# Modified Parametric Multichannel Wiener Filter for Low-latency Enhancement of Speech Mixtures with Unknown Number of Speakers

Ning Guo[1,*,†], Tomohiro Nakatani*, Shoko Araki*, and Takehiro Moriya*

* NTT Corporation, Japan

E-mail: tnak@ieee.org Tel: +81-774935134

† International Audio Laboratories Erlangen, Germany

E-mail: ning.guo@audiolabs-erlangen.de

*Abstract*—This paper introduces a novel low-latency online beamforming (BF) algorithm, named Modified Parametric Multichannel Wiener Filter (Mod-PMWF), for enhancing speech mixtures with unknown and varying number of speakers. Although conventional BFs such as linearly constrained minimum variance BF (LCMV BF) can enhance a speech mixture, they typically require such attributes of the speech mixture as the number of speakers and the acoustic transfer functions (ATFs) from the speakers to the microphones. When the mixture attributes are unavailable, estimating them by low-latency processing is challenging, hindering the application of the BFs to the problem. In this paper, we overcome this problem by modifying a conventional Parametric Multichannel Wiener Filter (PMWF). The proposed Mod-PMWF can adaptively form a directivity pattern that enhances all the speakers in the mixture without explicitly estimating these attributes. Our experiments will show the proposed BF's effectiveness in interference reduction ratios and subjective listening tests.

## I. INTRODUCTION

Microphone array processing has grown in importance due to the increasing need for remote communication. Speech signals captured by microphones often suffer from degradation caused by background noise and reverberation. Microphone array processing can improve the quality of speech signals by reducing noise and reverberation, thus enhancing speech perception and the performance of applications such as automatic speech recognition (ASR).

Researchers have successfully applied adaptive beamformers (BFs) [1]–[3] to enhance a single speech signal in the captured signal by low-latency processing. Typically, they use the speech's direction-of-arrival (DOA) to estimate the acoustic transfer function (ATF) from the speaker to the microphones (the steering vector) based on a plane-wave assumption and use the estimated ATF to optimize the BF [3], [4]. A method based on generalized eigenvalue decomposition (GEV) has been developed to estimate the ATF accurately in reverberation (i.e., multipath environments) using the spatial covariance matrices (SCMs) of the signals [5], [6]. Then, a Parametric Multichannel Wiener Filter (PMWF) has been proposed; It allows us to optimize a BF directly from the SCMs without estimating the ATF [7], [8]. Furthermore, it can flexibly control the speech distortion and noise reduction tradeoff.

In contrast, enhancing a mixture of speech signals, such as a conversation, with low-latency beamforming is still challenging. This is because the mixture may contain an unknown and varying number of speakers, where their ATFs are unknown in general, and two or more speakers may talk simultaneously. In such a scenario, a BF needs to adapt its directivity pattern quickly and accurately to the unknown and varying number of speaker locations. If the number of speakers and the ATFs are given, a Linearly Constrained Minimum Variance (LCMV) BF can form a directivity pattern that enhances all the speakers in a mixture [1], [6], [9]–[11]. Similarly, a blind source separation (BSS) technique, e.g., Independent Vector Extraction (IVE) [12]–[14], can enhance a mixture by first extracting each speaker and then remixing them, when the number of speakers is known. However, the requirements by these methods on the number of speakers and their ATFs are not acceptable in the assumed scenario. Alternatively, an Unknown Reference Multichannel Wiener Filter (UR-MWF) [15] can enhance a mixture without knowing the number of speakers or the ATFs. However, our experiments show that it has limited denoising and dereverberation performance.

To enable low-latency speech mixture enhancement, we propose to use a modified version of PMWF, referred to as Mod-PMWF. A Mod-PMWF is defined by replacing the spatial covariance matrix (SCM) of a single target source in a conventional PMWF with the SCM of the mixture of the desired signals. Our mathematical analysis reveals that a Mod-PMWF has several desirable characteristics for low-latency speech mixture enhancement:

1) A Mod-PMWF realizes a weighted sum of Minimum-Variance Distortionless Response (MVDR) BFs, each enhancing each source in the mixture. It adapts the weights of the MVDR BFs quickly to emphasize sources with larger powers at each time-frequency (TF) point. Thus, a Mod-PMWF forms a directivity pattern that enhances active sources at each TF point in the mixture.
2) We can reliably estimate a Mod-PMWF without knowing

the number of sources in the mixture using a practical approximation technique. In addition, similar to a conventional PMWF, we do not need to know the individual sources' ATFs for the estimation.

3) A Mod-PMWF has a low computational time complexity, comparable to conventional adaptive BFs such as a PMWF.

In experiments, we pick up three conventional BFs, GEV-based MVDR (GEV-MVDR) BF [6], Independent Single Component Extraction (ISCE) [13], and UR-MWF as ones that we think may perform mixture enhancement in certain senses without knowing the number of sources or their ATFs. We also pick up Independent Multiple Component Extraction (IMCE) [12] as one that can perform mixture enhancement with the prior knowledge of the number of sources. Comparison of our proposed Mod-PMWF with these conventional BFs using noisy reverberant speech mixtures confirms that the Mod-PMWF outperforms the other BFs in terms of interference reduction ratios and MUltiple Stimuli with Hidden Reference and Anchor (MUSHRA) subjective listening tests [16].

In the following, we introduce the signal model and present the conventional PMWF in Sections II and III. Section IV presents the proposed Mod-PMWF, and Section V describes other conventional BFs that we think may be used for mixture enhancement. Sections VI provide the experiment setup and results. Conclusions are drawn in Section VII.

## II. Signal Model and Enhancement Goal

Suppose $N$ speech signals are captured by $M$ microphones with reverberation and diffuse noise. Let $x_{m,t,f}$ be the captured signal at the $m$th microphone and a TF point $(t, f)$ in the short-time Fourier transform domain for $1 \leq t \leq T$ and $1 \leq f \leq F$, where $T$ and $F$ are the numbers of time frames and frequency bins, and let $(\cdot)^{\top}$ denote a non-conjugate transpose. Then the captured signal at all the microphones, $\mathbf{x}_{t,f} = [x_{1,t,f}, \ldots, x_{M,t,f}]^{\top} \in \mathbb{C}^{M}$, is modeled:

$$\mathbf{x}_{t,f} = \mathbf{d}_{t,f} + \mathbf{v}_{t,f}, \tag{1}$$

$$\mathbf{d}_{t,f} = \sum_{n=1}^{N} \mathbf{d}_{t,f}^{(n)} = \sum_{n=1}^{N} \mathbf{h}_{f}^{(n)} s_{t,f}^{(n)}, \tag{2}$$

where $\mathbf{d}_{t,f} \in \mathbb{C}^{M}$ is a mixture of speech signals $\mathbf{d}_{t,f}^{(n)}$ for all $n$. We assume that each speech signal, $\mathbf{d}_{t,f}^{(n)}$, comprises the direct signal plus the early reflections of the $n$th source, and can be modeled by a product of a time-invariant ATF $\mathbf{h}_{f}^{(n)} \in \mathbb{C}^{M}$ and the $n$th clean source signal $s_{t,f}^{(n)} \in \mathbb{C}$. $\mathbf{v}_{t,f}$ is the sum of all other signals, comprising the late reverberation of all sources and the additive diffuse noise. In this paper, we deal with $\mathbf{d}_{t,f}$ as the desired signal to be obtained and $\mathbf{v}_{t,f}$ as the interference signal to be reduced. In addition, we assume that $\mathbf{d}_{t}$ for all $n$ and $\mathbf{v}_{t}$ are mutually uncorrelated.

