

NEURAL NETWORK-BASED VIRTUAL MICROPHONE ESTIMATOR

Tsubasa Ochiai, Marc Delcroix, Tomohiro Nakatani, Rintaro Ikeshita, Keisuke Kinoshita, Shoko Araki

NTT corporation

ABSTRACT

Developing microphone array technologies for a small number of microphones is important due to the constraints of many devices. One direction to address this situation consists of virtually augmenting the number of microphone signals, e.g., based on several physical model assumptions. However, such assumptions are not necessarily met in realistic conditions. In this paper, as an alternative approach, we propose a neural network-based virtual microphone estimator (NN-VME). The NN-VME estimates virtual microphone signals directly in the time domain, by utilizing the precise estimation capability of the recent time-domain neural networks. We adopt a fully supervised learning framework that uses actual observations at the locations of the virtual microphones at training time. Consequently, the NN-VME can be trained using only multi-channel observations and thus directly on real recordings, avoiding the need for unrealistic physical model-based assumptions. Experiments on the CHiME-4 corpus show that the proposed NN-VME achieves high virtual microphone estimation performance even for real recordings and that a beamformer augmented with the NN-VME improves both the speech enhancement and recognition performance.

Index Terms— virtual microphone, time-domain network, supervised learning, beamforming, array signal processing

1. INTRODUCTION

Microphone array signal processing [1], which uses spatio-temporal information obtained with multiple microphones, has been an active research field for several decades and has been essential for the development of many applications such as noise reduction, source separation, and source localization. The performance of array signal processing-based approaches depends heavily on the number of microphones available, and a sufficient number of microphones is required to achieve a high level of performance. However, it is not always possible to equip commercial devices with many microphones due to structural constraints and cost restrictions. For example, most smartphones have one or two microphones and only few high-end models may have up to four. Therefore, developing array signal processing schemes that can operate with a small number of microphones is an important research topic.

Recently, researchers [2] have proposed virtually increasing the number of microphones in an array by including virtual microphone signals generated by the interpolation between two real microphone observations. They derived the phase component of a virtual microphone signal as a linear interpolation of the phases of real microphone signals, by introducing several assumptions such as a physical model of the plane wave propagation, W-disjoint orthogonality of the sources [3], and limiting array size small enough to avoid spatial aliasing. The phase was combined with an estimate of the amplitude obtained, e.g., by β divergence-based estimation [2]. This approach could increase the number of microphones in an array and would improve the array processing performance in an under-determined con-

dition. However, this approach relies on the above strong assumptions that may not always hold in real conditions, and the method has only been tested on simulated array recordings.

In this paper, we propose an alternative approach to estimate virtual microphone signals based on a supervised learning framework using a time-domain neural network, which does not explicitly rely on the physical model assumptions. Recently, neural networks operating directly in the time-domain, e.g., a time-domain audio separation network (TasNet) [4, 5], have demonstrated a high level of performance for speech separation. Moreover, it is confirmed in [6] that, when applying TasNet to microphone array recordings, the spatial information (i.e., phase) could be preserved, so that beamforming could be successfully constructed based on the output of TasNet. These works have revealed the potential of neural networks to accurately estimate time-domain signals, i.e., predict both amplitude and phase information. Motivated by these studies, in this paper, we investigated the potential of time-domain neural networks to estimate virtual microphone signals directly in the time-domain by utilizing the spatial information inferred from a few observed real microphone signals. We call the method neural network-based virtual microphone estimator (NN-VME).

To enable the NN-VME to estimate the virtual microphone signals, we adopt a supervised learning framework, assuming that during training we have access to recordings at the location of the virtual microphone. This assumption is reasonable in many circumstances, since we may have fewer constraints on the number of microphones during the system development (collection of training data) than during the actual deployment. By adopting the supervised learning framework, we do not explicitly rely on the physical model assumptions, such as the plane wave propagation, unlike [2]. Moreover, the neural network requires only microphone observations as training data. Consequently, the network can be naturally trained and deployed on real recordings.

