

On Semi-blind Source Separation-based Approaches to Nonlinear Echo Cancellation Based on Bilinear Alternating Optimization

Xianrui Wang, *Graduate Student Member, IEEE*, Yichen Yang, *Graduate Student Member, IEEE*, Andreas Brendel, *Member, IEEE*, Tetsuya Ueda, *Student Member, IEEE*, Shoji Makino, *Life Fellow, IEEE*, Jacob Benesty, Walter Kellermann, *Life Fellow, IEEE*, and Jingdong Chen, *Fellow, IEEE*

Abstract—Acoustic echo cancellation (AEC) is a crucial task in full duplex communications. As conventional linear filtering approaches are ineffective to deal with double-talk, various semi-blind source separation (SBSS)-based AEC algorithms are devised, most of which are formulated and implemented in the frequency domain based on the multiplicative transfer function (MTF) model for computational efficiency. To avoid large latency and in order to deal with loudspeaker nonlinearities, the convolutive transfer function (CTF) model and odd power series expansion are leveraged, which are employed by numerous SBSS-based nonlinear AEC (SBSS-NAEC) algorithms. Conventional SBSS-NAEC methods estimate the series expansion coefficients and the CTF filter simultaneously making the number of free parameters to estimate large. Hence, the corresponding algorithms are computationally expensive and are difficult to optimize. In this work, we propose to decouple the series expansion coefficients and the CTF filters into a bilinear form and present a bilinear alternating optimization framework for estimating the model parameters. An alternating iterative projection (AIP) algorithm and an alternating element-wise iterative source steering (AEISS) algorithm are proposed. As the bilinear representation consists of less parameters compared to the conventional methods, the proposed algorithms not only improve the AEC performance but also reduce the computational complexity, which is validated by comprehensive simulations and experiments.

Index Terms—Semi-blind source separation, nonlinear acoustic echo cancellation, odd power series expansion, convolutive transfer function model, bilinear, alternating optimization.

I. INTRODUCTION

Acoustic echoes, which are caused by coupling between the loudspeakers and microphones, are detrimental to full duplex voice communication [1]–[3]. Consequently, acoustic echo cancellation (AEC), a process to estimate and eliminate echoes, has to be used in full duplex voice communication

This work was supported in part by the National Key Research and Development Program of China under Grant No. 2021ZD0201502 and in part by the Major Program of National Science Foundation of China (NSFC) Grant No. 62192713.

Xianrui Wang and Yichen Yang are with the Center of Intelligent Acoustics and Immersive Communications, Northwestern Polytechnical University, China. They are also currently visiting PhD students at Waseda University, Japan. Andreas Brendel is with Fraunhofer IIS, Fraunhofer Institute for Integrated Circuits IIS in Erlangen, Germany. Tetsuya Ueda and Shoji Makino are with Waseda University, Japan. Jacob Benesty is with INRS-EMT, University of Quebec, Montreal, Canada. Walter Kellermann is with the Chair of Multimedia Communications and Signal Processing, University of Erlangen-Nuremberg, Germany. Jingdong Chen is with the Center of Intelligent Acoustics and Immersive Communications, Northwestern Polytechnical University, China.

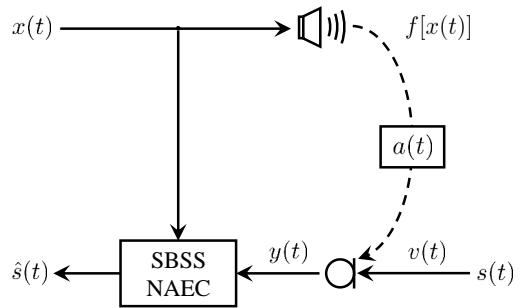


Fig. 1: Illustration of NAEC in full duplex communications.

systems [1]–[6]. Generally, AEC models the acoustic impulse response (AIR) from the loudspeaker to the microphone as a finite impulse response (FIR) filter. The problem is then transformed into one of FIR channel identification [3], [7]–[13]. Since AIRs are generally time-varying, adaptive filters have to be used and many adaptive algorithms have been developed over the last few decades, such as the least-mean-square (LMS) [14]–[17], normalized LMS (NLMS) [18]–[21], recursive least-squares (RLS) [22]–[25], and Kalman filters [26]–[28]. Those algorithms can achieve good AEC performance in the single-talk scenario, where there is no near-end speech; but their performance often degrades significantly in the presence of double-talk, where both the far- and near-end speech signals coexist. The traditional way to deal with this issue is through using a double-talk detector (DTD) [29], [30] and the adaptation process stops whenever double-talk is detected. While it has been proved to be a viable way to handle double-talk, this approach is unable to track the dynamic acoustic system when the adaptation is stalled and, as a result, often leads to severe performance degradation during double-talk. To address this issue, some Kalman filter-like algorithms were developed, which model both the near-end and echo signals as uncorrelated complex Gaussian distributed processes [31], [32]. While improvement is observed, the performance of such algorithms is still by far not satisfactory and do not meet the requirements of practical systems.

Alternatively, the double-talk issue can be addressed from the perspective of blind source separation (BSS) [33]–[35], and several semi-blind source separation-based AEC (SBSS-AEC) algorithms have been developed [36]–[39], which are

formulated in the frequency domain and can be implemented efficiently thanks to the fast Fourier transform (FFT). Early such algorithms adopted the so-called multiplicative transfer function (MTF) model in which the time-domain convolution is represented by the frequency-domain bin-wise multiplication [36], [39]. To achieve good performance, the analysis window for short-time Fourier transform (STFT) in this model has to be long enough to cover the effective part of the AIR, i.e., the components that significantly contribute to the echo. This, however, leads to large algorithmic latency, particularly when the system operates in large rooms or enclosures with high level of reverberation [40]. To circumvent this issue, the convolutive transfer function (CTF) model was developed, which represents the time-domain convolution with a frequency-domain convolution using signal components from multiple consecutive frames [41]–[44] and is structurally equivalent to frequency subband filtering [45], [46]. This framework offers flexibility to choose shorter analysis window, which enables to reduce the algorithmic latency. The CTF filter estimation can be formulated as a simplified hybrid exact-approximate diagonalization (HEAD) problem [47] based on the assumption that the far- and near-end signals are mutually independent, which has been well studied in the literature of BSS [48]–[55]. Iterative projection (IP) [48] is one of the most widely-used algorithms to solve this problem for SBSS-AEC [56], [57]. However, since it involves matrix inversion, the IP algorithm is computationally very expensive. With the inspiration of the iterative source steering (ISS) principle [53], a more computationally efficient algorithm was developed in [58], which is called the element-wise iterative source steering (EISS). As shown in [57], [58], both the IP and EISS algorithms can achieve significant performance improvement in comparison with the aforementioned Kalman filter for AEC in double-talk scenarios [31], [32].

Besides double-talk, another challenging issue to deal with in AEC is the loudspeaker nonlinearity, which makes the linear mixing model no longer appropriate [3], [5], i.e., the input of the linear mixing system is a nonlinear transform of the far-end signal. To deal with this problem, Volterra filters are often used [59]–[61]; but they require large memory and are computationally expensive and difficult to implement in most low-cost devices. Another simple yet efficient approach to model loudspeaker nonlinearities is through adding the so-called odd power series expansion [62]–[64] into the CTF framework, based on which several SBSS-based nonlinear AEC (SBSS-NAEC) algorithms have been proposed [58], [65]. These methods merge the CTF filter and the series expansion coefficients into a long vector, which represents the parameters to be identified [58], [65]. We will refer to this filter as the merged near-end signal extraction (MNE) filter in the rest of this work. The IP and EISS algorithms have been developed for estimating MNE filters, which have demonstrated promising NAEC performance [58], [65]. However, the computational complexity of such algorithms is quite high as the length of the MNE filter is generally large, which is proportional to the product of the CTF filter length and the number of series expansion coefficients.

In order to reduce the computational complexity and to

TABLE I: Notation of important acronyms.

NAEC	Nonlinear acoustic echo cancellation
AIR	Acoustic impulse response
SBSS	Semi-blind source separation
MTF	Multiplicative transfer function
CTF	Convolutive transfer function
MNE	Merged near-end extraction
HEAD	Hybrid exact-approximate diagonalization
MM	Majorize-minimization
LCQP	Linear constrained quadratic programming
IP	Iterative projection
AIP	Alternating iterative projection
EISS	Element-wise iterative source steering
AEISS	Alternating element-wise iterative source steering

improve the performance of IP and EISS algorithms, we propose to represent the nonlinear echo in a bilinear form [66], [67] in this work, which decouples the CTF filter and the series expansion coefficients. To the best of our knowledge, this is the first time that such a form is explored in the frequency-domain SBSS-NAEC, though a similar model was investigated for time-domain adaptive filtering [68]. We then present a bilinear alternating optimization framework for estimating the model parameters. Based on this framework and also following the principles in [57], [58], [65], we derive an alternating iterative projection (AIP) algorithm and an alternating element-wise iterative source steering (AEISS) algorithm. Since the number of free parameters to estimate is less than those in the conventional IP and EISS methods, the derived two algorithms cannot only improve the AEC performance but also reduce the computational complexity.

The contribution of this work are as follows. First, we leverage a bilinear form to decouple the CTF filter and the nonlinear expansion coefficients and reformulate the estimation of them as two SBSS problems. As the CTF filter and the nonlinear coefficients are intertwined, we proposed a bilinear alternating optimization, based on which two efficient algorithms, i.e., AIP and AEISS, are developed. Finally, simulations and experiments are carried out to validate the superior performance of AIP and AEISS. The notation of important acronyms is given in Table I for later reference.

II. SIGNAL MODEL

Consider a full-duplex communication scenario. The far-end signal is played back through a loudspeaker with some unknown nonlinearity. The loudspeaker signal is then convolved with the AIR as it travels from the loudspeaker to the near-end microphone, forming the so-called echo. Finally, the echo signal, along with the near-end signal, is captured by the near-end microphone and transmitted back to the far end. Mathematically, this process can be expressed at time t as

$$\begin{aligned} y(t) &= s(t) + v(t) \\ &= s(t) + a(t) \star f[x(t)], \end{aligned} \quad (1)$$

where $y(t)$ is the observed microphone signal, $s(t)$ is the near-end signal, $v(t) = a(t) \star f[x(t)]$ is the nonlinear acoustic echo, $a(t)$ stands for the AIR, \star denotes the linear convolution, and $f[\cdot]$ represents the loudspeaker responses, including both linear and nonlinear effects.

A. Odd Power Series Expansion

The Maclaurin series expansion is an efficient mathematical tool to approximate nonlinear functions. As memoryless nonlinearities typical for AEC applications are generally odd symmetric, the Maclaurin series expansion degenerates into the odd power series expansion, which has been successfully used to model memoryless loudspeaker nonlinearities [65], [68]. Instead of directly estimating the nonlinear system, several linear systems are used to approximate the nonlinearity and then the problem is transformed into one of linear system identification [63], [64]. According to [57], [58], [62]–[65], a typical memoryless nonlinearity can be expanded as

$$f[x(t)] = \sum_{n=1}^N b_n x^{2n-1}(t), \quad (2)$$

where N is the expansion order and b_n is the n th order odd power series expansion coefficient. It is worth noting that non-symmetric nonlinearities can also be addressed using Maclaurin series expansion by doubling the number of parameters. In such cases, the superiority of the proposed algorithms becomes even more significant.