This paper assumes that both late reverberation and additive diffuse noise share the same spatial characteristics. We model their sum collectively as an interference signal $\mathbf{v}_{t,f}$ that follows

stationary complex Gaussian distribution with a mean vector $\mathbf{0}$ and an SCM $\Phi_{\mathbf{v},f} \in \mathbb{C}^{M \times M}$:

$$p(\mathbf{v}_{t,f}) = \mathcal{N}(\mathbf{v}; \mathbf{0}, \Phi_{\mathbf{v},f}) \quad \text{where} \quad \Phi_{\mathbf{v},f} = E\{\mathbf{v}_{t,f} \mathbf{v}_{t,f}^{\mathsf{H}}\}, \tag{3}$$

$E\{\cdot\}$ is an expectation function, and $(\cdot)^{\mathsf{H}}$ is a conjugate transpose.

Our goal in this paper is to estimate a BF $\mathbf{W}_{t,f} \in \mathbb{C}^{M \times M}$ at each TF point that keeps the speech mixture $\mathbf{d}_{t,f}$ unchanged while minimizing $\mathbf{v}_{t,f}$. With the BF, we obtain an enhanced speech mixture, $\mathbf{y}_{t,f} (\approx \mathbf{d}_{t,f}) \in \mathbb{C}^{M}$:

$$\mathbf{y}_{t,f} = \mathbf{W}_{t,f}^{\mathsf{H}} \mathbf{x}_{t,f}. \tag{4}$$

## III. Conventional BF for single source extraction

This section gives a brief overview of a conventional BF, PMWF, for single source extraction. It will be modified in the next section to the Mod-PMWF for mixture enhancement. Hereafter, we omit the index $f$ in symbols as we apply the same processing separately in each frequency.

### A. PMWF

When the number of sources in the observed signal is one (i.e., $N = 1$), a PMWF $\mathbf{w}_{m,t} \in \mathbb{C}^{M}$ [7] that enhances $\mathbf{d}_{t}^{(1)}$ ($= \mathbf{h}^{(1)} s_{t}^{(1)}$) in Eq. (2) at the $m$th microphone is defined as one that minimizes the interference while constraining the speech distortion not exceeding an allowable level $\sigma$:

$$\mathbf{w}_{m,t} = \arg\min_{\mathbf{w}} E\{|\mathbf{w}^{\mathsf{H}} \mathbf{v}_{t}|^{2}\} \quad \text{s.t.} \quad E\{|\mathbf{w}^{\mathsf{H}} \mathbf{d}_{t}^{(1)} - d_{m,t}^{(1)}|^{2}\} < \sigma.$$

Letting $\Phi_{\mathbf{d},t}^{(1)} = E\{\mathbf{d}_{t}^{(1)} (\mathbf{d}_{t}^{(1)})^{\mathsf{H}}\}$ be the SCM of $\mathbf{d}_{t}^{(1)}$ at $t$, the solution $\mathbf{W}_{t} = [\mathbf{w}_{1,t}, \ldots, \mathbf{w}_{M,t}]$ at all microphones is obtained:

$$\mathbf{W}_{t} = \left(\Phi_{\mathbf{d},t}^{(1)} + \gamma \Phi_{\mathbf{v}}\right)^{-1} \Phi_{\mathbf{d},t}^{(1)}. \tag{5}$$

Here $\gamma (\geq 0)$ is a weight controlling the tradeoff between the noise reduction and the speech distortion. This solution is also known as the Speech Distortion-Weighted MWF (SDW-MWF) [17]. Then, considering that $\Phi_{\mathbf{d},t}^{(1)}$ in Eq. (5) is rank-1 because it can be rewritten as $\Phi_{\mathbf{d},t}^{(1)} = E\{|s_{t}^{(1)}|^{2}\} \mathbf{h}^{(1)} (\mathbf{h}^{(1)})^{\mathsf{H}}$, we obtain a PMWF [7]:

$$\mathbf{W}_{t}^{\mathrm{PMWF}} = \frac{\Phi_{\mathbf{v}}^{-1} \Phi_{\mathbf{d},t}^{(1)}}{\gamma + \lambda_{\mathbf{d},t}} \quad \text{where} \quad \lambda_{\mathbf{d},t} = \mathrm{tr}\{\Phi_{\mathbf{v}}^{-1} \Phi_{\mathbf{d},t}^{(1)}\}, \tag{6}$$

and $\mathrm{tr}\{\Phi\}$ is a matrix trace. Here, we can use $\gamma$ as a parameter that controls the noise reduction level. For example, a PMWF with $\gamma = 0$ becomes an MVDR BF:

$$\mathbf{W}_{t}^{\mathrm{MVDR}} = \frac{\Phi_{\mathbf{v}}^{-1} \Phi_{\mathbf{d},t}^{(1)}}{\lambda_{\mathbf{d},t}}. \tag{7}$$

It has also been shown that a PMWF with $\gamma = 1$ becomes a MWF.

An important advantage of a PMWF is that we can estimate it without estimating the ATF from the source to the microphones, i.e., $\mathbf{h}^{(1)}$, once we obtain $\Phi_{\mathbf{v}}$ and $\Phi_{\mathbf{d},t}^{(1)}$. Techniques to estimate $\Phi_{\mathbf{v}}$ and $\Phi_{\mathbf{d},t}^{(1)}$ have been proposed using voice activity detection [7] and TF masks estimated by a neural network [8].

## IV. Proposed BF for Mixture Enhancement

This section presents the modified version of PMWF (Mod-PMWF) for mixture enhancement. We first present the definition of the Mod-PMWF, and then analyze its characteristics focusing on how it can perform mixture enhancement. Then, we describe its implementation using a practical approximation technique.