We tested the effectiveness of the proposed method both in terms of 1) virtual microphone estimation performance and 2) speech enhancement performance when combined with a beamformer. To demonstrate the potential of the proposed method to deal with real recordings, we conducted these experiments with the well-known noisy speech benchmark, CHiME-4 corpus [7], which contains real recordings from noisy public environments. Through the experiments, we confirm that the proposed NN-VME achieves high virtual microphone estimation performance and that a beamformer augmented with the proposed NN-VME improves the speech enhancement and automatic speech recognition (ASR) performance.

In this paper, as a typical use-case for the proposed method, we evaluate the combination of the NN-VME with beamforming for noise reduction. Besides the combination of the NN-VME with array processing techniques, as the other potential application, the NN-VME could also be used for sound reproduction systems to generate virtual microphone signals at user-desired locations, where actual microphones could not be placed, e.g., due to physical constraints.

2. PROPOSED METHOD: NEURAL NETWORK-BASED VIRTUAL MICROPHONE ESTIMATOR

2.1. Network architecture

Figure 1 shows the network architecture and the training procedure of the proposed NN-VME. In the figure, for simplicity, the network receives two input channels corresponding to the observed real microphone, and it generates one output channel corresponding to the estimated virtual microphone. In the following, we will explain the general case where we predict simultaneously more than one virtual microphone.

Let \mathbf{r}_c be the \mathcal{T} -length time-domain waveform of the observed signal for the c -th real microphone, and $\hat{\mathbf{v}}_{c'}$ denotes the estimated signal for the c' -th virtual microphone. Given the real microphone signals $\mathbf{r} = \{\mathbf{r}_{c=1}, \dots, \mathbf{r}_{c=C_r}\}$ as an input, the proposed NN-VME module estimates the virtual microphone signals $\hat{\mathbf{v}} = \{\hat{\mathbf{v}}_{c'=1}, \dots, \hat{\mathbf{v}}_{c'=C_v}\}$ as:

$$\hat{\mathbf{v}} = \text{NN-VME}(\mathbf{r}), \quad (1)$$

where, C_r denotes the number of observed channels, i.e., real microphones, and C_v denotes the number of virtually estimated channels, i.e., virtual microphones, and $\text{NN-VME}(\cdot)$ is a neural network.

We adopt an architecture for $\text{NN-VME}(\cdot)$ inspired by Conv-TasNet [5], as it was shown to be able to estimate time-domain signals with high accuracy for speech separation. Similar to the original Conv-TasNet, the network is mainly composed of a 1d-convolution encoder layer, internal convolution blocks, and a 1d-deconvolution decoder layer. First, the encoder layer directly maps the time-domain signals to an intermediate representation, and then the intermediate representation is further processed by several convolution blocks. Finally, the decoder layer converts the intermediate representation back to time-domain signals.

Note that the original TasNet estimates the *separated signals* at the location of the real microphone, while the proposed NN-VME estimates the *observed signals* at the location of the virtual microphone.

2.2. Supervised training

The proposed NN-VME adopts the supervised learning framework in order to enable the NN-VME module to estimate the virtual microphone signals. Thus, in the training, we use actual microphone signals at the location of the virtual microphones as training targets.

For supervised training of the proposed NN-VME, we assume that a set of input and target signals $\{\mathbf{r}, \mathbf{t}\}$ is available. Here, $\mathbf{t} = \{\mathbf{t}_{c'=1}, \dots, \mathbf{t}_{c'=C_v}\}$ and $\mathbf{t}_{c'}$ denotes the target signal for the c' -th virtual microphone. Figure 1 illustrates this situation; a subset of the microphones (e.g., channels 1 and 3) is assigned as network input \mathbf{r} , and another subset (e.g., channel 2) is used as network target \mathbf{t} .

We train the network based on the time-domain loss between the estimated and actual signals at the location of the virtual microphones. As the training loss, we adopt the scale-dependent signal-to-noise ratio (SNR) [8] as follows:

$$\mathcal{L} = \sum_{c'=1}^{C_v} 10 \log_{10} \left(\frac{\|\mathbf{t}_{c'}\|^2}{\|\mathbf{t}_{c'} - \hat{\mathbf{v}}_{c'}\|^2} \right), \quad (2)$$

where $\hat{\mathbf{v}} = \text{NN-VME}(\mathbf{r})$, as described in Eq. (1).