B. Convulsive Transfer Function Model

To accelerate SBSS-NAEC algorithms, the MTF model representing time-domain convolution with a frequency-domain multiplication is often adopted. However, the MTF model requires an STFT analysis window to be at least as long as the effective part of the linear echo path, which will inevitably introduce significant algorithmic delay, especially in highly reverberant environments [40]. In order to address this limitation, CTF models which represent the time-domain convolution with a frequency-domain convolution is adopted [41]–[43]. CTF models are not restricted by the length of STFT window and, therefore, can achieve a compromise between the computational complexity and algorithmic delay [41]. With this model [42], the microphone signal can be expressed as [58], [65]

$$Y_{i,j} = S_{i,j} + \underbrace{\sum_{l=1}^L \sum_{n=1}^N b_n A_{i,j,l} X_{n,i,j-l+1}}_{V_{i,j}}, \quad (3)$$

where $i \in \{1, \dots, I\}$ and $j \in \{1, \dots, J\}$ are the frequency and time-frame indices, respectively, L is the length of the CTF subband filter, $Y_{i,j}$, $S_{i,j}$, $V_{i,j}$, and $X_{n,i,j}$ are the STFT-domain representations of $y(t)$, $s(t)$, $v(t)$, and $x^{2n-1}(t)$, respectively, and $A_{i,j,l}$ is a CTF filter coefficient representing the linear echo path in the CTF domain. Generally, for a given reverberation level, if shorter algorithmic delay is desired, a shorter STFT window has to be used. Consequently, a larger

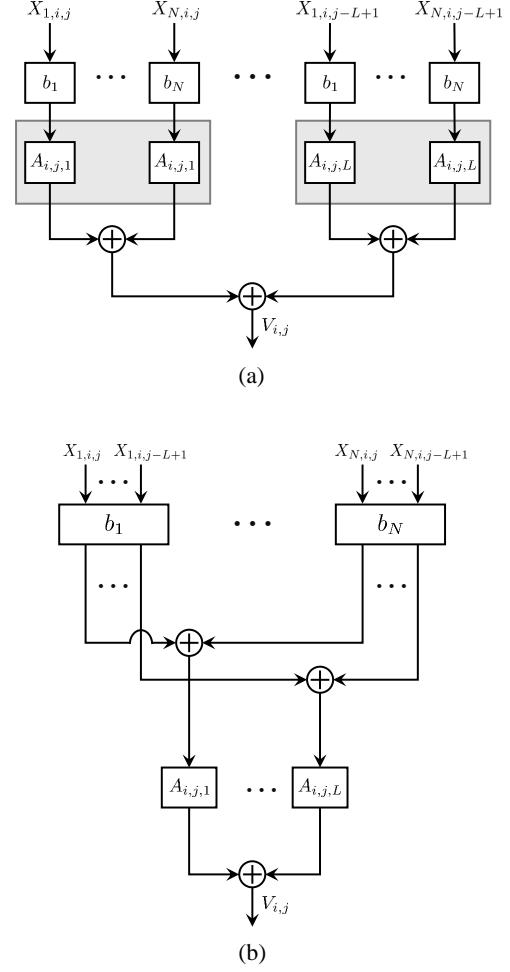


Fig. 2: Illustration of nonlinear echo model: (a) merged nonlinear echo model, (b) bilinear nonlinear echo model.

CTF filter length L has to be chosen [41], [43], which however inevitably leads to increase of the computational complexity.

C. Merged Nonlinear Echo Model

In the conventional SBSS-NAEC model [58], [65], the nonlinear coefficients are merged into the linear echo path as illustrated in Fig. 2 (a). Mathematically, this process can be expressed as

$$\begin{aligned} V_{i,j} &= \sum_{l=1}^L \sum_{n=1}^N b_n A_{i,j,l} X_{n,i,j-l+1} \\ &= \sum_{l=1}^L \sum_{n=1}^N A_{n,i,j,l}^m X_{n,i,j-l+1} \\ &= \sum_{l=1}^L (\mathbf{a}_{i,j,l}^m)^T \mathbf{x}_{i,j-l+1} \\ &= (\tilde{\mathbf{a}}_{i,j}^m)^T \tilde{\mathbf{x}}_{i,j}, \end{aligned} \quad (4)$$

where $A_{n,i,j,l}^m = b_n A_{i,j,l}$, the superscripts m and T denote, respectively, the merged model and transpose operator, and

$$\mathbf{a}_{i,j,l}^m = [A_{1,i,j,l}^m \ A_{2,i,j,l}^m \ \dots \ A_{N,i,j,l}^m]^T, \quad (5)$$

$$\tilde{\mathbf{a}}_{i,j}^m = [(\mathbf{a}_{i,j,1}^m)^T \ (\mathbf{a}_{i,j,2}^m)^T \ \dots \ (\mathbf{a}_{i,j,L}^m)^T]^T, \quad (6)$$

$$\mathbf{x}_{i,j-l} = [X_{1,i,j-l} \ X_{2,i,j-l} \ \dots \ X_{N,i,j-l}]^T, \quad (7)$$

$$\tilde{\mathbf{x}}_{i,j} = [\mathbf{x}_{i,j}^T \ \mathbf{x}_{i,j-1}^T \ \dots \ \mathbf{x}_{i,j-L+1}^T]^T. \quad (8)$$

As seen, in the merged model, the nonlinear echo is represented with $N \times I \times L$ parameters.

D. Bilinear Nonlinear Echo Model

Alternatively, as can be seen from (3), the nonlinear echo can be reformulated as a bilinear form of the odd power series expansion coefficients vector \mathbf{b} and the CTF filter vector $\mathbf{a}_{i,j}$:

$$\begin{aligned} V_{i,j} &= \mathbf{a}_{i,j}^T \mathbf{X}_{i,j} \mathbf{b}, \\ &= \mathbf{a}_{i,j}^T \mathbf{x}_{i,j}^a, \end{aligned} \quad (9)$$

or, equivalently,

$$\begin{aligned} V_{i,j} &= \mathbf{b}^T \mathbf{X}_{i,j}^T \mathbf{a}_{i,j} \\ &= \mathbf{b}^T \mathbf{x}_{i,j}^b, \end{aligned} \quad (10)$$

where

$$\mathbf{a}_{i,j} = [A_{i,j,1} \ A_{i,j,2} \ \dots \ A_{i,j,L}]^T, \quad (11)$$

$$\mathbf{b} = [b_1 \ b_2 \ \dots \ b_N]^T, \quad (12)$$

and

$$\mathbf{X}_{i,j} = \begin{bmatrix} X_{1,i,j} & X_{2,i,j} & \dots & X_{N,i,j} \\ X_{1,i,j-1} & X_{2,i,j-1} & \dots & X_{N,i,j-1} \\ \vdots & \vdots & \ddots & \vdots \\ X_{1,i,j-L+1} & X_{2,i,j-L+1} & \dots & X_{N,i,j-L+1} \end{bmatrix}. \quad (13)$$

The two transformed signal vectors $\mathbf{x}_{i,j}^a$ and $\mathbf{x}_{i,j}^b$ are of the following form:

$$\mathbf{x}_{i,j}^a = \mathbf{X}_{i,j} \mathbf{b}, \quad (14)$$

$$\mathbf{x}_{i,j}^b = \mathbf{X}_{i,j}^T \mathbf{a}_{i,j}. \quad (15)$$

By decoupling the series expansion parameters \mathbf{b} and the CTF filter coefficients $\mathbf{a}_{i,j}$ from each other by the bilinear nonlinear echo model, the number of free parameters for every time frame is reduced to $N + I \times L$. The bilinear nonlinear echo model is illustrated in Fig. 2 (b).

E. Probabilistic Signal Model

In the rest of this paper, we model the near-end signal with a generalized Gaussian distribution, which has been widely used in the literature [48], [49], [57], [58], [65]:

$$p(\mathbf{s}_j) \propto \exp \left[- \left(\frac{\|\mathbf{s}_j\|_2}{\gamma} \right)^\beta \right], \quad (16)$$

where

$$\mathbf{s}_j = [S_{1,j} \ S_{2,j} \ \dots \ S_{I,j}]^T \quad (17)$$

and $\|\cdot\|_2$ denotes the ℓ_2 norm. We assume that $\gamma > 0$ and $0 < \beta \leq 2$ to obtain supergaussian distributions in order to use the well-known majorize-minimization (MM) method [69].

III. CONVENTIONAL ALGORITHMS

In this section, we give a brief overview of the conventional IP and EISS algorithms, which are based on the merged nonlinear echo model. To formulate the NAEC problem as one of SBSS, the following mixing model is considered

$$\tilde{\mathbf{y}}_{i,j}^m = \mathbf{H}_{i,j}^m \tilde{\mathbf{s}}_{i,j}^m, \quad (18)$$

where

$$\tilde{\mathbf{y}}_{i,j}^m = [Y_{i,j} \ \tilde{\mathbf{x}}_{i,j}^T]^T, \quad (19)$$

$$\tilde{\mathbf{s}}_{i,j}^m = [S_{i,j} \ \tilde{\mathbf{x}}_{i,j}^T]^T \quad (20)$$

are two signal vectors, and

$$\mathbf{H}_{i,j}^m = \begin{bmatrix} 1 & (\tilde{\mathbf{a}}_{i,j}^m)^T \\ \mathbf{0}_{NL \times 1} & \mathbf{I}_{NL} \end{bmatrix} \quad (21)$$

is the mixing matrix, with $\mathbf{0}_{NL \times 1}$ being an all-zero column vector of length NL and \mathbf{I}_{NL} being the identity matrix of size $NL \times NL$. Now the problem can be solved from a signal separation perspective. The near-end signal can then be recovered through the following inverse system:

$$\mathbf{W}_{i,j}^m = \begin{bmatrix} 1 & -(\hat{\mathbf{a}}_{i,j}^m)^T \\ \mathbf{0}_{NL \times 1} & \mathbf{I}_{NL} \end{bmatrix}, \quad (22)$$

where $\hat{\mathbf{a}}_{i,j}^m$ is the estimate of $\tilde{\mathbf{a}}_{i,j}^m$. The MNE filter defined as

$$\mathbf{w}_{i,j}^m = [1 \ -(\hat{\mathbf{a}}_{i,j}^m)^T]^H, \quad (23)$$

where $(\cdot)^H$ denotes conjugate transpose, is used to extract the near-end signal, i.e.,

$$\hat{S}_{i,j} = (\mathbf{w}_{i,j}^m)^H \tilde{\mathbf{y}}_{i,j}^m. \quad (24)$$

It is worth to note that this signal separation problem is an SBSS problem as the far-end signal is precisely known. Now exploiting the mutual independence between the near-end and far-end signals, one can derive the following recursive negative log-likelihood function [58]:

$$\begin{aligned} \mathcal{L}_{i,j}^m &= -\frac{1}{\sum_{j'=1}^j (\eta^m)^{j-j'}} \sum_{j'=1}^j (\eta^m)^{j-j'} \log p(\mathbf{s}_{j'}) \\ &\quad - 2 \sum_{i=1}^I \log |\det \mathbf{W}_{i,j}^m|, \end{aligned} \quad (25)$$

where $\eta^m \in (0, 1)$ is a forgetting factor. With the MM method [48], [49], the following auxiliary function can be constructed

$$\mathcal{L}_{i,j}^{m,+} = \sum_{i=1}^I (\mathbf{w}_{i,j}^m)^H \mathbf{G}_{i,j}^m \mathbf{w}_{i,j}^m - 2 \sum_{i=1}^I \log |\det \mathbf{W}_{i,j}^m|, \quad (26)$$

where

$$\mathbf{G}_{i,j}^m = \eta^m \mathbf{G}_{i,j-1}^m + (1 - \eta^m) \varphi(\sigma_{s,j}^m) \tilde{\mathbf{y}}_{i,j}^m (\tilde{\mathbf{y}}_{i,j}^m)^H \quad (27)$$

is the auxiliary matrix and

$$\varphi(\sigma_{s,j}^m) = (\sigma_{s,j}^m)^{\beta-2}, \quad (28)$$

$$\sigma_{s,j}^m = \sqrt{\sum_{i=1}^I |(\mathbf{w}_{i,j-1}^m)^H \tilde{\mathbf{y}}_{i,j}^m|^2}, \quad (29)$$

is a weighting function. It can be shown that [48], [49]

$$\mathcal{L}_{i,j}^{m,+} \geq \mathcal{L}_{i,j}^m, \quad (30)$$

with equality if and only if $\mathbf{w}_{i,j} = \mathbf{w}_{i,j-1}$. The derivation of the auxiliary function is presented in Appendix A. Due to the structure of $\mathbf{W}_{i,j}^m$ given in (22), (26) is a simplified HEAD problem [47], which has been well studied in the literature of BSS [48], [52], [53]. Accordingly, the MNE filter (23) can be optimized adaptively by constructing and minimizing (26). Then the near-end signal is extracted according to (24).

A. Iterative Projection

By equating the Wirtinger derivative of (26) with respect to $(\mathbf{w}_{i,j}^m)^*$ to 0, where $(\cdot)^*$ denotes complex conjugation, one obtains

$$\mathbf{G}_{i,j}^m \mathbf{w}_{i,j}^m = (\mathbf{W}_{i,j}^m)^{-1} \mathbf{e}_{NL+1}, \quad (31)$$

where \mathbf{e}_{NL+1} is a unitary vector of length $NL + 1$, with the first element being 1. Then the IP-based update rule is obtained as [48], [49]

$$\mathbf{w}_{i,j}^m = (\mathbf{W}_{i,j}^m \mathbf{G}_{i,j}^m)^{-1} \mathbf{e}_{NL+1}. \quad (32)$$

Taking the structure of $\mathbf{W}_{i,j}^m$ into account, the update rule can be simplified as [57], [65]

$$\mathbf{w}_{i,j}^m = (\mathbf{G}_{i,j}^m)^{-1} \mathbf{e}_{NL+1}. \quad (33)$$

Now, we ensure the structure of the MNE filter (23), to be specific, the first element should be 1, by

$$\mathbf{w}_{i,j}^m := \frac{\mathbf{w}_{i,j}^m}{W_{i,j,1}^m}, \quad (34)$$

where $W_{i,j,1}^m$ is the first element of $\mathbf{w}_{i,j}^m$ and $:=$ denotes assignment.