### A. Definition of Mod-PMWF

For defining the Mod-PMWF, we replace the SCM of the target signal of a PMWF, i.e., $\Phi_{\mathbf{d},t}^{(1)}$ in Eq. (6), with the SCM of the desired mixture, i.e., $\Phi_{\mathbf{d},t} = E\{\mathbf{d}_t \mathbf{d}_t^{\mathsf{H}}\}$, at each time frame. The formula of the Mod-PMWF is:

$$\mathbf{W}_t^{\text{Mod-PMWF}} = \frac{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{d},t}}{\gamma + \lambda_{\mathbf{d},t}} \quad \text{where} \quad \lambda_{\mathbf{d},t} = \text{tr}\{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{d},t}\}. \tag{8}$$

This BF is different from a PMWF in that $\Phi_{\mathbf{d},t}$ contains characteristics of not only a single target source but all the sources included in the mixture.

In the following, we explain how a Mod-PMWF can perform mixture enhancement. Then, we present how we can reliably estimate a Mod-PMWF using an approximation in Sections IV-B and IV-C.

*1) Analysis of Mod-PMWF:* Based on the assumption that $\mathbf{d}_t^{(n)}$ for all $n$ are mutually uncorrelated, the SCM $\Phi_{\mathbf{d},t}$ can be expanded:

$$\Phi_{\mathbf{d},t} = \sum_{n=1}^{N} \Phi_{\mathbf{d},t}^{(n)}, \tag{9}$$

where $\Phi_{\mathbf{d},t}^{(n)} = E\{\mathbf{d}_t^{(n)}(\mathbf{d}_t^{(n)})^{\mathsf{H}}\}$ is an unknown SCM of $\mathbf{d}_t^{(n)}$. Substituting Eq. (9) into Eq. (8) followed by a certain mathematical manipulation yields:

$$\mathbf{W}_t^{\text{Mod-PMWF}} = \sum_{n=1}^{N} \mu_{n,t}\mathbf{W}_t^{\text{MVDR}(n)}. \tag{10}$$

where $\mathbf{W}_t^{\text{MVDR}(n)}$ is an MVDR BF that enhances the $n$th source $\mathbf{d}_t^{(n)}$ by reducing the interference $\mathbf{v}_t$. It is defined based on a PMWF with $\gamma = 0$ in Eq. (7) as

$$\mathbf{W}_t^{\text{MVDR}(n)} = \frac{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{d},t}^{(n)}}{\lambda_{\mathbf{d},t}^{(n)}} \quad \text{where} \quad \lambda_{\mathbf{d},t}^{(n)} = \text{tr}\{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{d},t}^{(n)}\}. \tag{11}$$

$\mu_{n,t}$ in Eq. (10) is a time-varying weight defined:

$$\mu_{n,t} = \frac{\lambda_{\mathbf{d},t}^{(n)}}{\gamma + \lambda_{\mathbf{d},t}} \quad \text{where} \quad \lambda_{\mathbf{d},t} = \sum_{n=1}^{N} \lambda_{\mathbf{d},t}^{(n)}. \tag{12}$$

According to Eq. (10), a Mod-PMWF is a weighted sum of the MVDR BFs, each of which is designed to enhance each source in the mixture. The weights are determined by $\lambda_{\mathbf{d},t}^{(n)}$, which are roughly proportional to the power of $\mathbf{d}_t^{(n)}$.

In the following, we look into more details of the characteristics of the BF for the cases with $\gamma = 0$ and $\gamma > 0$.

*2) Behavior with $\gamma = 0$:* When $\gamma = 0$, the Mod-PMWF has the following characteristics.

- The weights sum to one according to Eqs. (12), thus Eq. (10) becomes a weighted average of the MVDR BFs. Because all the MVDR BFs reduce $\mathbf{v}_t$, their weighted average, i.e., the Mod-PMWF, should also reduce $\mathbf{v}_t$.
- MVDR BFs composing the weighted average cover all the sources in $\mathbf{d}_t$ given $\Phi_{\mathbf{d},t}$, even without explicitly specifying the number of sources.
- Each MVDR BF follows the unknown ATF of each source included in $\mathbf{d}_t$. Thus, it can preserve its corresponding source without distortion.
- As the weights of the BFs are roughly proportional to the powers of the respective sources, the Mod-PMWF can quickly adapt the weights to focus on active speakers at each TF point, i.e., sources with larger powers.

In summary, each MVDR BF composing the Mod-PMWF preserves each source in the mixture while reducing the interference. The Mod-PMWF rapidly controls the weights of the MVDR BFs depending on the change in the relative powers of the sources. In this sense, the Mod-PMWF with $\gamma = 0$ can achieve mixture enhancement.

*3) Behavior with $\gamma > 0$:* With $\gamma > 0$, the Mod-PMWF is equivalent to multiplying the following factor $\eta_t$ to the output of the Mod-PMWF with $\gamma = 0$.

$$\eta_t = \frac{\lambda_{\mathbf{d},t}}{\gamma + \lambda_{\mathbf{d},t}}. \tag{13}$$

Based on analogy from the conventional PMWF, we expect that this factor works as a single channel Wiener filter that can further reduce the interference remaining in the output of the Mod-PMWF with $\gamma = 0$. Also, we confirmed such a behavior of the Mod-PMWF with $\gamma > 0$ based on our preliminary experiments. More thorough analysis on this behavior should be included in future work.

### B. Practical approximation for Mod-PMWF

Although accurate estimation of $\Phi_{\mathbf{d},t}$ is crucial for Mod-PMWF, it is challenging under general recording conditions. To avoid this difficulty, this paper proposes to approximate $\Phi_{\mathbf{d},t}$ by the SCM of the captured signal $\Phi_{\mathbf{x},t} = E\{\mathbf{x}_t \mathbf{x}_t^H\}$, disregarding $\Phi_{\mathbf{v}}$ in $\Phi_{\mathbf{x},t}$, as

$$\Phi_{\mathbf{d},t} \simeq \Phi_{\mathbf{x},t}(= \Phi_{\mathbf{d},t} + \Phi_{\mathbf{v}}). \tag{14}$$

Unlike $\Phi_{\mathbf{d},t}$, $\Phi_{\mathbf{x},t}$ can be obtained easily and accurately from given $\mathbf{x}_t$. In addition, this approximation does not significantly affect the performance of Mod-PMWF, as discussed below.