Note that, in contrast to the speech separation task, we simply need to estimate the virtual microphone signals and there is no permutation ambiguity, and thus permutation invariant loss is not necessary. Unlike the training for the speech enhancement techniques, the

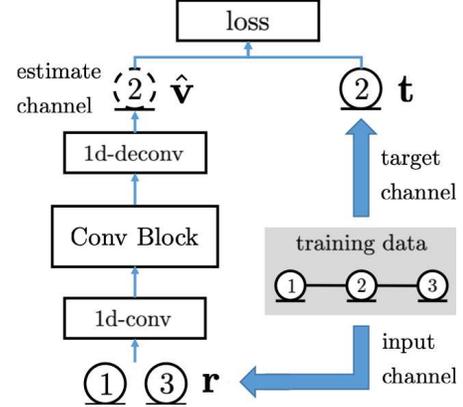


Fig. 1. Overview of network architecture and training procedure of neural network-based virtual microphone estimator

training for the proposed method does not require parallel noisy and clean signals and it only requires the multi-channel noisy observation signals, which would be relatively easy and low-cost to prepare. That is, the proposed model can be trained on real recordings, not simulated ones. By exploiting a large amount of training data, we expect that the powerful neural network will be able to provide fine modeling of real recordings.

3. APPLICATION EXAMPLE: BEAMFORMING WITH VIRTUAL MICROPHONE ESTIMATOR

The proposed NN-VME enables the generation of virtual microphone signals, and thus could be used with various array processing techniques. In this paper, to confirm the validity of the estimated virtual microphone signals for noise reduction applications, we investigated combining the NN-VME with a frequency domain beamformer.

3.1. General procedure

First, using the proposed NN-VME, we estimate the virtual microphone signals $\hat{\mathbf{v}} \in \mathbb{R}^{\mathcal{T} \times C_v}$ given the real microphone signals $\mathbf{r} \in \mathbb{R}^{\mathcal{T} \times C_r}$, as described in Eq. (1), and obtain the *augmented microphone signals* $\mathbf{y} = [\mathbf{r}, \hat{\mathbf{v}}] \in \mathbb{R}^{\mathcal{T} \times C}$, where $C = C_r + C_v$. Then, we apply a frequency-domain beamformer on top of the augmented microphone signals in frequency-domain representation, i.e., short-time Fourier transform (STFT), to obtain the enhanced speech signal. Finally, we reconstruct the enhanced time-domain waveform using the inverse STFT.

The enhanced speech signal in the STFT domain $\hat{X}_{t,f} \in \mathbb{C}$, is obtained as $\hat{X}_{t,f} = \mathbf{w}_f^H \mathbf{Y}_{t,f}$. Here, $\mathbf{Y}_{t,f} \in \mathbb{C}^C$ is a vector comprising the C -channel STFT coefficients of the augmented microphone signals at time-frequency bin (t, f) , $\mathbf{w}_f \in \mathbb{C}^C$ is a vector comprising the beamforming filter coefficients, and H represents the conjugate transpose.

3.2. Mask-based MVDR Formalization

In this paper, we adopt the minimum variance distortionless response (MVDR) formalization of [9], and we compute the time-invariant filter coefficients \mathbf{w}_f as follows:

$$\mathbf{w}_f = \frac{(\Phi_f^N)^{-1} \Phi_f^S}{\text{Tr}((\Phi_f^N)^{-1} \Phi_f^S)} \mathbf{u}, \quad (3)$$

where $\Phi_f^S \in \mathbb{C}^{C \times C}$ and $\Phi_f^N \in \mathbb{C}^{C \times C}$ are the spatial covariance (SC) matrices for speech and noise signals, respectively. $\mathbf{u} \in \mathbb{R}^C$ is a one-hot vector representing the reference microphone.

Based on [10], we approximately estimate the SC matrices using time-frequency masks as follows:

$$\Phi_f^\nu = \frac{1}{\sum_{t=1}^T m_{t,f}^\nu} \sum_{t=1}^T m_{t,f}^\nu \mathbf{Y}_{t,f} \mathbf{Y}_{t,f}^H, \quad (4)$$

where $\nu \in \{S, N\}$. $m_{t,f}^S \in [0, 1]$ and $m_{t,f}^N \in [0, 1]$ are the time-frequency masks for speech and noise, respectively.