One may notice that the IP algorithm updates the MNE filter in a two-step manner, i.e., update and normalization, which does not fully take advantage of the structure of the unit upper triangular structure of the demixing matrix (22).

B. Element-wise Iterative Source Steering

A more efficient update strategy is given by the EISS algorithm [58], which is modified from the rank-1 update rule proposed in [52]. For EISS, the MNE filter is updated element-wisely as

$$W_{i,j,k}^m = \begin{cases} W_{i,j-1,k}^m - U_{i,j,k}^m, & k = 1 \\ (1 - U_{i,j,1}^m) W_{i,j-1,k}^m - U_{i,j,k}^m, & \text{else} \end{cases}, \quad (35)$$

where $W_{i,j,k}^m$ is the k -th element of the MNE filter $\mathbf{w}_{i,j}^m$ and $U_{i,j,k}^m$ is a steering stepsize yet to be determined. Substituting (35) into (26) gives the following auxiliary function:

$$\begin{aligned} \mathcal{L}_{i,j}^{m,+} = & -2 \log \left| 1 - (U_{i,j,1}^m)^* \right| \\ & + \left(\mathbf{w}_{i,j}^m - \tilde{\mathbf{u}}_{i,j}^m \right)^H \mathbf{G}_j^m \left(\mathbf{w}_{i,j}^m - \tilde{\mathbf{u}}_{i,j}^m \right), \end{aligned} \quad (36)$$

where

$$\tilde{\mathbf{u}}_{i,j}^m = \begin{bmatrix} U_{i,j,1}^m & U_{i,j,1}^m W_{i,j-1,2}^m + U_{i,j,2}^m & \dots \\ U_{i,j,1}^m W_{i,j-1,NL+1}^m + U_{i,j,NL+1}^m \end{bmatrix}^T \quad (37)$$

is a vector containing the parameters to be estimated. Then, forcing the derivative with respect to $(U_{i,j,k}^m)^*$ to 0 gives

$$U_{i,j,k}^m = \begin{cases} 1 - \left[(\mathbf{w}_{i,j}^m)^H \mathbf{G}_{i,j}^m \mathbf{w}_{i,j}^m \right]^{-\frac{1}{2}}, & k = 1 \\ \frac{(\mathbf{g}_{i,j,k}^m)^H \mathbf{w}_{i,j}^m}{G_{i,j,k}^m}, & \text{else} \end{cases}, \quad (38)$$

where $\mathbf{g}_{i,j,k}^m$ is the k -th column of matrix $\mathbf{G}_{i,j}^m$ and $G_{i,j,k}^m$ is the k -th diagonal element of $\mathbf{G}_{i,j}^m$. After updating the MNE filter with (38) and (35), the normalization given in (34) is performed.

Again, the EISS algorithm updates the MNE filter coefficients in the two-step manner, which does not take the unit upper triangular structure of the demixing matrix (22) into full consideration.

IV. PROPOSED METHODS

As can be seen, the MNE filter length is $NL + 1$. Consequently, the computational complexity for identifying the merged nonlinear echo model is high. To address this issue, we propose two improved algorithms based on the previously discussed bilinear nonlinear echo model, which take advantage of the structure of the demixing system to update the filter coefficients with a one-step strategy. Note that the two-step update strategy in the IP and EISS algorithms can also be adopted to the presented methods; but the dimension of the related matrices is higher, which makes the update process computationally more expensive.

Firstly, if the series expansion coefficients \mathbf{b}_j are fixed, we can formulate the identification of the CTF filter $\mathbf{a}_{i,j}$ as

$$\tilde{\mathbf{y}}_{i,j}^a = \mathbf{H}_{i,j}^a \tilde{\mathbf{s}}_{i,j}^a, \quad (39)$$

where

$$\tilde{\mathbf{y}}_{i,j}^a = \begin{bmatrix} Y_{i,j} & (\mathbf{x}_{i,j}^a)^T \end{bmatrix}^T, \quad (40)$$

$$\tilde{\mathbf{s}}_{i,j}^a = \begin{bmatrix} S_{i,j} & (\mathbf{x}_{i,j}^a)^T \end{bmatrix}^T \quad (41)$$

are two signal vectors, and

$$\mathbf{H}_{i,j}^a = \begin{bmatrix} 1 & \mathbf{a}_{i,j}^T \\ \mathbf{0}_{L \times 1} & \mathbf{I}_L \end{bmatrix} \quad (42)$$

is the mixing matrix. The near-end signal is extracted by

applying the demixing matrix:

$$\mathbf{W}_{i,j}^a = \begin{bmatrix} 1 & -\hat{\mathbf{a}}_{i,j}^T \\ \mathbf{0}_{L \times 1} & \mathbf{I}_L \end{bmatrix}, \quad (43)$$

where $\hat{\mathbf{a}}_{i,j}^T$ is the estimate of $\mathbf{a}_{i,j}^T$. With the filter

$$\mathbf{w}_{i,j}^a = [1 \ -\hat{\mathbf{a}}_{i,j}^T]^H, \quad (44)$$

the near-end signal can be extracted

$$\hat{S}_{i,j} = (\mathbf{w}_{i,j}^a)^H \tilde{\mathbf{y}}_{i,j}^a. \quad (45)$$

To estimate $\mathbf{a}_{i,j}$, the following cost function is considered

$$\begin{aligned} \mathcal{L}_{i,j}^a = -\frac{1}{\sum_{j'=1}^j (\eta^a)^{j-j'}} \sum_{j'=1}^j (\eta^a)^{j-j'} \log p(\mathbf{s}_{j'}) \\ - 2 \sum_{i=1}^I \log |\det \mathbf{W}_{i,j}^a|, \end{aligned} \quad (46)$$

where $\eta^a \in (0, 1)$ is a forgetting factor. As the first element of $\mathbf{w}_{i,j}^a$ is fixed to 1, we obtain due to the one-step strategy:

$$\det \mathbf{W}_{i,j}^a = 1. \quad (47)$$

Additionally, with the majorize-minimization (MM) method, the following auxiliary function can be constructed

$$\mathcal{L}_{i,j}^{a,+} = \sum_{i=1}^I (\mathbf{w}_{i,j}^a)^H \mathbf{G}_{i,j}^a \mathbf{w}_{i,j}^a, \quad (48)$$

where the auxiliary matrix $\mathbf{G}_{i,j}^a$ can be partitioned as

$$\mathbf{G}_{i,j}^a = \begin{bmatrix} (\sigma_{y,i,j}^a)^2 & (\mathbf{q}_{i,j}^a)^H \\ \mathbf{q}_{i,j}^a & \mathbf{R}_{i,j}^a \end{bmatrix}, \quad (49)$$

where

$$(\sigma_{y,i,j}^a)^2 = \eta^a (\sigma_{y,i,j-1}^a)^2 + (1 - \eta^a) \varphi(\sigma_{s,j}^a) |Y_{i,j}|^2, \quad (50)$$

$$\mathbf{q}_{i,j}^a = \eta^a \mathbf{q}_{i,j-1}^a + (1 - \eta^a) \varphi(\sigma_{s,j}^a) Y_{i,j}^* \mathbf{x}_{i,j}^a, \quad (51)$$

$$\mathbf{R}_{i,j}^a = \eta^a \mathbf{R}_{i,j-1}^a + (1 - \eta^a) \varphi(\sigma_{s,j}^a) \mathbf{x}_{i,j}^a (\mathbf{x}_{i,j}^a)^H, \quad (52)$$

with

$$\varphi(\sigma_{s,j}^a) = (\sigma_{s,j}^a)^{\beta-2}, \quad (53)$$

and

$$\sigma_{s,j}^a = \sqrt{\sum_{i=1}^I |(\mathbf{w}_{i,j-1}^a)^H \tilde{\mathbf{y}}_{i,j}^a|^2}. \quad (54)$$

Analogously to (30), one can verify that

$$\mathcal{L}_{i,j}^{a,+} \geq \mathcal{L}_{i,j}^a. \quad (55)$$

Now, $\hat{\mathbf{a}}_{i,j}$ can be optimized by solving

$$\hat{\mathbf{a}}_{i,j} = \operatorname{argmin}_{\mathbf{a}_{i,j}} \mathcal{L}_{i,j}^{a,+}, \quad \text{s.t. } (\mathbf{w}_{i,j}^a)^H \mathbf{e}_{L+1} = 1. \quad (56)$$

After $\hat{\mathbf{a}}_{i,j}$ are updated, one can construct $\mathbf{x}_{i,j}^b$ according to (15). The series expansion coefficients can then be identified

as

$$\tilde{\mathbf{y}}_{i,j}^b = \mathbf{H}^b \tilde{\mathbf{s}}_{i,j}^b, \quad (57)$$

where

$$\tilde{\mathbf{y}}_{i,j}^b = [Y_{i,j} \ (\mathbf{x}_{i,j}^b)^T]^T, \quad (58)$$

$$\tilde{\mathbf{s}}_{i,j}^b = [S_{i,j} \ (\mathbf{x}_{i,j}^b)^T]^T, \quad (59)$$

$$\mathbf{H}^b = \begin{bmatrix} 1 & \mathbf{b}^T \\ \mathbf{0}_{N \times 1} & \mathbf{I}_N \end{bmatrix}. \quad (60)$$

The estimated demixing matrix is

$$\mathbf{W}_j^b = \begin{bmatrix} 1 & -\hat{\mathbf{b}}_j^T \\ \mathbf{0}_{N \times 1} & \mathbf{I}_N \end{bmatrix}, \quad (61)$$

where $\hat{\mathbf{b}}_j$ is the estimate of \mathbf{b} at time-frame j . Using the following filter:

$$\mathbf{w}_j^b = [1 \ -\hat{\mathbf{b}}_j^T]^H, \quad (62)$$

one can extract the near-end signal as

$$\hat{S}_{i,j} = (\mathbf{w}_j^b)^H \tilde{\mathbf{y}}_{i,j}^b. \quad (63)$$

Now, the recursive negative log-likelihood function for $\hat{\mathbf{b}}_j$ is written as

$$\begin{aligned} \mathcal{L}_j^b = -\frac{1}{\sum_{j'=1}^j (\eta^b)^{j-j'}} \sum_{j'=1}^j (\eta^b)^{j-j'} \log p(\mathbf{s}_{j'}) \\ - 2I \log |\det \mathbf{W}_j^b|, \end{aligned} \quad (64)$$

where $\eta^b \in (0, 1)$ is a forgetting factor for estimating the series expansion coefficients. Analogously, with the one-step strategy, we have

$$\det \mathbf{W}_j^b = 1. \quad (65)$$

Therefore, with the MM method, the following auxiliary function is constructed

$$\mathcal{L}_j^{b,+} = (\mathbf{w}_j^b)^H \bar{\mathbf{G}}_j^b \mathbf{w}_j^b. \quad (66)$$

The auxiliary matrix $\bar{\mathbf{G}}_j^b$ can again be partitioned as

$$\bar{\mathbf{G}}_j^b = \begin{bmatrix} (\bar{\sigma}_{y,j}^b)^2 & (\bar{\mathbf{q}}_j^b)^H \\ \bar{\mathbf{q}}_j^b & \bar{\mathbf{R}}_j^b \end{bmatrix}, \quad (67)$$

where

$$(\bar{\sigma}_{y,j}^b)^2 = \eta^b (\bar{\sigma}_{y,j-1}^b)^2 + \frac{(1 - \eta^b) \varphi(\sigma_{s,j}^b)}{I} \sum_{i=1}^I |Y_{i,j}|^2, \quad (68)$$

$$\bar{\mathbf{q}}_j^b = \eta^b \bar{\mathbf{q}}_{j-1}^b + \frac{(1 - \eta^b) \varphi(\sigma_{s,j}^b)}{I} \sum_{i=1}^I Y_{i,j}^* \mathbf{x}_{i,j}^b, \quad (69)$$

$$\bar{\mathbf{R}}_j^b = \eta^b \bar{\mathbf{R}}_{j-1}^b + \frac{(1 - \eta^b) \varphi(\sigma_{s,j}^b)}{I} \sum_{i=1}^I \mathbf{x}_{i,j}^b (\mathbf{x}_{i,j}^b)^H, \quad (70)$$

with

$$\varphi(\sigma_{s,j}^b) = (\sigma_{s,j}^b)^{\beta-2}, \quad (71)$$

and

$$\sigma_{s,j}^b = \sqrt{\sum_{i=1}^I |(\mathbf{w}_{j-1}^b)^H \tilde{\mathbf{y}}_{i,j}^b|^2}. \quad (72)$$

Then the nonlinear coefficients can be estimated by solving

$$\hat{\mathbf{b}}_j = \operatorname{argmin}_{\mathbf{L}_j^{b,+}} \mathcal{L}_j^{b,+}, \quad \text{s.t. } (\mathbf{w}_j^b)^H \mathbf{e}_{N+1} = 1. \quad (73)$$

Now, the entire system can be identified by iteratively and alternately updating $\hat{\mathbf{w}}_{i,j}^a$ and $\hat{\mathbf{w}}_j^b$. Finally, the near-end signal can be extracted by applying (63).