With this approximation, the Mod-PMWF in Eq. (8) is rewritten:

$$\mathbf{W}_t^{\text{Approx-Mod-PMWF}} = \frac{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{x},t}}{\gamma + \lambda_{\mathbf{x},t}} \quad \text{where} \quad \lambda_{\mathbf{x},t} = \text{tr}\{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{x},t}\}. \tag{15}$$

Similar to Eq. (10), the approximated Mod-PMWF can be expanded using $\Phi_{\mathbf{x},t} = \Phi_{\mathbf{d},t} + \Phi_{\mathbf{v}}$:

$$\mathbf{W}_t^{\text{Approx-Mod-PMWF}} = \mu'_{\mathbf{d},t}\mathbf{W}_t^{\text{Mod-PMWF}} + \mu'_{\mathbf{v},t}\frac{\mathbf{I}_M}{M}, \tag{16}$$

where $\mathbf{I}_M \in \mathbb{R}^{M \times M}$ is an identity matrix. $\mu'_{\mathbf{d},t}$ and $\mu'_{\mathbf{v},t}$ are time-varying weights roughly proportional to the power of $\mathbf{d}_t$ and $\mathbf{v}_t$, and defined as

$$\mu'_{\mathbf{d},t} = \frac{\gamma + \lambda_{\mathbf{d},t}}{\gamma + \lambda_{\mathbf{d},t} + \lambda_{\mathbf{v}}} \quad \text{and} \quad \mu'_{\mathbf{v},t} = \frac{\lambda_{\mathbf{v}}}{\gamma + \lambda_{\mathbf{d},t} + \lambda_{\mathbf{v}}}. \quad (17)$$

where $\lambda_{\mathbf{v}} = M (= \text{tr}\{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{v}}\})$.

Equation (16) is a weighted average of the Mod-PMWF (with no approximation) and a filter that reduces the gain of the captured signal by a factor of $M$. We note that $\mathbf{W}_t^{\text{Mod-PMWF}}$ in Eq. (16) is accurate because it is based on accurate $\Phi_{\mathbf{d},t}$, which is included in $\Phi_{\mathbf{x},t}$. This means that using the approximated Mod-PMWF in Eq. (15) is equivalent to applying an accurate Mod-PMWF and then adding the captured signal with a reduced level. Although this approximation slightly reduces the interference reduction effect of Mod-PMWF, it can precisely preserve the speech mixture. We use this approximation thoughout experiments in this paper.

### C. Implementation of Mod-PMWF

To estimate a Mod-PMWF with the approximation, we only need to estimate $\Phi_{\mathbf{x},t}$ and $\Phi_{\mathbf{v},t}$ and calculate Eq. (15) at each TF point. First, according to the convention of low-latency online adaptive BFs, we can estimate $\Phi_{\mathbf{x},t}$ by online processing as the time average of $\mathbf{x}_t\mathbf{x}_t^{\mathsf{H}}$ using a forgetting factor $\beta$ ($0 < \beta \leq 1$). It can be recursively updated:

$$\Phi_{\mathbf{x},t} = \beta\Phi_{\mathbf{x},t-1} + (1-\beta)\mathbf{x}_t\mathbf{x}_t^{\mathsf{H}}. \quad (18)$$

Note that active sources may differ at different time frames, and sources that are not included in the current frame $t$ may be included in the past captured signals in Eq. (18). Then, their influence remains in $\Phi_{\mathbf{x},t-1}$. To minimize such influence, it is desirable to use a small $\beta$.

As for $\Phi_{\mathbf{v},t}$, $\Phi_{\mathbf{v},t}$ can be estimated, e.g., from noise signals recorded separately in advance, or from the captured signal during speech absent periods. Also, we can use the same techniques proposed for conventional PMWF [7], [8].

## V. CONVENTIONAL BFs THAT MAY BE APPLICABLE TO MIXTURE ENHANCEMENT

This section describes several conventional BFs that we think may be used for mixture enhancement in certain senses. We also give a comparison of the computational time complexity of the BFs with our proposed Mod-PMWF. In the next section, we will compare the BFs with our proposed Mod-PMWF by experiments.

### A. GEV-MVDR BF

The first conventional BF is a GEV-MVDR BF. We first present it as a technique for single source enhancement, and then explain how we can use it for mixture enhancement.

When the number of sources in the observed signal is one (i.e., $N = 1$), an ATF at each time frame $\mathbf{h}^{(1)}$ for the source can be estimated based on GEV [6]:

$$\mathbf{h}_t^{(1)} = \Phi_{\mathbf{v}}\mathbf{u}_t \quad \text{where} \quad \mathbf{u}_t = \text{MaxEig}\{\Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{x},t}\}, \quad (19)$$

where $\text{MaxEig}\{\Phi\}$ extracts an eigenvector of $\Phi$ corresponding to the maximum eigenvalue. Considering $\Phi_{\mathbf{d},t}^{(1)} = E\{|s_t^{(1)}|^2\}\mathbf{h}_t^{(1)}(\mathbf{h}_t^{(1)})^{\mathsf{H}}$ in Eq. (7) and substituting Eq. (19) into Eq. (7), we obtain a GEV-MVDR BF:

$$\mathbf{W}_t^{\text{GEV-MVDR}} = \frac{\mathbf{u}_t\mathbf{u}_t^{\mathsf{H}}\Phi_{\mathbf{v}}}{\text{tr}\{\mathbf{u}_t\mathbf{u}_t^{\mathsf{H}}\Phi_{\mathbf{v}}\}}. \quad (20)$$

For low-latency processing, we can calculate $\Phi_{\mathbf{x},t}$ in Eq. (19) using Eq. (18). Then, we can obtain $\mathbf{u}_t$ in Eq. (20) in a computationally efficient way using a power method for calculating $\text{MaxEig}\{\cdot\}$ in Eq. (19) [18], [19]. This iterates the following update after first setting $\mathbf{u}_t = \mathbf{u}_{t-1}$ at each time frame:

$$\mathbf{u}_t \leftarrow \Phi_{\mathbf{v}}^{-1}\Phi_{\mathbf{x},t}\mathbf{u}_t / |\mathbf{u}_t|, \quad (21)$$

In general, the power method converges very quickly, and only one iteration was sufficient in our experiments.

*1) Application of GEV-MVDR BF to mixture enhancement:* As shown in [20], a GEV-MVDR BF can be viewed also as a realization of a Maximum Signal-to-Noise Ratio (MaxSNR) BF [2], [21], [22]. This property allows us to use it for mixture enhancement. We explain this in the following.