3.3. Virtual Microphone Loading

Our preliminary experiments showed that while use of virtual microphones in beamforming was effective for increasing signal-to-distortion ratios (SDR) [11], it does not necessarily contribute to improving the ASR performance. This is probably because processing artifacts are introduced by the virtual microphone estimation. To reduce the effect of such artifacts, we introduced a virtual microphone loading term $\mathbf{Z} \in \mathbb{R}^C$ to the noise SC matrix Φ_f^N :

$$\Phi_f^N \leftarrow \Phi_f^N + \epsilon \mathbf{Z}, \quad (5)$$

where $\mathbf{Z} = \{z_{c,c'}\}_{c=1,c'=1}^{C,C}$ is a matrix of zeros excepts for the diagonal elements corresponding to the virtual microphone, i.e., $z_{c_v,c_v} = 1$, c_v denotes the channel index corresponding to the virtual microphone, and ϵ is a loading hyperparameter that controls the contribution of the virtual microphone in the construction of the beamformer. For example, when we set a larger value to ϵ , it means that the virtual microphone is contaminated by larger noise that is not correlated with other microphones. As a result, the estimated beamformer puts less weight on the virtual microphone channels.

4. RELATION TO PRIOR WORK

Our work is related to [2, 12–14], which proposed to generate virtual microphone signals by interpolation/extrapolation of signals observed at real microphones. In these methods, the phase components are estimated based on several physical model assumptions, which would make it difficult to use in practice. For example, for the CHiME-4 corpus, the microphone spacing is too large to assume no spatial aliasing and thus it cannot be applied to our experimental settings.

In [14], a supervised learning-based framework is also adopted to estimate the amplitude components of the virtual microphone. This framework minimizes the loss function between the output of the beamformer and the clean signal, but it causes overfitting problem on the open test set. In contrast, our work assumes the availability of recordings at the location of the virtual microphone during training, and it minimizes the loss function between the output of the NN-VME and actual microphone signals at the locations. We experimentally confirmed the generalization capability for the open test set and even for real recordings as shown in the experimental results of Section 5.

5. EXPERIMENT

To evaluate the proposed NN-VME, we conducted two types of evaluations: 1) evaluation of the virtual microphone estimation performance by the proposed NN-VME and 2) evaluation of the enhancement performance by the beamformer with the estimated virtual microphone. In the experiment, we report the results that estimate one

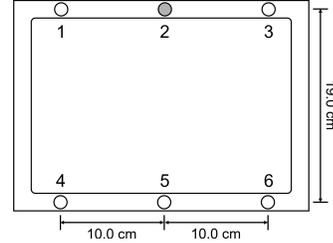


Fig. 2. Microphone array geometry for CHiME-4 corpus. All microphones face forward except for microphone 2.

virtual microphone, but our method can be naturally extended to estimate multiple virtual microphones.

5.1. Experimental conditions

We evaluated the proposed NN-VME on the CHiME-4 corpus [7]. The CHiME-4 corpus consists of speech recorded using a tablet device equipped with a rectangular microphone array with 6 channels, as illustrated in Figure 2. The corpus contains not only simulated data but also real recordings from noisy public environments.

The training set consisted of 3 hours of real speech data uttered by 4 speakers and 15 hours of simulation speech data uttered by 83 speakers. The evaluation set consisted of 1320 utterances of real and simulated noisy speech data uttered by 4 speakers, respectively. Out of these utterances, the CHiME-4 organizers reported that 12% of all the real data were affected by microphone failures, mainly for channels 4 and 5 [15]. Dealing with such microphone failures is out of the scope of this study. Therefore, to remove the utterances with microphone failures, we excluded microphone signals whose minimum cross-correlation score among channels 4, 5, and 6 was less than 0.9, using the cross-correlation coefficients provided by the organizers [16]. The resultant evaluation set consisted of 1149 utterances (i.e., approximately 13% of the data were rejected).