A. Alternating Iterative Projection

We now adopt the IP algorithm for our bilinear alternating optimization scheme by directly solving (56) and (73). Note that as $\det(\mathbf{W}_{i,j}^a) = \det(\mathbf{W}_j^b) = 1$, the HEAD problem degenerates into one of linear constrained quadratic programming (LCQP) [70]–[72], which has been widely studied in minimum variance distortionless response (MVDR) filter [70], [73], [74] and linearly constrained minimum variance (LCMV) filter [70], [75]–[77] optimization.

Firstly, we use the estimate $\hat{\mathbf{b}}_j$ from the previous frame to construct $\mathbf{x}_{i,j}^a$ and update the associated statistics. The solution of (56) is given as [73]

$$\mathbf{w}_{i,j}^a = \frac{(\mathbf{G}_{i,j}^a)^{-1} \mathbf{e}_{L+1}}{\mathbf{e}_{L+1}^T (\mathbf{G}_{i,j}^a)^{-1} \mathbf{e}_{L+1}}. \quad (74)$$

Assuming that both $\mathbf{G}_{i,j}^a$ and $\mathbf{R}_{i,j}^a$ are invertible, one can express the inverse of $\mathbf{G}_{i,j}^a$ as

$$(\mathbf{G}_{i,j}^a)^{-1} = \begin{bmatrix} (\mathbf{S}_{i,j}^a)^{-1} & -(\mathbf{S}_{i,j}^a)^{-1} (\mathbf{c}_{i,j}^a)^H \\ -(\mathbf{S}_{i,j}^a)^{-1} \mathbf{c}_{i,j}^a & (\mathbf{R}_{i,j}^a)^{-1} + (\mathbf{S}_{i,j}^a)^{-1} \mathbf{D}_{i,j}^a \end{bmatrix}, \quad (75)$$

where

$$\mathbf{S}_{i,j}^a = (\sigma_{y,i,j}^a)^2 - (\mathbf{q}_{i,j}^a)^H (\mathbf{R}_{i,j}^a)^{-1} \mathbf{q}_{i,j}^a \quad (76)$$

is the Schur complement of $\mathbf{R}_{i,j}^a$ in $\mathbf{G}_{i,j}^a$ [78]–[80] and

$$\mathbf{c}_{i,j}^a = (\mathbf{R}_{i,j}^a)^{-1} \mathbf{q}_{i,j}^a, \quad (77)$$

$$\mathbf{D}_{i,j}^a = \mathbf{c}_{i,j}^a (\mathbf{c}_{i,j}^a)^H. \quad (78)$$

Substituting (75) into (74), we have

$$(\mathbf{S}_{i,j}^a)^{-1} \begin{bmatrix} 1 \\ -\mathbf{a}_{i,j}^* \end{bmatrix} = \begin{bmatrix} (\mathbf{S}_{i,j}^a)^{-1} \\ -(\mathbf{S}_{i,j}^a)^{-1} \mathbf{c}_{i,j}^a \end{bmatrix}. \quad (79)$$

Therefore, the CTF filter can be updated as

$$\hat{\mathbf{a}}_{i,j} = \left[(\mathbf{R}_{i,j}^a)^{-1} \mathbf{q}_{i,j}^a \right]^*. \quad (80)$$

Algorithm 1 Alternating iterative projection algorithm

1: Setting forgetting factors η^a and η^b

Initial parameters

$$\mathbf{X}_{i,0} = \mathbf{0}_{L \times N}, \quad \hat{\mathbf{a}}_{i,0} = \mathbf{0}_{L \times 1}, \quad \hat{\mathbf{b}}_0 = \begin{bmatrix} 1 & \mathbf{0}_{(N-1) \times 1}^T \end{bmatrix}^T, \\ \mathbf{q}_{i,0}^a = \mathbf{0}_{L \times 1}, \quad \bar{\mathbf{q}}_0^b = \mathbf{0}_{N \times 1},$$

$$\mathbf{R}_{i,0}^a = 10^{-4} \mathbf{I}_L, \quad \bar{\mathbf{R}}_0^b = 10^{-4} \mathbf{I}_N$$

2: **for** $j = 1; j < J; j = j + 1$ **do**

3: Insert $X_{n,i,j}$ into $\mathbf{X}_{i,j-1}$ to get $\mathbf{X}_{i,j}$

4: Update CTF filter associated statistics with Eq. (14), (40), (51), (52), (53), (54)

$$\mathbf{x}_{i,j}^a = \mathbf{X}_{i,j} \hat{\mathbf{b}}_{j-1},$$

$$\tilde{\mathbf{y}}_{i,j}^a = \begin{bmatrix} Y_{i,j} & (\mathbf{x}_{i,j}^a)^T \end{bmatrix}^T$$

$$\sigma_{s,j}^a = \sqrt{\sum_{i=1}^I |(\mathbf{w}_{i,j-1}^a)^H \tilde{\mathbf{y}}_{i,j}^a|^2}$$

$$\varphi(\sigma_{s,j}^a) = (\sigma_{s,j}^a)^{\beta-2}$$

$$\mathbf{q}_{i,j}^a = \eta^a \mathbf{q}_{i,j-1}^a + (1 - \eta^a) \varphi(\sigma_{s,j}^a) Y_{i,j}^* \mathbf{x}_{i,j}^a$$

$$\mathbf{R}_{i,j}^a = \eta^a \mathbf{R}_{i,j-1}^a + (1 - \eta^a) \varphi(\sigma_{s,j}^a) \mathbf{x}_{i,j}^a (\mathbf{x}_{i,j}^a)^H$$

5: Update CTF filter with Eq. (80)

$$\hat{\mathbf{a}}_{i,j} = \left[(\mathbf{R}_{i,j}^a)^{-1} \mathbf{q}_{i,j}^a \right]^*$$

6: Update series expansion coefficients associated statistics with Eq. (15), (58), (69), (70), (71), (72)

$$\mathbf{x}_{i,j}^b = \mathbf{X}_{i,j}^T \hat{\mathbf{a}}_{i,j}$$

$$\tilde{\mathbf{y}}_{i,j}^b = \begin{bmatrix} Y_{i,j} & (\mathbf{x}_{i,j}^b)^T \end{bmatrix}^T$$

$$\sigma_{s,j}^b = \sqrt{\sum_{i=1}^I |(\mathbf{w}_{i,j-1}^b)^H \tilde{\mathbf{y}}_{i,j}^b|^2}$$

$$\varphi(\sigma_{s,j}^b) = (\sigma_{s,j}^b)^{\beta-2}$$

$$\bar{\mathbf{q}}_j^b = \eta^b \bar{\mathbf{q}}_{j-1}^b + \frac{(1-\eta^b)\varphi(\sigma_{s,j}^b)}{I} \sum_{i=1}^I Y_{i,j}^* \mathbf{x}_{i,j}^b$$

$$\bar{\mathbf{R}}_j^b = \eta^b \bar{\mathbf{R}}_{j-1}^b + \frac{(1-\eta^b)\varphi(\sigma_{s,j}^b)}{I} \sum_{i=1}^I \mathbf{x}_{i,j}^b (\mathbf{x}_{i,j}^b)^H$$

7: Update series expansion coefficients with Eq. (87)

$$\hat{\mathbf{b}}_j = \left[(\bar{\mathbf{R}}_j^b)^{-1} \bar{\mathbf{q}}_j^b \right]^*$$

8: Extract near-end signal

$$\hat{S}_{i,j} = (\mathbf{w}_j^b)^H \tilde{\mathbf{y}}_{i,j}^b$$

9: **end for**

After updating the CTF filter $\hat{\mathbf{a}}_{i,j}$, we use it to construct $\mathbf{x}_{i,j}^b$. The solution of (73) is given by

$$\mathbf{w}_j^b = \frac{(\bar{\mathbf{R}}_j^b)^{-1} \mathbf{e}_{N+1}}{\mathbf{e}_{N+1}^T (\bar{\mathbf{R}}_j^b)^{-1} \mathbf{e}_{N+1}}. \quad (81)$$

Following the previous analysis, one can express the inverse of $\bar{\mathbf{G}}_j^b$ as

$$(\bar{\mathbf{G}}_j^b)^{-1} = \begin{bmatrix} (\bar{S}_j^b)^{-1} & -(\bar{S}_j^b)^{-1} (\bar{\mathbf{c}}_j^b)^H \\ -(\bar{S}_j^b)^{-1} \bar{\mathbf{c}}_j^b & (\bar{\mathbf{R}}_j^b)^{-1} + (\bar{S}_j^b)^{-1} \bar{\mathbf{D}}_j^b \end{bmatrix}, \quad (82)$$

where

$$\bar{S}_j^b = (\bar{\sigma}_{y,j}^b)^2 - (\bar{\mathbf{q}}_j^b)^H (\bar{\mathbf{R}}_j^b)^{-1} \bar{\mathbf{q}}_j^b \quad (83)$$

is the Schur complement of $\bar{\mathbf{R}}_j^b$ in $\bar{\mathbf{G}}_j^b$ and

$$\bar{\mathbf{c}}_j^b = (\bar{\mathbf{R}}_j^b)^{-1} \bar{\mathbf{q}}_j^b, \quad (84)$$

$$\bar{\mathbf{D}}_j^b = \bar{\mathbf{c}}_j^b (\bar{\mathbf{c}}_j^b)^H. \quad (85)$$

Substituting (82) into (81) gives

$$(\bar{S}_j^b)^{-1} \begin{bmatrix} 1 \\ -\mathbf{b}_j^* \end{bmatrix} = \begin{bmatrix} (\bar{S}_j^b)^{-1} \\ -(\bar{S}_j^b)^{-1} \bar{\mathbf{c}}_j^b \end{bmatrix}. \quad (86)$$

Therefore, the CTF filter can be updated as

$$\hat{\mathbf{b}}_j = \left[(\bar{\mathbf{R}}_j^b)^{-1} \bar{\mathbf{q}}_j^b \right]^*. \quad (87)$$

We summarize the AIP algorithm in Algorithm 1. Based on (80) and (87), RLS-like algorithms can be derived, which, however, will be left to the reader's investigation.

B. Alternating Element-wise Iterative Source Steering

In the following, we extend a previously proposed EIIS [58] algorithm for the use in our bilinear alternating optimization. We first construct $\mathbf{x}_{i,j}^a$ with $\hat{\mathbf{b}}_{j-1}$ and update $\hat{\mathbf{a}}_{i,j}$. As we fix the first element in $\mathbf{w}_{i,j}^a$ to 1, the original EIIS can be further simplified by skipping calculating the first steering stepsize. Therefore, one can directly update every element in the CTF filter. The update rule is given by

$$\hat{A}_{i,j,l} = \hat{A}_{i,j-1,l} + U_{i,j,l}^a, \quad l = 1, 2, \dots, L, \quad (88)$$

where $\hat{A}_{i,j,l}$ is the estimate of $A_{i,j,l}$ and $U_{i,j,l}^a$ is a steering stepsize yet to be determined. Substituting (88) into (48) gives

$$\mathcal{L}^{a,+} = \left[\mathbf{w}_{i,j}^a - \tilde{\mathbf{u}}_{i,j}^a \right]^H \mathbf{G}_{i,j}^a \left[\mathbf{w}_{i,j}^a - \tilde{\mathbf{u}}_{i,j}^a \right], \quad (89)$$

where

$$\tilde{\mathbf{u}}_{i,j}^a = \begin{bmatrix} 0 & (U_{i,j,1}^a)^* & \dots & (U_{i,j,L}^a)^* \end{bmatrix}^T. \quad (90)$$

Forcing the derivative of (89) with respect to $(U_{i,j,l}^a)^*$ to be 0, one can determine the steering stepsize as

$$U_{i,j,l}^a = \frac{(\mathbf{w}_{i,j-1}^a)^H \mathbf{g}_{i,j,l+1}^a}{G_{i,j,l+1,l+1}^a}, \quad (91)$$

where $\mathbf{g}_{i,j,l+1}^a$ denotes the $(l+1)$ th column of $\mathbf{G}_{i,j}^a$ and $g_{i,j,l+1,l+1}^a$ is the $(l+1)$ th diagonal element of $\mathbf{G}_{i,j}^a$. Following