Regardless of whether the captured signal contains a single source or multiple sources, we can define a MaxSNR BF $\mathbf{w}_t^{\text{MaxSNR}} \in \mathbb{C}^M$ as one that maximizes the SNR of the signal, i.e., the power ratio of the desired signal $\mathbf{d}_t$ to the noise $\mathbf{v}_t$:

$$\mathbf{w}_t^{\text{MaxSNR}} = \arg\max_{\mathbf{w}} \frac{E\{|\mathbf{w}^{\mathsf{H}}\mathbf{d}_t|^2\}}{E\{|\mathbf{w}^{\mathsf{H}}\mathbf{v}_t|^2\}}, \quad (22)$$

$$= \arg\max_{\mathbf{w}} \frac{E\{|\mathbf{w}^{\mathsf{H}}\mathbf{x}_t|^2\}}{E\{|\mathbf{w}^{\mathsf{H}}\mathbf{v}_t|^2\}}, \quad (23)$$

where we used $E\{|\mathbf{w}^{\mathsf{H}}\mathbf{x}_t|^2\} = E\{|\mathbf{w}^{\mathsf{H}}\mathbf{d}_t|^2\} + E\{|\mathbf{w}^{\mathsf{H}}\mathbf{v}_t|^2\}$ to obtain the above second line. The solution to Eq. (23) is obtained based on GEV as $\mathbf{w}_t^{\text{MaxSNR}} = b^*\mathbf{u}_t$, where $b$ is an indefinite constant, $(\cdot)^*$ is a complex conjugate, and $\mathbf{u}_t$ is from Eq. (19). Let $\mathbf{b} \in \mathbb{C}^M$ be a vector composed of $M$ different realizations of $b$ and let $\mathbf{W}_t = \mathbf{u}_t\mathbf{b}^{\mathsf{H}}$ be a matrix realization of the MaxSNR BF according to $\mathbf{b}$. Then, the MaxSNR BF with $\mathbf{b}$ defined below is identical to the GEV-MVDR BF in Eq. (20).

$$\mathbf{W}_t^{\text{MaxSNR}} = \mathbf{u}_t\mathbf{b}^{\mathsf{H}} \quad \text{where} \quad \mathbf{b} = \frac{\Phi_{\mathbf{v}}\mathbf{u}_t}{\text{tr}\{\mathbf{u}_t\mathbf{u}_t^{\mathsf{H}}\Phi_{\mathbf{v}}\}}. \quad (24)$$

The above analysis indicates that a GEV-MVDR BF in Eq. (20) works as a MaxSNR BF even when the observed signal contains two or more sources. In other words, the GEV-MVDR BF may enhance the mixture of speech signals, i.e., $\mathbf{d}_t$ in Eq. (1), in the MaxSNR sense regardless of the number of sources.

### B. Two versions of Independent Component Extraction (ICE): ISCE and IMCE

Next, we pick up Independent Component Extraction (ICE) as one that is applicable to mixture enhancement. ICE is a variation of a popular BSS technique, IVE [12]–[14]. Even

without prior knowledge of ATFs of the sources, both IVE and ICE can extract a given number ($N \geq 1$) of sources based on an assumption that the sources are mutually independent. We here adopt ICE because our preliminary experiments showed the superiority of ICE over IVE for mixture enhancement. The difference between ICE and IVE is the source models they use. While IVE's model specifies the source's spectral characteristics over all frequencies, ICE's model specifies the characteristics separately at each frequency. We defined the ICE's model at each frequency by a complex Gaussian with a mean of 0 and a time-varying variance. The optimization algorithm for online ICE can be easily obtained from that of online IVE [14].

We present two different usages of ICE, ISCE and IMCE, for mixture enhancement in the following.

*1) Application of ISCE to mixture enhancement:* When we set the number of sources to extract as $N = 1$ for ICE, we call it Independent Single Component Extraction (ISCE). This paper uses ISCE for mixture enhancement by utilizing speech sparseness. With the speech sparseness, we assume that speakers are not simultaneously active at each TF point. This assumption is based on the spectral characteristics of speech and has proven effective for source separation purposes [23]–[25]. With this assumption, we can simplify the mixture enhancement task to a single source extraction task at each TF point. Thus, we may use ISCE for mixture enhancement by simply applying ISCE to a mixture with quickly adapting ISCE over time frames.

*2) Application of IMCE to mixture enhancement:* When we assume that the number of sources is given ($N \geq 2$) and use ICE to extract all the sources from a mixture, we call it Independent Multiple Component Extraction (IMCE). By remixing the extracted sources, we can obtain an enhanced mixture. Unlike the other methods, this method requires the prior knowledge of the number of sources. Thus, we use IMCE just for reference.

*C. UR-MWF*

We can also use a version of conventional MWF, called UR-MWF [15], for mixture enhancement. It is defined as a BF $\mathbf{W}_t^{\mathrm{MWF}}$ that gives the minimum mean square error (MMSE) estimate of the desired mixture $\mathbf{d}_t$:

$$\mathbf{W}_t^{\mathrm{MWF}} = \arg\min_{\mathbf{W}} E\{\|\mathbf{d}_t - \mathbf{W}^{\mathsf{H}}\mathbf{x}_t\|_2^2\}. \tag{25}$$

Assuming $\Phi_{\mathbf{x},t} = E\{\mathbf{d}_t\mathbf{d}_t^{\mathsf{H}}\} + \Phi_{\mathbf{v}}$, we obtain UR-MWF:

$$\mathbf{W}_t^{\mathrm{UR\text{-}MWF}} = \Phi_{\mathbf{x},t}^{-1}(\Phi_{\mathbf{x},t} - \Phi_{\mathbf{v}}). \tag{26}$$

Based on Eq. (26), UR-MWF can enhance a mixture given $\Phi_{\mathbf{x},t}$ and $\Phi_{\mathbf{v}}$ without knowing the number of sources or their ATFs.

*D. Comparison of computational time complexity*

Table I shows the computational time complexity of the BFs and the Mod-PMWF. All the BFs except for IMCE have the same complexity $O(M^2)$. The complexity increases to $O(NM^2)$ for IMCE to extract $N$ sources.

TABLE I: Computational time complexity for adapting and applying a BF to each TF point. The complexity of Mod-PMWF, GEV-MVDR, and UR-MWF is estimated assuming that we recursively update $\Phi_{\mathbf{v}}$ and its inverse, e.g., during the speech absent periods, as in [19].

| Mod-PMWF | GEV-MVDR | ISCE | IMCE | UR-MWF |
|----------|----------|------|------|--------|
| $O(M^2)$ | $O(M^2)$ | $O(M^2)$ | $O(NM^2)$ | $O(M^2)$ |

## VI. Experiments

We evaluated our proposed Mod-PMWF in comparison with the other conventional BFs using signal level metrics and a MUSHRA listening test.

*A. Dataset*

For the evaluation, we prepared a simulated noisy reverberant mixture dataset containing 100 samples of two speakers from the ATR digital speech database set B [26]. To reverberate each speaker's utterance, we used room impulse responses (RIRs) measured with a cubic array of eight microphones in a room with a T60 of 380 ms. Each edge of the cubic array was 4 cm, and a microphone was equipped at each vertex of the array. We used 5 out of 8 microphones in the experiments. The two speakers were located 1 m from the array on opposite sides. We then mixed two speech signals starting simultaneously with babble noise recorded by the same array in the same room. We created two versions of the dataset by setting the Reverberant speech Mixture to Noise Ratio (RMNR) to 20 dB and 10 dB. The sampling frequency was set at 16 kHz.