As the evaluation metrics, we used the SDR of BSSEval [11] and the word error rate (WER). To evaluate the virtual microphone estimation performance, we computed the SDR between the estimated virtual microphone signals and the observed real microphone signals at the channel corresponding to the virtual microphone. To evaluate the enhancement performance of the beamformer, we used the clean reverberant signals at the fourth channel as reference. Since we required access to clean signals, this evaluation could only be performed on the simulated data. We used Kaldi’s CHiME-4 recipe [17, 18] to evaluate the ASR performance, which consisted of a deep neural network-hidden Markov model hybrid acoustic model [19, 20] trained with the lattice-free maximum mutual information criterion [21]. We used a trigram language model for decoding.

5.2. Experimental configurations

We adopted the Conv-TasNet-based network architecture for the network configuration of our proposed NN-VME. By following the notations of [5], we set the hyperparameters as follows: $N = 256$, $L = 20$, $B = 256$, $H = 512$, $P = 3$, $X = 8$, and $R = 4$.

We trained the proposed NN-VME using both simulated and real data of the training set, by adopting the Adam algorithm [22] with an initial learning rate of 0.0001 and the gradient clipping [23]. We stopped the training procedure after 200 epochs.

For the MVDR beamformer, we used the trained mask estimation model [10] provided in the GitHub repository [24], which was used in the Kaldi’s CHiME-4 recipe. For the STFT computation,

Table 1. SDR [dB] for virtual microphone estimator, in which noisy observed signal is used as reference signal

mic type	eval ch	ref ch	simu	real
RM	4	5	12.1	8.8
VM	5 (4,6)	5	16.6	13.8
VM	5 (3,4,6)	5	16.5	13.8
RM	5	6	8.3	7.8
VM	6 (4,5)	6	12.3	11.8
VM	6 (3,4,5)	6	12.5	11.9

we used a Blackman window with a length and shift set at 64 ms and 16 ms, respectively. In the ASR experiment, we set the loading hyperparameter ϵ of Eq. (5) to 0.05.

5.3. Experimental results

5.3.1. Evaluation of virtual microphone estimation performance

Table 1 shows the SDR scores of the evaluated NN-VME on the simulated and real recordings, where *RM* refers to real microphone, while *VM* refers to virtual microphone estimated by the proposed NN-VME. Here, note that the reference signal for the SDR computation is the noisy observation signal at the channel corresponding to the virtual microphone rather than the clean signal, and thus the virtual microphone estimation performance can be evaluated even for the real recordings.

In the table, the first column (“eval ch”) shows the channel index of the virtual or real microphone signal used as the estimated signal in the SDR calculation. The second column (“ref ch”) shows the channel index of the real microphone signal used as the reference signal. Here, the notation “5 (4,6)” indicates that the virtual microphone signal at channel 5 was estimated using the real microphone signals at channels 4 and 6. As a baseline, we compare the scores with the SDR obtained with the closest real microphones (i.e., the closest in terms of SDR). These results are shown in the first (eval ch 4, ref ch 5) and fourth rows (eval ch 5, ref ch 6) of Table 1.

Table 1 shows that the estimated signals by the proposed NN-VME modules (e.g., “5 (4,6)”) achieved much higher SDR scores than the observed signal recorded at the close microphone (e.g., “4”). These results demonstrate that, even for real recordings, the proposed NN-VME, i.e., the supervised learning-based virtual microphone estimation framework, has a potential to estimate the virtual microphone signals, which are not really observed by a microphone, by utilizing the spatial information inferred from a few observed real microphone signals.

The table shows results for interpolation, i.e., the virtual microphone is located between the real microphones (e.g., “5 (4,6)”), and extrapolation in the horizontal direction (e.g., “6 (4,5)”). In both cases, the NN-VME could predict the virtual microphone with an SDR of more than about 12 dB. In addition, we observed that increasing the number of the input channels (e.g., “5 (3,4,6)”) did not contribute to a significant improvement of the virtual microphone estimation performance. Investigating the behavior of the proposed method for various microphone array configurations is an interesting direction for future work.