Algorithm 2 Alternating element-wise iterative source steering algorithm

- 1: Setting forgetting factors η^a and η^b
Initial parameters
 $\mathbf{X}_{i,0} = \mathbf{0}_{L \times N}$, $\hat{\mathbf{a}}_{i,0} = \mathbf{0}_{L \times 1}$, $\hat{\mathbf{b}}_0 = [1 \quad \mathbf{0}_{(N-1) \times 1}]^T$,
 $\mathbf{q}_{i,0}^a = \mathbf{0}_{L \times 1}$, $\bar{\mathbf{q}}_0^b = \mathbf{0}_{N \times 1}$,
 $\mathbf{R}_{i,0}^a = 10^{-4} \mathbf{I}_L$, $\bar{\mathbf{R}}_0^b = 10^{-4} \mathbf{I}_N$
 - 2: **for** $j = 1; j < J; j = j + 1$ **do**
 - 3: Insert $X_{n,i,j}$ into $\mathbf{X}_{i,j-1}$ to get $\mathbf{X}_{i,j}$
 - 4: Update CTF filter associated statistics with Eq. (14), (40), (51), (52), (53), (54)
 - 5: **for** $l = 1; l < L; l = l + 1$ **do**
 - 6: Calculate the steering stepsize for the l -th element in the CTF filter with Eq. (92)
 $U_{i,j,l}^a = \frac{(q_{i,j,l}^a)^* - \hat{\mathbf{a}}_{j-1}^T \mathbf{r}_{i,j,l}^a}{R_{i,j,l,l}^a}$
 - 7: Update the l -th element in the CTF filter with Eq. (88)
 $\hat{A}_{i,j,l} = \hat{A}_{i,j-1,l} + U_{i,j,l}^a$
 - 8: **end for**
 - 9: Update series expansion coefficients associated statistics with Eq. (15), (58), (69), (70), (71), (72)
 - 10: **for** $n = 1; n < N; n = n + 1$ **do**
 - 11: Calculate the steering stepsize for the n -th series expansion coefficient with Eq. (97)
 $u_{j,n}^b = \frac{(\bar{q}_{j,n}^b)^* - \hat{\mathbf{b}}_{j-1}^T \bar{\mathbf{r}}_{j,n}^b}{\bar{r}_{j,n,n}^b}$
 - 12: Update the n -th series expansion coefficient with (93)
 $\hat{b}_{j,n} = \hat{b}_{j-1,n} + u_{j,n}^b$
 - 13: **end for**
 - 14: Extract near-end signal with Eq. (63)
 $\hat{S}_{i,j} = (\mathbf{w}_j^b)^H \tilde{\mathbf{y}}_{i,j}^b$
 - 15: **end for**
-

the structure of (44) and (49), we have

$$U_{i,j,l}^a = \frac{(q_{i,j,l}^a)^* - \hat{\mathbf{a}}_{j-1}^T \mathbf{r}_{i,j,l}^a}{R_{i,j,l,l}^a}, \quad (92)$$

where $\mathbf{r}_{i,j,l}^a$ denotes the l -th column of $\mathbf{R}_{i,j}^a$, $q_{i,j,l}^a$ is the l -th element of $\mathbf{q}_{i,j}^a$ and $R_{i,j,l,l}^a$ is the l -th diagonal element of $\mathbf{R}_{i,j}^a$. In comparison with (92), (91) needs to update $\mathbf{R}_{i,j}^a$ rather than $\mathbf{G}_{i,j}^a$, which helps reduce the computation cost. Then the CTF filter $\hat{\mathbf{a}}_{i,j}$ is updated with (88).

Note that the coefficients b_n are independent of frequency. Therefore, we use lower case to denote the associated scalars. Now, we construct $\mathbf{x}_{i,j}^b$ with the estimated $\hat{\mathbf{a}}_{i,j}$ and update $\hat{\mathbf{b}}_j$.

TABLE II: Main equations and dominant computational complexity for each step of all SBSS-NAEC algorithms.

	IP	EISS	AIP	AEISS
Step 1	Eq. (27)-(29) $\mathcal{O}[(NL+1)^2]$	Eq. (27)-(29) $\mathcal{O}[(NL+1)^2]$	Eq. (51)-(54), (69)-(72) $\max[\mathcal{O}(N^2), \mathcal{O}(L^2)]$	Eq. (51)-(54), (69)-(72) $\max[\mathcal{O}(N^2), \mathcal{O}(L^2)]$
Step 2	Eq. (33), (34) $\mathcal{O}[(NL+1)^3]$	Eq. (34), (35), (38) $\mathcal{O}[(NL+1)^2]$	Eq. (80), (87) $\max[\mathcal{O}(N^3), \mathcal{O}(L^3)]$	Eq. (88), (92), (93), (97) $\max[\mathcal{O}(N^2), \mathcal{O}(L^2)]$
Step 3	Eq. (24) $\mathcal{O}(NL)$	Eq. (24) $\mathcal{O}(NL)$	Eq. (63) $\max[\mathcal{O}(N), \mathcal{O}(L)]$	Eq. (63) $\max[\mathcal{O}(N), \mathcal{O}(L)]$
Overall	$\mathcal{O}[(NL+1)^3]$	$\mathcal{O}[(NL+1)^2]$	$\max[\mathcal{O}(N^3), \mathcal{O}(L^3)]$	$\max[\mathcal{O}(N^2), \mathcal{O}(L^2)]$

Similarly, the nonlinear coefficients are updated as

$$\hat{b}_{j,n} = \hat{b}_{j-1,n} + u_{j,n}^b, \quad n = 1, 2, \dots, N, \quad (93)$$

where $\hat{b}_{j,n}$ is the estimate of b_n in time frame j and $u_{j,n}^b$ is a steering stepsize to be determined. Substituting (93) into (66) gives the following expression:

$$\mathcal{L}^{b,+} = (\mathbf{w}_j^b - \tilde{\mathbf{u}}_j^b)^H \bar{\mathbf{G}}_j^b (\mathbf{w}_j^b - \tilde{\mathbf{u}}_j^b), \quad (94)$$

where

$$\tilde{\mathbf{u}}_j^b = [0 \quad (u_{j,1}^b)^* \quad \dots \quad (u_{j,N}^b)^*]^T. \quad (95)$$

By optimizing (94) with respect to $(u_{j,n}^b)^*$, the steering stepsize is determined as

$$u_{j,n}^b = \frac{(\mathbf{w}_{j-1}^b)^H \bar{\mathbf{g}}_{j,n+1}^b}{\bar{g}_{j,n+1,n+1}^b}, \quad (96)$$

where $\bar{\mathbf{g}}_{j,n+1}^b$ is the $(n+1)$ th column of $\bar{\mathbf{G}}_j^b$ and $\bar{g}_{j,n+1,n+1}^b$ is the $(n+1)$ th diagonal element of $\bar{\mathbf{G}}_j^b$. By considering the structure of (44) and (49), we have

$$u_{j,n}^b = \frac{(\bar{q}_{j,n}^b)^* - \hat{b}_{j-1}^T \bar{\mathbf{r}}_{j,n}^b}{\bar{r}_{j,n,n}^b}, \quad (97)$$

where $\bar{\mathbf{r}}_{j,n}^b$, $\bar{q}_{j,n}^b$ and $\bar{r}_{j,n,n}^b$ are the n th column of $\bar{\mathbf{R}}_j^b$, n th element of $\bar{\mathbf{q}}_j^b$ and n -th diagonal element of $\bar{\mathbf{R}}_j^b$, respectively.

Then, the nonlinear coefficients \hat{b}_j are updated using (93). The AEISS algorithm is summarized in Algorithm 2.

V. COMPLEXITY ANALYSIS

We now analyze the computational complexity of the proposed AIP and AEISS algorithms and compare them with the conventional IP and EISS algorithms. All algorithms consist of the following three fundamental steps: 1) updating the associated statistics, 2) updating the corresponding filters, and 3) extracting the near-end signal. The computational complexity depends on the nonlinear echo model, the filter length, and the chosen optimization method. The dominant complexity and the associated equation for each step to process a single time-frequency component is presented in Table II. For IP-

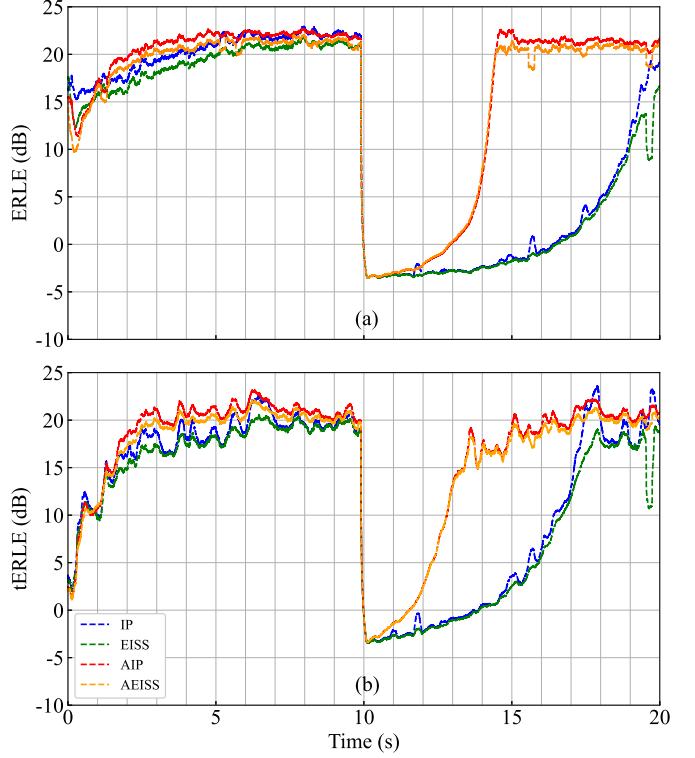


Fig. 3: Tracking ability of all the compared SBSS-NAEC algorithms. The far-end signal is an AR(1) signal. The AIR changes at the 10 second. (a) ERLE over time in the single-talk case and (b) tERLE over time in double-talk.

based methods, the dominant computational cost comes from the inversion of the auxiliary covariance matrix. Therefore, the complexity of the original IP algorithm is $\mathcal{O}[(NL+1)^3]$. With the bilinear nonlinear echo model and one-step strategy, AIP successfully reduces it to $\max[\mathcal{O}(N^3), \mathcal{O}(L^3)]$. As for EISS-based methods, the computational cost is dominated by updating the associated statistics and by the calculation of the steering stepsizes. The complexity of the original EISS algorithm is $\mathcal{O}[(NL+1)^2]$ while AEISS reduces it to $\max[\mathcal{O}(N^2), \mathcal{O}(L^2)]$. It is clear, that the computational com-

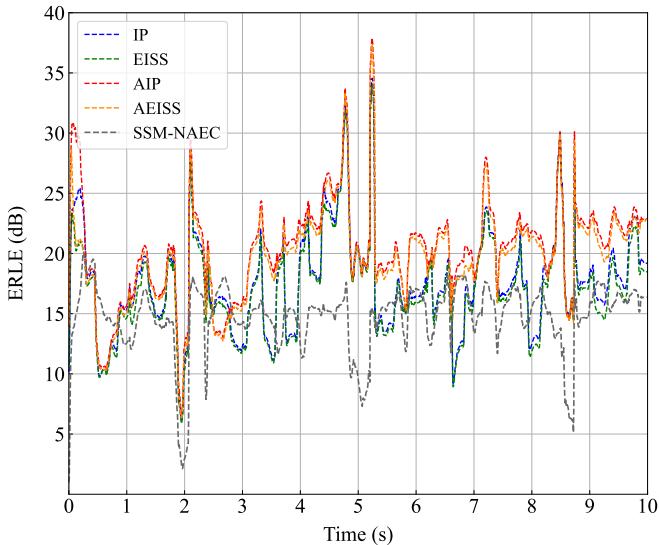


Fig. 4: ERLE performance achieved by all compared algorithms in a single-talk situation.

plexity of EISS-based methods is one order lower than the IP-based counterpart. Obviously, for large P and L , the AIP and AEISSL algorithms are much more computationally efficient than the conventional counterparts, indicating that they are much more suitable for real-time communications (RTC) in low-resource hardware.

