, R. Cutler, DNSMOS P.835: A non-intrusive perceptual objective speech quality metric to evaluate noise suppressors. (2021). [arXiv:2110.01763](https://arxiv.org/abs/2110.01763)
- [52] H. G. Hirsch, D. Pearce, The AURORA experimental framework for the performance evaluations of speech recognition systems under noisy conditions, in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Beijing, China (2000) pp. 181-188
- [53] C. K. A. Reddy *et al.*, The Interspeech 2020 deep noise suppression challenge: Datasets, subjective speech quality and testing framework (2020)

**Correntropy based Multi-objective Multi-channel Speech Enhancement  
(ID: CSSP-D-21-00332R1)**

Xingyue Cui, Zhe Chen, Fuliang Yin, Xianfa Xu

**The Point-by-point Responses to the Reviewer's Comments**

We would like to express our appreciation to the reviewers for providing us with valuable comments for improving this manuscript. In the following, we present the point-by-point replies to the reviewers' comments.

**Reviewer1:**

**Comments No.1:** *Still, in Abstract, line number 24-29 is not clear, may have to reframe it.*

**Answer:**

Thanks for the reviewer's advice. We have reframed the sentences in line number 24-29. Please refer to the Abstract of the revised manuscript.

**Comments No.2:** *Typo, page-7, line 52, lager-->larger*

**Answer:**

Sorry for our carelessness. The "lager" has been modified as "larger" with red font in Page 7 in the revised manuscript.

**Comments No.3:** *The authors should be consistent with their notations, either IRM or LIRM can be used, If LIRM is used have to change in fig.1 and fig 3 accordingly.*

**Answer:**

Thanks for the reviewer's valuable advice. We have modified the notations to make them consistent, and Fig. 1 and Fig. 3 have also been changed accordingly.

**Comments No.4:** *Still, the computation of LIRM in the intermediate hidden layer is not clear. Are the authors passing both the noisy and clean LPS signal to the intermediate layer to compute the ratio? (clearly mention the same in the respective section)*

**Answer:**

Sorry for our unclear description. As shown in Fig. 3, LIRM feature is one of the intermediate outputs, which is estimated or computed through the BiLSTM and DNN modules by supervised learning. Both the noisy and clean LPS signals are not used during the computation of LIRM in the intermediate layer. In fact, the ultimate goal of the proposed method is to obtain the LPS feature of the clean speech, thus, according



to Eq. (4), the noisy LPS signal is passed and added with the corresponding LIRM to obtain the indirectly estimated LPS of clean speech in each channel. And the multi-channel indirectly estimated LPS are then fused into one channel and incorporated with directly estimated LPS to form a single-channel LPS feature for further estimating the final LPS of clean speech.

According to the reviewer's advice, we have added the corresponding description in the revised manuscript. Please refer to the first paragraph in Section 2.6.

**Comments No.5:** *In the experiment section, the author showed MRCAE outperforms the proposed approach in low SNR scenarios in the case of untrained noise and dataset conditions. The author should justify it appropriately, and should also mention which attribute of MRCAE helps to perform better in low SNR scenarios. The author can include this in the conclusion section and can provide appropriate motivation for future work.*

**Answer:**

Thanks for the reviewer's comment. Indeed, as shown in Tables 9 and 11, MRCAE outperforms the proposed method in low SNR scenarios. This may because MRCAE is a waveform-based learning method, while the proposed method is a feature-based learning method. Generally, the waveform-based methods process speech in time domain, which causes less damage to frequency bins and could retain more information, including useful clean speech and residual noise. Thus, when the input SNR is low, these retain information make the utterances more smoothly. But when the input SNR is high, the excessive residual noise will significantly degrade the performance, especially for speech intelligibility, as shown in Figs.10 and 11. In contrast, the feature-based learning methods are easy to remove the useful frequency domain information of clean speech during denoising and dereverberation, especially under strong noisy conditions, and lead to speech distortion in time domain.