5.3.2. Evaluation of beamformer enhancement performance

Table 2 shows the SDR scores of the evaluated beamformers on the simulated data. Here, *VM BF* refers to the beamformer with the es-

Table 2. SDR [dB] (higher is better) and WER [%] (lower is better) for beamformer, in which clean signal is used as reference signal

Method	used ch		SDR (simu)	WER (real)
	real	virtual		
(1) no process	-	-	8.6	15.8
(2) RM BF	4,6	-	10.8	12.0
(3) RM BF	4,5,6	-	14.2	9.4
(4) VM BF	4,6	5	13.4	11.1
(5) RM BF	3,4,6	-	12.7	10.0
(6) RM BF	3,4,5,6	-	15.2	8.5
(7) VM BF	3,4,6	5	14.2	9.5

timated virtual microphone, while *RM BF* refers to the beamformer only with the real microphone. In the table, the columns “real” and “virtual” in “used ch” denote the channel indices corresponding to the real and virtual microphones, which are used for constructing the beamformer, respectively. For example, the “VM BF” of row (4) is constructed using the two real microphone signals (i.e., channels 4 and 6) and one virtual microphone signal (i.e., channel 5).

Table 2 shows that the proposed VM BF (e.g., row (4)) successfully achieved higher SDR score compared to the RM BF constructed with the same real microphone signals (e.g., row (2)). Here, another RM BF (e.g., row (3)) would correspond to the upper-bound performance of the VM BF.

In addition to the above SDR-based evaluation, we conducted an ASR evaluation to evaluate the beamformer’s performance on the real recordings. Table 2 also shows the WER of the evaluated RM and VM BFs on the real data. From the table, we confirmed that, even for the real recordings, the proposed VM BF (e.g., row (4)) reduced the WER compared to its corresponding RM BF (e.g., row (2)) by up to 0.9 %. A similar trend is observed when using more microphones (i.e., rows (5)-(7)).

These results demonstrate that the estimated virtual microphone signals can contribute to improved enhancement performance when combined with a beamformer.

In the table, the results of the VM BF with the virtual microphone loading, as described in Section 3.3, are reported. The WER scores of the VM BF without the loading are 15.0 % for row (4) and 13.4 % for row (7), respectively. This confirms the effectiveness of the virtual microphone loading technique to improve the ASR performance for the VM BFs.

6. CONCLUSION

This paper proposed a novel virtual microphone estimation scheme, i.e., NN-VME, that employs a time-domain neural network architecture to directly predict the waveform of a virtual microphone signal. The proposed method relies on supervised learning framework and the high modeling capability of the network to build a model that can accurately predict virtual microphone signals based on a few real microphone observations.

We showed experimentally that the proposed NN-VME could precisely predict virtual microphone signals even for real recordings, and that these predicted signals could help boost beamforming performance.

Future work will include, evaluation of the effectiveness of the proposed method for 1) various microphone array configurations, 2) various acoustic conditions, and 3) other array processing techniques, e.g., source separation [25] and source localization [26].