VI. SIMULATIONS AND RESULTS

For consistency, we use the same speech signals as in the paper [65]. The sampling rate of all signals is 16 kHz. We randomly picked two AIRs from RWCP_E2A [82], which correspond to a reverberation time of approximately 300 ms. To model loudspeaker distortions, we consider the hard-clipping [83] function, which is of the following form:

$$f[x(t)] = \begin{cases} -\rho, & x(t) < -\rho \\ x(t), & |x(t)| \leq \rho \\ \rho, & x(t) > \rho \end{cases}, \quad (98)$$

where $\rho > 0$ is the maximum output amplitude of the loudspeaker. In our simulation, we set the value of ρ consistent to that in the paper [65], i.e., $\rho = 0.2x_{\max}$, where $x_{\max} = \max(|x(t)|)$ is the maximum amplitude of the original signal. To model the nonlinearity, we set the order of odd power series expansion to $N = 5$. A von Hann window is used and the window length is 1024 samples. The overlap between consecutive frames is 75%. Since the STFT window is much shorter than the AIR, we set the CTF filter length to $L = 5$. All experiments are conducted on a laptop with i7-12700H CPU. For the IP and EISS methods, we set the forgetting factors to $\eta^m = 0.992$. To make the algorithms robust, the auxiliary matrix $\tilde{\mathbf{G}}_{i,0}^m$ is set to $10^{-3} \times \mathbf{I}_{NL+1}$. For the AIP and AEISSL methods, we set $\eta^a = \eta^b = 0.98$. The auxiliary matrices $\mathbf{G}_{i,0}^a$ and $\tilde{\mathbf{G}}_0^b$ are set to be $10^{-4} \times \mathbf{I}_L$ and $10^{-4} \times \mathbf{I}_N$, respectively, and the two auxiliary vectors $\mathbf{q}_{i,0}^a$ and \mathbf{q}_0^b are set to be $\mathbf{0}_{L \times 1}$ and $\mathbf{0}_{N \times 1}$, respectively. The CTF filter and nonlinear coefficients are initialized as $\hat{\mathbf{a}}_{i,0} = \mathbf{0}_{L \times 1}$ and $\hat{\mathbf{b}}_0 = [1 \quad \mathbf{0}_{(N-1) \times 1}]^T$. The

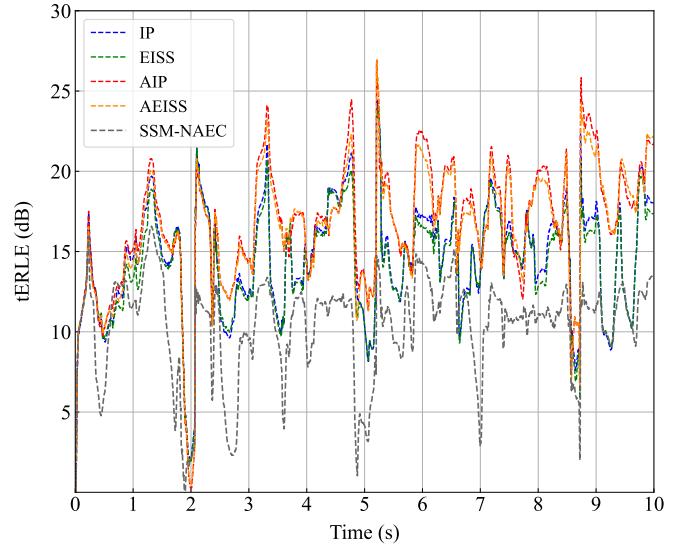


Fig. 5: tERLE performance achieved by all compared algorithms in a double-talk situation.

shape parameter β is set to $\beta = 0.4$ to be consistent with [65]. Note that with this configuration, AIP and AEISSL algorithms have similar steady-state performance as their conventional counterparts as illustrated in Sec. VI-A. To evaluate their performance, we use the echo return loss enhancement (ERLE) [31], [84] for the single-talk case and the true echo return loss enhancement (tERLE) [31], [39] for the double-talk case. The implementation of these two measures involves smoothing, utilizing samples within 0.2s vicinity. Besides, an experiment with real recordings is carried out where the Perceptual Evaluation of Speech Quality (PESQ) [85] and the Short-Time Objective Intelligibility (STOI) [86] metrics are used to measure the quality of the obtained near-end signals.

We compare the proposed algorithms with the original SBSS-NAEC algorithms and a single channel version of [32], i.e., a state-of-the-art state-space model-based NAEC (SSM-NAEC) algorithm.

A. Tracking Ability

A 20-second-long AR(1) signal with $x_{\max}=1$, generated by filtering a white Gaussian noise with the system $1/(1 - 0.8z^{-1})$, is utilized as the far-end signal. Subsequently, it undergoes the hard clipping function to produce the nonlinear loudspeaker signal. For the first 10 seconds, the nonlinearly distorted signal is convolved with the first AIR and it is then convolved with the second AIR. In other words, the AIR changes at $t=10$ s. We consider both single-talk and double-talk situations. For the double-talk case, the female speech signal from [65] is used as near-end signal and the corresponding signal-to-noise ratio (SNR) is set to 0 dB. A white Gaussian noise is added as background noise corresponding with a signal-to-noise ratio (SNR) of 60 dB. Note that we do not compare the tracking ability with SSM-NAEC as it cannot achieve comparable AEC performance with SBSS-NAEC algorithms in double-talk, which will be validated in Sec. (VI-C) and Sec. (VI-D). As seen in Fig. 3, AIP and

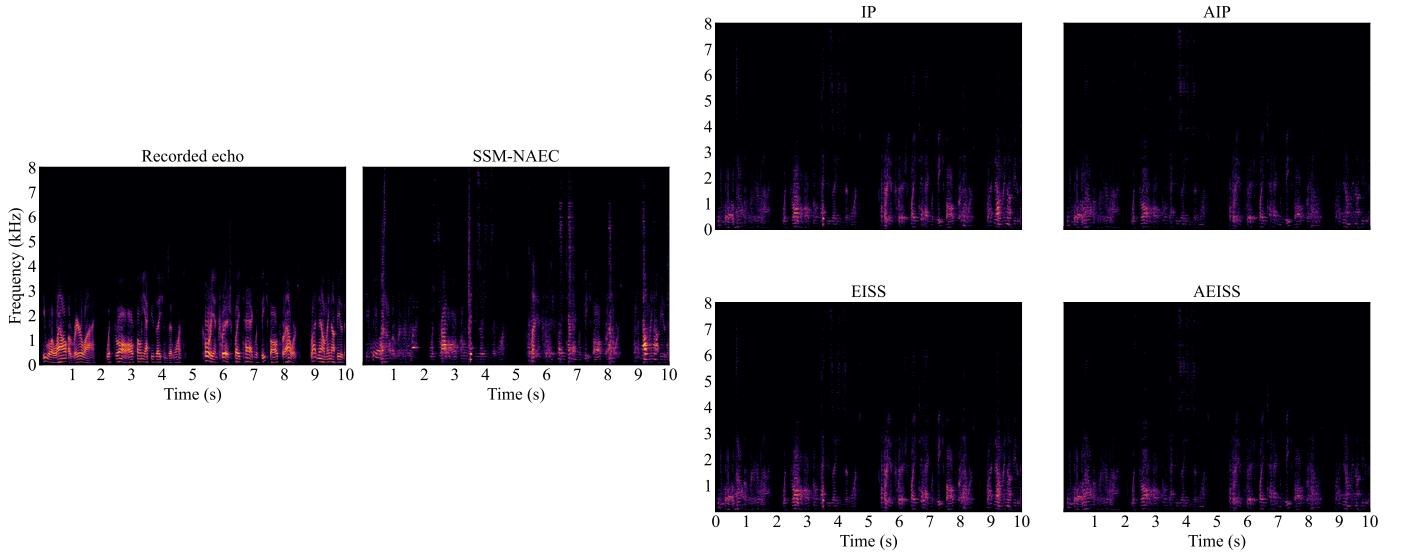


Fig. 6: Spectrograms of recorded echo and error signals generated by all compared algorithms.

AEISSL algorithms achieved similar steady-state performance with their conventional counterparts in the first 10 seconds. For the second 10 seconds, the proposed AIP and AEISSL have significantly better tracking ability since they have less number of free parameters to estimate.

B. Performance in Single-talk Case

Now, we compare the performance of all SBSS-NAEC algorithms during single talk using the data from [65]. A 10-second long male speech signal is used as the far-end signal. This signal is passed through the hard clipping function and is then convolved with the first AIR to generate the nonlinear echo. White Gaussian noise is added as background noise with an SNR of 60 dB. The ERLE performance achieved by all compared algorithms are plotted in Fig. 4. As can be seen, SSM-NAEC achieves similar performance as IP and EISSL. Moreover, the AIP and AEISSL algorithms yield much better performance in the single-talk case.

C. Performance in Double-talk Case

We now compare the performance of all compared algorithms in the double-talk case, again, using the data from [65]. The nonlinear echo is generated following the similar process as in Sec. VI-B. A female speech signal of 10-second long is used as the near-end signal with an SER of 0 dB. White Gaussian noise is added as background noise with an SNR of 60 dB. In Fig. 5, we plot the tERLE performance of all compared algorithms. Obviously, the SSM-NAEC algorithm cannot achieve comparable AEC performance compared with the SBSS-NAEC algorithms in double-talk situations. Moreover, the AIP and AEISSL algorithms yield much better performance than IP and EISSL, demonstrating the superiority of AIP and AEISSL in dealing with double-talk in AEC applications.

TABLE III: PESQ and STOI of obtained near-end signals with real recordings.

Algorithms	PESQ	STOI
Unprocessed	1.22	0.73
SSM-NAEC	1.51	0.84
IP	1.81	0.92
EISSL	1.77	0.91
AIP	2.15	0.95
AEISSL	2.09	0.95

D. Experiment with Real Recordings

Now, we compare the AEC performance of all compared algorithms with real recorded echoes. The aforementioned 10-second speech signal from a male speaker is played by a low-cost loudspeaker and picked up by a mobile phone in an office environment, sampled at 16 kHz. Subsequently, a 10-second speech signal from a female speaker is introduced as the near-end signal at an SER of 0 dB. It is important to note that we refrain from adding additional background noise, as the office noise is already captured in the recordings. We define the error signal as

$$e(t) = s(t) - \hat{s}(t). \quad (99)$$

The spectrograms of the error signals generated by all compared algorithms are shown in Fig. 6. To enhance clarity, the lower bound in the figure is set to 55 dB below the highest power in the data. As can be seen, the error signal generated by SSM-NAEC has much higher power than those generated by the SBSS-NAEC algorithms. Besides, SSM-NAEC generates some new components in high frequency bins, indicating that it caused distortion. Moreover, it is also observed that the power of the error signals generated by AIP and AEISSL is lower than those generated by their conventional counterparts. The obtained near-end signals are also evaluated with PESQ and STOI, which are shown in Tab. III. Remarkably, the proposed

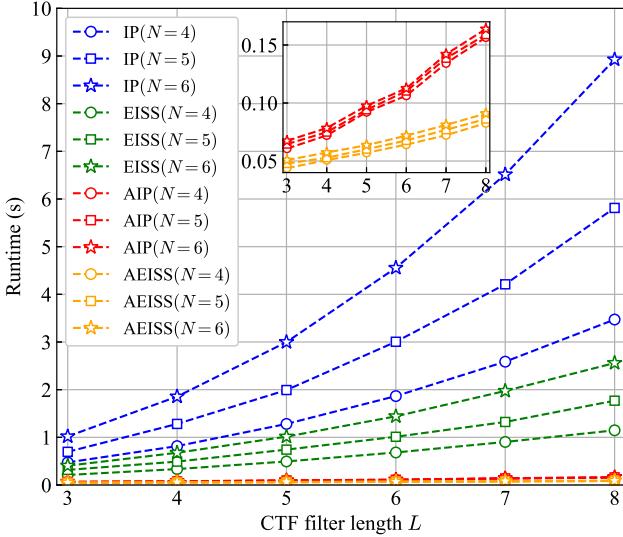


Fig. 7: The average runtime to process 1-second long signal with a 16 kHz sampling rate.

AIP and AEIIS algorithms achieve notably higher PESQ and STOI values compared to other methods, underscoring their ability to obtain high-quality near-end signals.

E. Runtime Comparison

Finally, we compare the runtime of SBSS-NAEC algorithms. Two 20-minute long white Gaussian noise signals are used as the far-end and near-end signals. We compare the computation time of the studied algorithms, which covers all the aforementioned three steps. The odd power series expansion order N is set, respectively, to 4, 5, and 6 while the CTF filter length L varies from 3 to 8. The average runtime for processing a 1-second long signal with all studied algorithms is plotted in Fig. 7. As seen, the runtime of all SBSS-NAEC algorithms increases with the value of N and L . Under the same configuration, EISS-based methods have less runtime than the IP-based methods. The proposed AIP and AEIIS algorithms are much more efficient than the conventional IP and EISS algorithms, and the difference is more significant as the value of L and N increases. This validates another property of the proposed AIP and AEIIS algorithms, i.e., besides being able to achieve better performance, it also has lower computational complexity as compared to their conventional counterparts.

VII. CONCLUSIONS

In this paper, we adopted a bilinear model to represent nonlinear echoes. To estimate the series expansion coefficients and the CTF filter in this model, we presented a bilinear alternating optimization framework. Under this framework, two algorithms, i.e., the AIP and AEIIS, were derived, both exploit the independence criteria to estimate the model parameters. We showed that the proposed AIP and AEIIS algorithms are capable to achieve nonlinear AEC in both the single-talk and double-talk scenarios. Since the bilinear representation consists of less parameters compared to a conventional

CTF model, the developed AIP and AEIIS algorithms have demonstrated interesting properties as compared to the conventional IP and EISS algorithms including improved NAEC performance, better tracking ability, and reduced complexity.