*B. Analysis condition*

Throughout the experiments, we set $\gamma = 0$ for Mod-PMWF to evaluate its MVDR behavior in this paper. The forgetting factor $\beta$ for updating $\Phi_{\mathbf{x},t}$ in Eq. (18) decides how fast the algorithm adapts to the change in the signals' spatial characteristics. We chose $\beta = 0.99$ for UR-MWF, 0.96 for ISCE and IMCE, and 0.85 for Mod-PMWF and GEV-MVDR based on our preliminary experiments. We set the window length and shift at 20 ms and 5 ms.

In Sections VI-C and VI-D, we assume that the interference SCMs $\Phi_{\mathbf{v}}$ can be estimated in advance. In concrete, we estimated them from a part of the measured babble noise not used for generating the mixtures. Then, we evaluated Mod-PMWF with no prior knowledge of $\Phi_{\mathbf{v}}$ in Section VI-E, where we estimated $\Phi_{\mathbf{v}}$ by online processing from observed signals [19] using estimated speech presence probability [27].

*C. Evaluation results using Signal Level Metrics*

We defined two signal level metrics for the evaluation [28], [29]: Segmental Desired to Interference Ratio (SegDIR) and Segmental Desired to Distortion Ratio (SegDDR). SegDIR was used for measuring the degree of interference reduction. SegDDR was used for measuring the degree of signal distortion introduced to the enhanced speech. Let $d_m(i)$ and $v_m(i)$ be the $i$th samples of the desired mixture and interference signal in
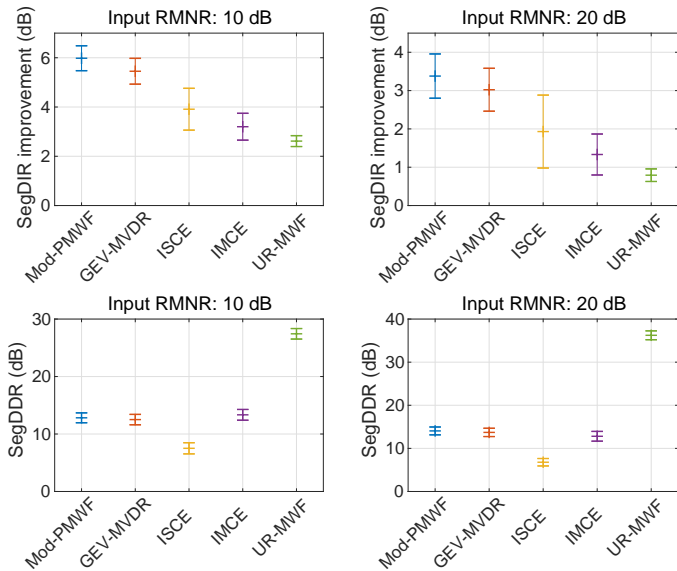
Fig. 1: SegDIR improvements and SegDDRs averaged over 100 samples. Error bars indicate standard deviation. Mod-MPDR, Gev-MVDR, and UR-MWF used interference SCMs separately estimated in advance.

the time domain at the $m$th microphone, and $\hat{d}_m(i)$ and $\hat{v}_m(i)$ be samples obtained by applying a BF to $d_m(i)$ and $v_m(i)$. To determine the desired signal $d_m(i)$, we set the duration of the early reflections at 50 ms. Then, the metrics are defined as:

$$\mathrm{SegDIR} = \frac{10}{|\mathcal{A}|} \sum_{\tau \in \mathcal{A}} \log_{10} \frac{\sum_{m=1}^{M} \sum_{i=\tau\delta+1}^{\tau\delta+\delta} |\hat{d}_m(i)|^2}{\sum_{m=1}^{M} \sum_{i=\tau\delta+1}^{\tau\delta+\delta} |\hat{v}_m(i)|^2},$$

$$\mathrm{SegDDR} = \frac{10}{|\mathcal{A}|} \sum_{\tau \in \mathcal{A}} \log_{10} \frac{\sum_{m=1}^{M} \sum_{i=\tau\delta+1}^{\tau\delta+\delta} |d_m(i)|^2}{\sum_{m=1}^{M} \sum_{i=\tau\delta+1}^{\tau\delta+\delta} |\hat{d}_m(i) - d_m(i)|^2},$$

where $\delta = 800$ samples ($= 50$ ms) is the segment length, $\tau$ is a segment index, $\mathcal{A}$ is a set of segments where at least one speech signal is active, and $|\mathcal{A}|$ is the number of segments in $\mathcal{A}$.

Figure 1 shows the evaluation results. In the figure, we show SegDIR as its relative improvement from the captured signal and SegDDR as its absolute value.

First, when we compare SegDIRs obtained by the BFs, the proposed Mod-PMWF was the best of all. In particular, Mod-PMWF significantly outperformed ISCE, IMCE, and UR-MWF. Next, when we compare SegDDRs obtained by the BFs, Mod-PMWF was better than ISCM, and comparable with GEV-MVDR and IMCE. In contrast, UR-MWF achieved significantly better SegDDR than all the others, however, its improvement of SegDIR was minimal.

We can draw the following implications from these results.
1) Although Mod-PMWF and GEV-MVDR did not use prior knowledge of the number of sources, they achieved comparable SegDDRs with IMCE that uses the knowledge of the number of sources. This suggests that Mod-PMWF and GEV-MVDR effectively enhanced the

mixture based on their unique mechanisms described in Sections IV-A1 and V-A1.
2) For improving SegDIR, the sparseness assumption used for ISCE was effective, to some extent, in comparison with IMCE, but insufficient in comparison with Mod-PMWF and GEV-MVDR. Also, the assumption introduced more speech distortion than the other BFs.
3) The inferior SegDIR obtained by IMCE is probably caused by remixing of two separated signals. The remixing sums up the interference that remained in each separated signal. In contrast, Mod-PMWF could effectively reduced the interference while enhancing the mixture. This is probably because Mod-PMWF adaptively controls the BF weights to enhance only active sources at each TF point.

In contrast, it is difficult to evaluate the effectiveness of UR-MWF in comparison with the other BFs based only on the above results.

To compare the proposed BF with the other BFs, including UR-MWF, more reliably, we conducted a subjective listening test.

*D. Evaluation Results with MUSHRA Listening Test*

We conducted a MUSHRA listening test [16] to evaluate the perceived speech quality. We picked up two speech mixtures composed of two male or two female speakers and tested them under each RMNR condition. Each sample had seven signals to rate, including the desired mixture as the (hidden) reference and the noisy reverberant observation as the anchor. Ten experienced listeners participated in the test. We asked them to concentrate on the overall quality of the audio including both the noise level and the degradation of speech.