7. REFERENCES

- [1] Jacob Benesty, Jingdong Chen, and Yiteng Huang, *Microphone array signal processing*, vol. 1, Springer Science & Business Media, 2008.
- [2] Hiroki Katahira, Nobutaka Ono, Shigeki Miyabe, Takeshi Yamada, and Shoji Makino, “Nonlinear speech enhancement by virtual increase of channels and maximum SNR beamformer,” *EURASIP Journal on Advances in Signal Processing*, vol. 2016, no. 1, pp. 1–8, 2016.
- [3] Ozgur Yilmaz and Scott Rickard, “Blind separation of speech mixtures via time–frequency masking,” *IEEE Transactions on Signal Processing*, vol. 52, no. 7, pp. 1830–1847, 2004.
- [4] Yi Luo and Nima Mesgarani, “TasNet: Time-domain audio separation network for real-time, single-channel speech separation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018, pp. 696–700.
- [5] Yi Luo and Nima Mesgarani, “Conv-TasNet: Surpassing ideal time–frequency magnitude masking for speech separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 27, no. 8, pp. 1256–1266, 2019.
- [6] Tsubasa Ochiai, Marc Delcroix, Rintaro Ikeshita, Keisuke Kinoshita, Tomohiro Nakatani, and Shoko Araki, “Beam-TasNet: Time-domain audio separation network meets frequency-domain beamformer,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 6384–6388.
- [7] Jon Barker, Ricard Marxer, Emmanuel Vincent, and Shinji Watanabe, “The third ‘CHiME’ speech separation and recognition challenge: Dataset, task and baselines,” in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2015, pp. 504–511.
- [8] Jonathan Le Roux, Scott Wisdom, Hakan Erdogan, and John R Hershey, “SDR – half-baked or well done?,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2019, pp. 626–630.
- [9] Mehrez Souden, Jacob Benesty, and Sofiene Affes, “On optimal frequency-domain multichannel linear filtering for noise reduction,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 2, pp. 260–276, 2009.
- [10] Jahn Heymann, Lukas Drude, and Reinhold Haeb-Umbach, “Neural network based spectral mask estimation for acoustic beamforming,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2016, pp. 196–200.
- [11] Emmanuel Vincent, Remi Gribonval, and Cedric Fevotte, “Performance measurement in blind audio source separation,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [12] Kouei Yamaoka, Shoji Makino, Nobutaka Ono, and Takeshi Yamada, “Performance evaluation of nonlinear speech enhancement based on virtual increase of channels in reverberant environments,” in *European Signal Processing Conference (EUSIPCO)*, 2017, pp. 2324–2328.
- [13] Ryoga Jinzai, Kouei Yamaoka, Mitsuo Matsumoto, Takeshi Yamada, and Shoji Makino, “Microphone position realignment by extrapolation of virtual microphone,” in *Asia-Pacific Signal and Information Processing Association Annual Summit and Conference (APSIPA ASC)*, 2018, pp. 367–372.
- [14] Kouei Yamaoka, Li Li, Nobutaka Ono, Shoji Makino, and Takeshi Yamada, “CNN-based virtual microphone signal estimation for MPDR beamforming in underdetermined situations,” in *European Signal Processing Conference (EUSIPCO)*, 2019, pp. 1–5.
- [15] “http://spandh.dcs.shef.ac.uk/chime_challenge/CHiME4/overview.html,”.
- [16] “http://spandh.dcs.shef.ac.uk/chime_challenge/CHiME4/data.html,”.
- [17] Daniel Povey, Arnab Ghoshal, Gilles Boulianne, Lukas Burget, Ondrej Glembek, Nagendra Goel, Mirko Hannemann, Petr Motlicek, Yanmin Qian, Petr Schwarz, Jan Silovsky, Georg Stemmer, and Karel Vesely, “The Kaldi speech recognition toolkit,” in *IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*, 2011.
- [18] “https://github.com/kaldi-asr/kaldi/tree/master/egs/chime4/s5_6ch,”.
- [19] Herve Boudlard and Nelson Morgan, “Connectionist speech recognition: A hybrid approach,” 1994.
- [20] Geoffrey Hinton, Li Deng, Dong Yu, George E Dahl, Abdelrahman Mohamed, Navdeep Jaitly, Andrew Senior, Vincent Vanhoucke, Patrick Nguyen, Tara N Sainath, and Brian Kingsbury, “Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups,” *IEEE Signal Processing Magazine*, vol. 29, no. 6, pp. 82–97, 2012.
- [21] Daniel Povey, Vijayaditya Peddinti, Daniel Galvez, Pegah Ghahremani, Vimal Manohar, Xingyu Na, Yiming Wang, and Sanjeev Khudanpur, “Purely sequence-trained neural networks for ASR based on lattice-free MMI,” in *Interspeech*, 2016, pp. 2751–2755.
- [22] Diederik P Kingma and Jimmy Ba, “Adam: A method for stochastic optimization,” in *International Conference on Learning Representations (ICLR)*, 2015.
- [23] Razvan Pascanu, Tomas Mikolov, and Yoshua Bengio, “On the difficulty of training recurrent neural networks,” in *International Conference on Machine Learning (ICML)*, 2013, pp. 1310–1318.
- [24] “<https://github.com/fgnt/nn-gev>,”.
- [25] Aapo Hyvarinen and Erkki Oja, “Independent component analysis: Algorithms and applications,” *Neural networks*, vol. 13, no. 4-5, pp. 411–430, 2000.
- [26] Ralph Schmidt, “Multiple emitter location and signal parameter estimation,” *IEEE Transactions on Antennas and Propagation*, vol. 34, no. 3, pp. 276–280, 1986.