Meanwhile, new experiments show that the latest work CP-NBDF has the best performance compared with other methods in low SNR scenarios. A possible reason is that CP-NBDF uncouples the inter-frequency dependency from full-band signal, and processes each frequency bin separately to make full use of the narrow-band information. However, MRCAE and the proposed method are all channel-wise methods, which processes the signal in full-band and could not focus on all the information from every frequency bin, and then lead to inferior performance compared to CP-NBDF.

According to the reviewer's advice, the above description about which attribute of MRCAE helps it perform well and the reason why CP-NBDF performs the best in low SNR scenarios have been added in the revised manuscript, respectively. Please refer to the last two paragraphs in Subsection 3.3.4. Besides, they are also included in Conclusion section to provide the motivation for future work.

1  
2 **Reviewer2:**  
3

4 Thanks for your previous comments that are very useful for improving our  
5 manuscript.  
6  
7

8  
9  
10 **Reviewer3:**  
11

12 **Comments No.1:** *The authors have now compared with MRCAE based method,*  
13 *which was published in 2018. Although its good to compare with that work, but the*  
14 *reviewer expected comparison with more recent work, at least those that are*  
15 *published after 2020 to have a good comparison.*  
16  
17

18  
19 *Following are some, which the authors could have tried:*  
20  
21

22 *(i) Panagiotis Tzirakis et al., "Multi-Channel Speech Enhancement using Graph*  
23 *Neural Networks", ICASSP 2021.*

24 *(ii) Papers under Multi-Channel Speech Enhancement session of Interspeech 2020*  
25 *and Interspeech 2021 for comparison.*  
26  
27

28  
29 [https://www.isca-speech.org/archive/interspeech\\_2020/index.html#Multi-Channel%20](https://www.isca-speech.org/archive/interspeech_2020/index.html#Multi-Channel%20Speech%20Enhancement)  
30 [Speech%20Enhancement](https://www.isca-speech.org/archive/interspeech_2020/index.html#Multi-Channel%20Speech%20Enhancement)  
31  
32

33  
34 **Answer:**  
35

36 According to the reviewer's suggestion, a more recent multi-channel method named  
37 channel-padding narrow band deep filtering (CP-NBDF) [A1] ([46] in revised  
38 manuscript) is employed for comparison, and the results have been added in all the  
39 experiments. Please refer to Tables 6-12 and Figs. 5-12 and the corresponding  
40 descriptions of Subsections 3.2 and 3.3 in the revised manuscript.  
41  
42

43  
44  
45 [A1] S. Zhang, X. Li. *Microphone array generalization for multichannel narrowband*  
46 *deep speech enhancement*, in: *Interspeech*, pp. 666-670 (2021)  
47  
48

49  
50 **Comments No.2:** *Although performance measure such as PESQ and STOI are used*  
51 *for speech enhancement, in the recent years the trend is more towards how speech*  
52 *enhancement can help some application. For instance, the current DNS challenge*  
53 *(<https://www.microsoft.com/en-us/research/academic-program/deep-noise-suppressio>*  
54 *n-challenge-icassp-2022/) focuses on a combined metric that is derived using*  
55 *DNSMOS P.835 (<https://arxiv.org/pdf/2110.01763.pdf>) and ASR word accuracy (as*  
56 *often it is found enhanced speech degrades the speech accuracy). So, I recommend*  
57 *authors to use such things for future works, which will make the works not only*  
58  
59  
60  
61  
62  
63  
64  
65

*interesting, but also useful for the applications.*

**Answer:**

Thanks for the reviewer's comment. DNSMOS P.835 is a perceptual objective metric that serves as a proxy for subjective scores. However, when we tried to connect to DNSMOS API provided by Microsoft, we found that the process speed is very slow and the connection often interrupts. Thus, considering the page and time limitations, the DNSMOS metric is only conducted on untrained noise types experiment. Please refer to Table 10 and the last paragraph of Subsection 3.3.4 in the revised manuscript. Moreover, ASR word accuracy experiment has also been added to further evaluate the proposed method. Please refer to Section 3.4 in the revised manuscript.