Figure 2 shows the result. The two sub-figures show the results under input RMNRs of 20 dB and 10 dB and both sub-figures exhibit almost the same tendency. First, the participants rated the high-quality reference as 100 points and test conditions between 40-80 points, confirming our choice of MUSHRA to be appropriate according to the ITU-R recommendation BS.1534. In addition, improvements obtained by all the BFs were statistically significant from the anchor, i.e., the noisy observation rated with the lowest scores. Thus our anchor can indicate how the methods compare to known audio quality levels.

Next, when comparing results between BFs, the proposed Mod-PMWF significantly outperformed all the other BFs, including GEV-MVDR and UR-MWF. The result was better correlated with SegDIR than with SegDDR. One reason might be that the speech distortion is, to some extent, masked by the active speech and noise residual and thus not as annoying to the listeners as the noise residual.

*E. Evaluation with no prior knowledge of interference SCM*

Finally, we evaluated Mod-PMWF without using prior knowledge of the interference SCM. In this experiment, we recursively updated the interference SCM $\Phi_{\mathbf{v},t,f}$ at each TF
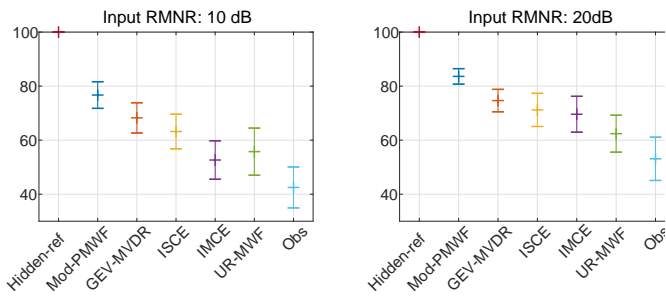
Fig. 2: MUSHRA test results. Error bars indicate 95% confidence based on t-distribution test. Mod-MPDR, Gev-MVDR and UR-MWF used interference SCMs separately estimated in advance.

point as

$$\Phi_{\mathbf{v},t,f} = \alpha' \Phi_{\mathbf{v},t-1,f} + (1-\alpha')\mathbf{x}_{t,f}\mathbf{x}_{t,f}^H, \quad (28)$$

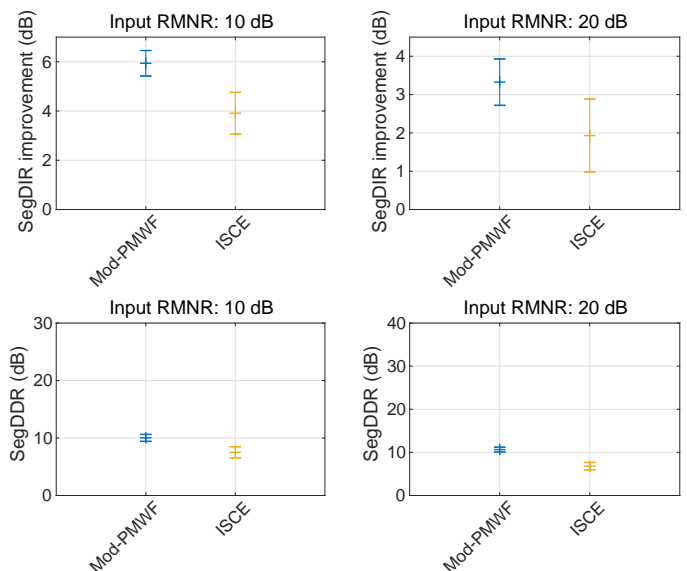$$\alpha' = \alpha + (1-\alpha)q_{t,f}, \quad (29)$$

where $\alpha = 0.9998$ is a forgetting factor and $q_{t,f}$ is speech presence probability estimated recursively by maximum likelihood estimation [27].

Figures 3 (a) and (b) show SegDIR improvements, SegDDRs, and MUSHRA test results. In this experiment, selected two methods, Mod-PMWF based on online interference SCM estimation and ISCE based on BSS, mainly elaborate reliability of the MUSHRA test. For the MUSHRA test, we also added a mixture sample composed of male and female speakers. Totally we had 3 samples for each RMNR condition. Eight experienced listeners participated in this test.
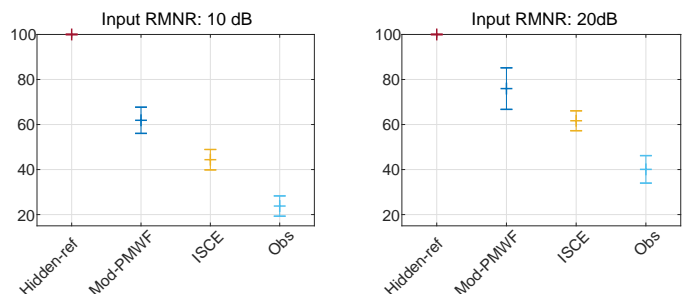
Mod-PMWF was significantly better than ISCE in terms of all the metrics. The overall tendency of the results in the figures was almost identical to that in Figs. 1 and 2 except the SegDDRs obtained by Mod-PMWF were slightly lower than Fig. 1. These results show that Mod-PMWF effectively performed mixture enhancement even when no prior knowledge on the interference SCMs was available.

## VII. CONCLUSION

This paper proposed a Mod-PMWF for low-latency online speech mixture enhancement in noisy reverberant environments. We showed mathematically that Mod-PMWF is equivalent to a weighted sum of MVDR BFs, where each MVDR BF can preserve each source included in the captured signal and the time-varying weights quickly adapt at each time frame to put larger weights on sources with larger powers. We can estimate the Mod-PMWF by low-latency online processing without relying on prior knowledge of the number of sources or the ATFs from the sources to microphones. We discussed that approximating the SCM of the desired mixture by the SCM of the captured signal for Mod-PMWF is advantageous to make it accurately preserve the speech mixture although it slightly reduces the interference reduction effect. We verified the effectiveness of the Mod-PMWF for mixture enhancement in terms of interference reduction ratios and subjective speech



(a) SegDIR improvements and SegDDRs averaged over 100 samples. Error bars indicate standard deviation.



(b) MUSHRA test results. Error bars indicate 95% confidence based on t-distribution test.

Fig. 3: Evaluation with no prior knowledge of interference SCM. The interference SCM was estimated by online processing for Mod-PMWF.

quality tests. The Mod-PMWF outperformed several conventional BFs, including GEV-MVDR BF that performs mixture enhancement in a MaxSNR sense, a BSS technique ISCE utilizing the speech sparseness, another BSS technique IMCE with prior knowledge of the number of sources, and UR-MWF that gives the MMSE estimates of the desired mixtures.

## REFERENCES

[1] B. D. V. Veen and K. M. Buckley, "Beamforming: A versatile approach to spatial filtering," *IEEE ASSP Magazine*, vol. 5, no. 2, pp. 4–24, 1988.

[2] H. L. V. Trees, *Optimum Array Processing, Part IV of Detection, Estimation, and Modulation Theory*. New York: Wiley-Interscience, 2002.

[3] S. Doclo, W. Kellermann, S. Makino, and S. Nordholm, "Multichannel signal ennhancement algorithms for assisted listening devices," *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 18–30, 2015.

[4] J. Xiao, Z.-Q. Luo, I. Merks, and T. Zhang, "A robust adaptive binaural beamformer for hearing devices," in *Proc. Asilomar Conference on Signals, Systems, and Computers*, 2017.

[5] S. Gannot, D. Burshtein, and E. Weinstein, "Signal enhancement using beamforming and non-stationarity with applications to speech," *IEEE Trans. Signal Processing*, vol. 49, no. 8, pp. 1614–1626, 2001.

[6] S. Markovich-Golan, S. Gannot, and I. Cohen, "Multi-channel eigenspace beamforming in a reverberant noisy environment with multiple interfering speech signals," *IEEE Trans. ASLP*, vol. 17, no. 6, pp. 1071–1086, 2009.

[7] M. Souden, J. Benesty, and S. Affes, "On optimal frequency-domain multichannel linear filtering for noise reduction," *IEEE Trans. Audio, Speech, and Language Processing*, vol. 18, no. 2, pp. 260–276, 2007.

[8] T. Ochiai, S. Watanabe, T. Hori, and J. R. Hershey, "Multichannel end-to-end speech recognition," in *Proc. International conference on machine learning*, 2017, pp. 2632–2641.

[9] M. Er and A. Cantoni, "Derivative constraints for broad-band element space antenna array processors," *IEEE Transactions on Acoustics, Speech, Signal Processing*, vol. 31, no. 6, pp. 1378–1393, 1983.

[10] O. Schwartz, S. Gannot, and E. A. P. Habets, "Multispeaker LCMV beamformer and postfilter for source separation and noise reduction," *IEEE/ACM Trans. Audio, Speech, and Language Processing*, vol. 25, no. 5, 2017.

[11] A. Herzog and E. A. P. Habets, "Direction and reverberation preserving noise reduction of ambisonics signals," *IEEE/ACM Trans. Audio, Speech, and Language Processing*, vol. 28, 2020.

[12] R. Scheibler and N. Ono, "Independent vector analysis with more microphones than sources," in *Proc. IEEE WASPAA*, 2019.

[13] R. Ikeshita, T. Nakatani, and S. Araki, "Block coordinate descent algorithms for auxiliary-function-based independent vector extraction," *IEEE Trans. Signal Processing*, vol. 69, pp. 3252–3267, 2021.

[14] T. Ueda, T. Nakatani, R. Ikeshita, K. Kinoshita, S. Araki, and S. Makino, "Low latency online source separation and noise reduction based on joint optimization with dereverberation," in *Proc. European Signal Processing Conference (EUSIPCO)*, 2021, pp. 1000–1004. DOI: 10.23919/EUSIPCO54536.2021.9616119.

[15] A. Spriet, M. Moonen, and J. Wouters, "Robustness analysis of multichannel Wiener filtering and generalized sidelobe cancellation for multimicrophone noise reduction in hearing aid applications," *IEEE Transactions on Speech and Audio Processing*, vol. 13, no. 4, pp. 487–503, Jul. 2005.

[16] *ITU-R recommendation BS.1534*.

[17] S. Doclo, A. Spriet, M. Moonen, and J. Wouters, "Frequency-domain criterion for the speech distortion weighted multichannel Wiener filter for robust noise reduction," *Speech Communication*, vol. 49, pp. 636–656, 2007.

[18] A. Krueger, E. Warsits, and R. Haeb-Umbach, "Speech enhancement with a GSC-like structure employing eigenvector-based transfer function ratio estimation," *IEEE Trans. Audio, Speech, and Language Processing*, vol. 19, no. 1, pp. 206–219, 2011.

[19] T. Nakatani and K. Kinoshita, "Simultaneous denoising and dereverberation for low-latency applications using frame-by-frame online unified convolutional beamformer," in *Proc. Interspeech 2019*, 2019, pp. 111–115. DOI: 10.21437/Interspeech.2019-1286.

[20] E. Warsitz and R. Haeb-Umbach, "Blind acoustic beamforming based on generalized eigenvalue decomposition," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 5, 2007.

[21] S. Araki, H. Sawada, and S. Makino, "Blind speech separation in a meeting situation with maximum SNR beamformer," in *Proc. IEEE ICASSP*, 2007, pp. 41–44.

[22] J. Heymann, L. Drude, and R. Haeb-Umbach, "Neural network based spectral mask estimation for acoustic beamforming," in *Proc. IEEE ICASSP*, 2016, pp. 196–200.

[23] O. Yilmaz and S. Rickard, "Blind separation of speech mixtures via time-frequency masking," *IEEE Transactions on Signal Processing*, vol. 52, no. 7, pp. 1830–1847, Jul. 2004.

[24] D. Wang, "On ideal binary mask as the computational goal of auditory scene analysis," in *Speech Separation by Humans and Machines*, 2005, pp. 181–197.

[25] M. Souden, S. Araki, K. Kinoshita, T. Nakatani, and H. Sawada, "A multichannel MMSE-based framework for speech source separation and noise reduction," *IEEE Trans. Audio, Speech, and Language Processing*, vol. 21, no. 9, pp. 1913–1928, 2013.

[26] A. Kurematsu, K. Takeda, Y. Sagisaka, S. Katagiri, H. Kuwabara, and K. Shikano, "ATR Japanese speech database as a tool of speech recognition and synthesis," *Speech communication*, vol. 9, no. 4, pp. 357–363, 1990.

[27] M. Souden, M. Delcroix, K. Kinoshita, T. Yoshioka, and T. Nakatani, "Noise power spectral density tracking: A maximum likelihood perspective," *IEEE Signal Processing Letters*, vol. 19, no. 8, pp. 495–498, 2012.

[28] J. H. L. Hansen and B. L. Pellom, "An effective quality evaluation protocol for speech enhancement algorithms," in *Proc. International conference on spoken language processing*, 1998, pp. 2819–2822.

[29] Y. Hu and P. C. Loizou, "Evaluation of objective quality measures for speech enhancement," *IEEE Trans. Audio, Speech, and Language Processing*, vol. 16, no. 1, pp. 229–238, 2008.