APPENDIX A DERIVATION OF THE AUXILIARY FUNCTION

With the definition of

$$\tilde{\sigma}_{s,j}^m = \sqrt{\sum_{i=1}^I |(\mathbf{w}_{i,j}^m)^H \tilde{\mathbf{y}}_{i,j}^m|^2}, \quad (100)$$

the log likelihood function can be expressed as

$$\mathcal{F}(\tilde{\sigma}_{s,j}^m) = -\log p(\mathbf{s}_j). \quad (101)$$

As \mathbf{s}_j follows a super Gaussian distribution, the following auxiliary function can be used [49]

$$\begin{aligned} \mathcal{F}^+(\tilde{\sigma}_{s,j}^m) &= \frac{\mathcal{F}'(\tilde{\sigma}_{s,j}^m)}{2\tilde{\sigma}_{s,j}^m} (\tilde{\sigma}_{s,j}^m)^2 \\ &\quad + \mathcal{F}(\sigma_{s,j}^m) - \frac{\sigma_{s,j}^m \mathcal{F}'(\sigma_{s,j}^m)}{2}, \end{aligned} \quad (102)$$

where $(\cdot)'$ represents the derivative and $\sigma_{i,j}^m$ is defined in (29). The above auxiliary function satisfies

$$\mathcal{F}^+(\tilde{\sigma}_{s,j}^m) \geq \mathcal{F}(\tilde{\sigma}_{s,j}^m). \quad (103)$$

with equality if and only if $\tilde{\sigma}_{s,j}^m = \sigma_{s,j}^m$, i.e., $\mathbf{w}_{i,j}^m = \mathbf{w}_{i,j-1}^m$. Therefore, instead of minimizing $\mathcal{F}(\tilde{\sigma}_{i,j}^m)$ at each time frame, we can minimize $\mathcal{F}^+(\tilde{\sigma}_{s,j}^m)$. Note that

$$\frac{\mathcal{F}'(\sigma_{s,j}^m)}{2\sigma_{s,j}^m} = \phi(\sigma_{s,j}^m), \quad (104)$$

which is defined in (28). Now substituting (100) and (103) into (25) and considering that $\det \mathbf{W}_{i,j}^m = 1$, we obtain (26), where we neglected irrelevant constant terms.

REFERENCES

- [1] A. Gilloire, E. Moulines, D. Slock, and P. Duhamel, “State of the art in acoustic echo cancellation,” in *Digital Signal Processing in Telecommunications*, A. R. Figueiras-Vidal, Eds., Berlin: Springer, 1996.
- [2] S. L. Gay and J. Benesty, Eds., *Acoustic Signal Processing for Telecommunication*. Boston, MA: Kluwer, 2000.
- [3] J. Benesty, T. Gänslor, D. R. Morgan, M. M. Sondhi, and S. L. Gay, *Advances in Network and Acoustic Echo Cancellation*. Berlin: Springer, 2001.
- [4] S. L. Gay, “Acoustic echo cancellation for duplex audio communication,” *J. Acoust. Soc. Am.*, vol. 99, no. 4, pp. 2504–2529, 1996.
- [5] J. Benesty, M. M. Sondhi, and Y. Huang, *Springer Handbook of Speech Processing*. New York, NY: Springer, 2007.
- [6] Y. Huang, J. Chen, and J. Benesty, “Immersive audio schemes,” *IEEE Signal Process. Mag.*, vol. 28, no. 1, pp. 20–32, Jan. 2011.
- [7] C. Paleologu, J. Benesty, and S. Ciochină, “Linear system identification based on a Kronecker product decomposition,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 26, no. 10, pp. 1793–1808, Oct. 2018.
- [8] X. Wang, G. Huang, J. Benesty, J. Chen, and I. Cohen, “Time difference of arrival estimation based on a Kronecker product decomposition,” *IEEE Signal Process. Lett.*, vol. 28, pp. 51–55, Dec. 2020.

- [9] S. Haykin, *Adaptive Filter Theory*. Englewood Cliffs, NJ: Prentice-Hall, 2002.
- [10] A. H. Sayed, *Fundamentals of Adaptive Filtering*. Hoboken, NJ: Wiley, 2003.
- [11] A. H. Sayed, *Adaptive Filters*. Hoboken, NJ: Wiley, 2008
- [12] C. Paleologu, J. Benesty, and S. Ciochină, “Linear system identification based on a Kronecker product decomposition,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 26, no. 10, pp. 1793–1808, Oct. 2018.
- [13] C.-L. Stanciu, J. Benesty, C. Paleologu, R.-L. Costea, L.-M. Dogariu, S. Ciochină, “Decomposition-based Wiener filter using the Kronecker product and conjugate gradient method,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 32, pp. 124–138, Oct. 2023.
- [14] E. R. Ferrara, “Fast implementation of LMS adaptive filters,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 28, no. 4, pp. 474–475, Aug. 1980.
- [15] G. Long, F. Ling, and J. G. Proakis, “The LMS algorithm with delayed coefficient adaptation,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 37, no. 9, pp. 1397–1405, Sep. 1989.
- [16] Y. Chen, Y. Gu, and A. O. Hero, “Sparse LMS for system identification,” in *Proc. IEEE ICASSP*, 2009, pp. 3125–3128.
- [17] G. Su, J. Jin, Y. Gu, and J. Wang, “Performance analysis of ℓ_0 -norm constraint least mean square algorithm,” *IEEE Trans. Signal Process.*, vol. 60, no. 5, pp. 2223–2235, May 2012.
- [18] A. E Albert and L. S. Gndrn Jr., *Stochastic Approximation and Nodinem Regreion*. Cambridge, MA: MIT Press, 1967.
- [19] J. I. Nagumo and A. Noda, “A learning method for system identification,” *IEEE Trans. Automat. Contr.*, vol. 12, no. 3, pp. 282–287, Jun. 1967.
- [20] D. L. Duttweiler, “Proportionate normalized least-mean-squares adaptation in echo cancelers,” *IEEE Trans. Speech Audio Process.*, vol. 8, no. 5, pp. 508–518, Sep. 2000.
- [21] C. Paleologu, J. Benesty, and S. Ciochină, “An improved proportionate NLMS algorithm based on the ℓ_0 norm,” in *Proc. IEEE ICASSP*, 2010, pp. 14–19.
- [22] V. Panuska, “An adaptive recursive-least-squares identification algorithm,” in *Proc. IEEE Symp. Adap. Process. Decis. Control.*, 1969, pp. 65–65.
- [23] C. Paleologu, J. Benesty, and S. Ciochină, “A robust variable forgetting factor recursive least-squares algorithm for system identification,” *IEEE Signal Process. Lett.*, vol. 15, pp. 597–600, Oct. 2008.
- [24] C. Elisei-IIiescu, C. Paleologu, J. Benesty, C. Stanciu, C. Anghel, and S. Ciochină, “Recursive least-squares algorithms for the identification of low-rank systems,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 5, pp. 903–918, May 2019.
- [25] L.-M. Dogariu, J. Benesty, C. Paleologu, and S. Ciochină, “Identification of room acoustic impulse responses via Kronecker product decompositions,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 30, no. 8, pp. 2828–2841, 2022.
- [26] G. Enzner and P. Vary, “Frequency-domain adaptive Kalman filter for acoustic echo control in hands-free telephones,” *Signal Process.*, vol. 86, no. 6, pp. 1140–1156, Jun. 2006.
- [27] C. Paleologu, J. Benesty, and S. Ciochină, “Study of the general Kalman filter for echo cancellation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 21, no. 8, pp. 1539–1549, Aug. 2013.
- [28] C. Paleologu, J. Benesty, and S. Ciochină, “Study of the optimal and simplified Kalman filters for echo cancellation,” in *Proc. IEEE ICASSP*, 2013, pp. 580–584.
- [29] J. Benesty, D. R. Morgan, and J. H. Cho, “A new class of doubletalk detectors based on cross-correlation,” *IEEE Trans. Speech, Audio Process.*, vol. 8, no. 2, pp. 168–172, Mar. 2000.
- [30] H. Buchner, J. Benesty, T. Gansler, and W. Kellermann, “Robust extended multidelay filter and double-talk detector for acoustic echo cancellation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 14, no. 5, pp. 1633–1644, Sep. 2006.
- [31] S. Malik and G. Enzner, “State-space frequency-domain adaptive filtering for nonlinear acoustic echo cancellation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 7, pp. 2065–2079, Sep. 2012.
- [32] J. Park and J.-H. Chang, “State-space microphone array nonlinear acoustic echo cancellation using multi-microphone nearend speech covariance,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 27, no. 10, pp. 1520–1534, Oct. 2019.
- [33] P. Comon, “Independent component analysis, a new concept?”, *Signal Process.*, vol. 36, no. 3, pp. 287–314, Apr. 1994.
- [34] A. Hyvärinen, J. Karhunen, and E. Oja, *Independent Component Analysis*. Hoboken, NJ: Wiley, 2001.
- [35] S. Makino, *Audio source separation*. Berlin: Springer, 2018.
- [36] D. W. E. L. Schobben and C. W. Sommen, “A frequency domain blind signal separation method based on decorrelation,” *IEEE Trans. Signal Process.*, vol. 50, no. 8, pp. 1855–1865, Aug. 2002.
- [37] H. Buchner and W. Kellermann, “A fundamental relation between blind and supervised adaptive filtering illustrated for blind source separation and acoustic echo cancellation,” in *Proc. Joint Workshop Hands-free Speech Communication and Microphone Arrays (HSCMA)*, 2008, pp. 17–20.
- [38] T. S. Wada and B.-H. Juang, “Acoustic echo cancellation based on independent component analysis and integrated residual echo enhancement,” in *Proc. IEEE WASPAA*, 2009, pp. 205–208.
- [39] F. Nesta, T. S. Wada, and B.-H. Juang, “Batch-online semi-blind source separation applied to multi-channel acoustic echo cancellation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 19, no. 3, pp. 583–599, Mar. 2011.
- [40] S. Araki, R. Mukai, S. Makino, T. Nishikawa, and H. Saruwatari, “The fundamental limitation of frequency domain blind source separation for convolutive mixtures of speech,” *IEEE Trans. Speech, Audio Process.*, vol. 11, no. 2, pp. 109–116, Mar. 2003.
- [41] Y. Avargel and I. Cohen, “System identification in the short time Fourier transform domain with crossband filtering,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, no. 4, pp. 1305–1319, May 2007.
- [42] R. Talmon, I. Cohen, and S. Gannot, “Relative transfer function identification using convolutive transfer function approximation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 17, no. 4, pp. 546–555, May 2009.
- [43] R. Talmon, I. Cohen, and S. Gannot, “Convulsive transfer function generalized sidelobe canceler,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 17, no. 7, pp. 1420–1434, Sep. 2009.
- [44] X. Wang, A. Brendel, G. Huang, Y. Yang, W. Kellermann, and J. Chen, “Spatially informed independent vector analysis for source extraction based on the convolutive transfer function model,” in *Proc. IEEE ICASSP*, 2023, pp. 1–5.
- [45] W. Kellermann, “Analysis and design of multirate systems for cancellation of acoustical echoes,” in *Proc. IEEE ICASSP*, 1988, pp. 2570–2573.
- [46] S. L. Gay and R. J. Mammone, “Fast converging subband acoustic echo cancellation using RAP on the WE DSP16A,” in *Proc. IEEE ICASSP*, 1990, pp. 1141–1144.
- [47] A. Yeredor, “On hybrid exact-approximate joint diagonalization,” in *Proc. CAMSAP*, 2009, pp. 312–315.
- [48] N. Ono, “Stable and fast update rules for independent vector analysis based on auxiliary function technique,” in *Proc. IEEE WASPAA*, 2011, pp. 189–192.
- [49] N. Ono, “Auxiliary-function-based independent vector analysis with power of vector-norm type weighting functions,” in *Proc. APSIPA ASC*, 2012, pp. 1–4.
- [50] D. Kitamura, N. Ono, H. Sawada, H. Kameoka, and H. Saruwatari, “Determined blind source separation unifying independent vector analysis and nonnegative matrix factorization,”

- IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 9, pp. 1626–1641, Sep. 2016.
- [51] Y. Kubo, N. Takamune, D. Kitamura, and H. Saruwatari, “Blind speech extraction based on rank-constrained spatial covariance matrix estimation with multivariate generalized Gaussian distribution,” *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 28, pp. 1948–1963, Jun. 2020.
- [52] R. Scheibler and N. Ono, “Fast and stable blind source separation with rank-1 updates,” in *Proc. IEEE ICASSP*, 2020, pp. 236–240.
- [53] R. Scheibler, “Independent vector analysis via log-quadratically penalized quadratic minimization,” *IEEE Trans. Signal Process.*, vol. 69, pp. 2509–2524, Apr. 2021.
- [54] T. Nakashima and N. Ono, “Inverse-free online independent vector analysis with flexible iterative source steering,” in *Proc. APSIPA ASC*, 2022, pp. 749–753.
- [55] T. Nakashima, R. Ikeshita, N. Ono, S. Araki, and T. Nakatani, “Fast online source steering algorithm for tracking single moving source using online independent vector analysis,” in *Proc. IEEE ICASSP*, 2023, pp. 1–5.
- [56] Z. Wang, Y. Na, Z. Liu, B. Tian, and Q. Fu, “Weighted recursive least square filter and neural network based residual echo suppression for the AEC-challenge,” in *Proc. IEEE ICASSP*, 2021, pp. 141–145.
- [57] G. Cheng, L. Liao, H. Chen, and J. Lu, “Semi-blind source separation for nonlinear acoustic echo cancellation,” *IEEE Signal Process. Lett.*, vol. 28, pp. 474–478, 2021.
- [58] K. Lu, X. Wang, T. Ueda, S. Makino, and J. Chen, “A computationally efficient semi-blind source separation approach for nonlinear echo cancellation based on an element-wise iterative source steering,” 2023, *arXiv:2004.03926v1*.
- [59] T. Koh and E. Powers, “Second-order Volterra filtering and its application to nonlinear system identification,” *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 33, no. 6, pp. 1445–1455, Dec. 1985.
- [60] L. Tan and J. Jiang, “Adaptive Volterra filters for active control of nonlinear noise processes,” *IEEE Trans. Signal Process.*, vol. 49, no. 8, pp. 1667–1676, Aug. 2001.
- [61] A. Guerin, G. Faucon, and R. Le Bouquin-Jeannes, “Nonlinear acoustic echo cancellation based on Volterra filters,” *IEEE Trans. Speech Audio Process.*, vol. 11, no. 6, pp. 672–683, Nov. 2003.
- [62] A. Stenger and W. Kellermann, “Adaptation of a memoryless preprocessor for nonlinear acoustic echo cancelling,” *Signal Process.*, vol. 80, no. 9, pp. 1747–1760, Sep. 2000.
- [63] F. Kuech, A. Mitnacht, and W. Kellermann, “Nonlinear acoustic echo cancellation using adaptive orthogonalized power filters,” in *Proc. IEEE ICASSP*, 2005, pp. 105–108.
- [64] F. Kuech and W. Kellermann, “Orthogonalized power filters for nonlinear acoustic echo cancellation,” *Signal Process.*, vol. 86, no. 6, pp. 1168–1181, Jun. 2006.
- [65] G. Cheng, L. Liao, K. Chen, Y. Hu, C. Zhu, and J. Lu, “Semi-blind source separation using convolutive transfer function for nonlinear acoustic echo cancellation,” *J. Acoust. Soc. Am.*, vol. 153, no. 1, pp. 88–95, Jan. 2023.
- [66] J. Benesty, C. Paleologu, and S. Ciochină, “On the identification of bilinear forms with the Wiener filter,” *IEEE Signal Process. Lett.*, vol. 24, no. 5, pp. 653–657, May 2017.
- [67] L. Dogariu, C. Paleologu, S. Ciochină, J. Benesty, and P. Piantanida, “Identification of bilinear forms with the Kalman filter,” in *Proc. IEEE ICASSP*, 2018, pp. 4134–4138.
- [68] M. Z. Ikram, “Non-linear acoustic echo cancellation using cascaded Kalman filtering,” in *Proc. IEEE ICASSP*, 2014, pp. 1320–1324.
- [69] D. R. Hunter and K. Lange, “A tutorial on MM algorithms,” *Amer. Statist.*, vol. 58, no. 1, pp. 30–37, 2004.
- [70] H. L. Van Trees, “*Optimum Array Processing: Part IV of Detection, Estimation and Modulation Theory*”. New York: Wiley, 2002.
- [71] S. Boyd and L. Vandenberghe, *Convex Optimization*. Cambridge, U.K.: Cambridge Univ. Press, 2004.
- [72] D. P. Palomar and Y. C. Eldar, *Convex Optimization in Signal Processing and Communications*. Cambridge, U.K.: Cambridge Univ. press, 2010.
- [73] J. Capon, “High resolution frequency-wavenumber spectrum analysis,” *Proc. IEEE*, vol. 57, no. 8, pp. 1408–1418, Aug. 1969.
- [74] X. Wang, J. Benesty, G. Huang, and J. Chen, “A minimum variance distortionless response spectral estimator with Kronecker product filters”, in *Proc. EUSIPCO*, 2022, pp. 2261–2265.
- [75] O. L. Frost, “An algorithm for linearly constrained adaptive array processing,” *Proc. IEEE*, vol. 60, no. 8, pp. 926–935, Aug. 1972.
- [76] J. Dmochowski, J. Benesty, and S. Affes, “Linearly constrained minimum variance source localization and spectral estimation,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 16, no. 8, pp. 1490–1502, Nov. 2008.
- [77] K. M. Buckley and L. J. Griffiths, “An adaptive generalized sidelobe canceller with derivative constraints,” *IEEE Trans. Antennas Propagat.*, vol. 34, no. 3, pp. 311–319, Mar. 1986.
- [78] H. V. Henderson and S. R. Searle, “On deriving the inverse of a sum of matrices,” *SIAM Rev.*, vol. 23, no. 1, pp. 53–60, Jan. 1981.
- [79] R. A. Horn and C. R. Johnson, *Matrix Analysis*. Cambridge, U.K.: Cambridge univ. Press, 1985.
- [80] F. Zhang, *The Schur Complement and Its Applications*. New York, NY: Springer, 2005.
- [81] V. Panayotov, G. Chen, D. Povey, and S. Khudanpur, “Librispeech: An ASR corpus based on public domain audio books,” in *Proc. IEEE ICASSP*, 2015, pp. 5206–5210.
- [82] S. Nakamura, K. Hiyane, F. Asano, T. Nishiura, and T. Yamada, “Acoustical sound database in real environments for sound sceneunderstanding and hands-free speech recognition,” in *Proc. Lang. Resour. Eval.*, 2000, pp. 965–996.
- [83] R. Cutler, A. Saabas, T. Parnamaa, M. Purin, H. Gamper, S. Braun, K. Sørensen, and R. Aichner, “ICASSP 2022 acoustic echo cancellation challenge,” in *Proc. IEEE ICASSP*, 2022, pp. 9107–9111.
- [84] M. Zeller and W. Kellermann, “Fast and robust adaptation of DFT-domain Volterra filters in diagonal coordinates using iterated coefficient updates,” *IEEE Trans. Signal Process.*, vol. 58, no. 3, pp. 1589–1604, Mar. 2010.
- [85] *Mapping Function for Transforming Raw Results Scores to MOS-LQO*, ITU-T Rec. P. 862.1, 2003.
- [86] C. Taal, R. Hendriks, R. Heusdens, and J. Jensen, “An algorithm for intelligibility prediction of time-frequency weighted noisy speech,” *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 19, pp. 2125–2136, Sep. 2011.



Xianrui Wang received the Bachelor's degree in Control Science and Technology from Northwestern Polytechnical University, Xi'an, China, in 2019. He is currently working toward the Ph.D. degree in information and communication engineering with the Center of Intelligent Acoustics and Immersive Communications in Northwestern Polytechnical University, China. He is currently a visiting research fellow in Waseda university, Kitakyushu, Japan. His research interests include blind and semi-blind source separation and extraction, acoustic echo cancellation, speech dereverberation, and source localization. He is an active reviewer of IEEE Signal Processing Letters, IEEE Wireless Communications Letters, Sensors and Actuators A: Physical, and IEEE ICASSP.



Yichen Yang (Graduate Student Member, IEEE) received the Bachelor's and Master's degrees in 2018 and 2021, respectively, from Northwestern Polytechnical University, Xi'an, China. He is currently a Ph.D. student of information and communication engineering with the Center of Intelligent Acoustics and Immersive Communications, Northwestern Polytechnical University, Xi'an, China, and a Visiting Ph.D. student with Waseda University, Kitakyushu, Japan. His research interests include array signal processing, speech enhancement, blind source separation, and machine learning.



Andreas Brendel received the B.Sc. degree and the M.Sc. degree (with distinction) in electrical engineering from the Friedrich-Alexander Universität Erlangen-Nürnberg, Germany, in 2014 and 2016, respectively. From 2016 to 2022, he was with the chair of Multimedia Communications and Signal Processing of the Friedrich-Alexander Universität Erlangen-Nürnberg where he received the Dr.-Ing. degree (with distinction) in 2022. He was a visiting scientist at Tsukuba University, Japan (Sep.-Nov. 2016) and at Bar Ilan University, Israel (Sep.-Oct. 2018). In 2022, Dr. Brendel joined the Fraunhofer Institute for Integrated Circuits IIS in Erlangen as an associated researcher. He authored or coauthored more than 50 peer-reviewed papers in journals and conference proceedings as well as several patents. His research interests include acoustic source separation and deep generative modeling for speech synthesis and coding.



Tetsuya Ueda (Student Member, IEEE) received the B.Sc. and M.E. degrees in information engineering and engineering from the University of Tsukuba, Japan, in 2020 and 2022, respectively. He is currently working toward the Ph.D. degree with Waseda University, Kitakyushu, Japan. His research interests include acoustic signal processing, speech enhancement, and dereverberation.



Jacob Benesty received a Master degree in microwaves from Pierre & Marie Curie University, France, in 1987, and a Ph.D. degree in control and signal processing from Paris-Saclay University, France, in April 1991. During his Ph.D. (from Nov. 1989 to Apr. 1991), he worked on adaptive filters and fast algorithms at the Centre National d'Etudes des Telecommunications (CNET), Paris, France. From January 1994 to July 1995, he worked at Telecom Paris University on multichannel adaptive filters and acoustic echo cancellation. From October 1995 to May 2003, he was first a consultant and then a Member of the Technical Staff at Bell Laboratories, Murray Hill, NJ, USA. In May 2003, he joined the University of Quebec, INRS-EMT, in Montreal, Quebec, Canada, as a professor. He is also an Adjunct Professor with Aalborg University, Denmark, and a Guest Professor with Northwestern Polytechnical University, Xi'an, China. His research interests are in signal processing, acoustic signal processing, and multimedia communications. He is the inventor of many important technologies. In particular, he was the lead researcher at Bell Labs who conceived and designed the world-first real-time hands-free full-duplex stereophonic teleconferencing system. Also, he conceived and designed the world-first PC-based multi-party hands-free full-duplex stereo conferencing system over IP networks. He is the editor of the book series Springer Topics in Signal Processing. He was the general chair and technical chair of many international conferences and a member of several IEEE technical committees. Four of his journal papers were awarded by the IEEE Signal processing Society, in 2010 he received the Gheorghe Cartianu Award from the Romanian Academy, and in 2023 he received an Honorary Doctorate (Doctor Techniques Honoris Causa) from Aalborg University, Denmark, for his distinguished efforts in audio and acoustic signal processing. He has co-authored and co-edited/co-authored numerous books in the area of acoustic signal processing.



Walter Kellermann (Life Fellow, IEEE) is a professor for communications at the University of Erlangen-Nuremberg, Germany, since 1999. He received the Dipl.-Ing. (univ.) degree in Electrical Engineering from the University of Erlangen-Nuremberg, in 1983, and the Dr.-Ing. degree from the Technical University Darmstadt, Germany, in 1988. From 1989 to 1990, he was a postdoctoral Member of Technical Staff at AT&T Bell Laboratories, Murray Hill, NJ. In 1990, he joined Philips Kommunikations Industrie, Nuremberg, Germany, to work on hands-free communication in cars. From 1993 to 1999, he was a Professor at the Fachhochschule Regensburg, where he also became Director of the Institute of Applied Research in 1997. In 1999, he cofounded DSP Solutions, a consulting firm in digital signal processing, and he joined the University Erlangen-Nuremberg as a Professor and Head of the Audio Research Laboratory. He authored or coauthored 20+ book chapters, 400+ refereed papers in journals and conference proceedings, as well as 80+ patents, and is a co-recipient of ten best paper awards. His current research interests include speech signal processing, array signal processing and machine learning, especially for acoustic signal processing. Dr. Kellermann served as an Associate Editor and Guest Editor to various journals, including the IEEE Transactions on Speech and Audio Processing (2000 to 2004), the IEEE Signal Processing Magazine (2015), the IEEE Journal on Special Topics in Signal Processing (2019) and presently serves as Associate Editor to the EURASIP Journal on Applied Signal Processing. He was the General Chair of eight mostly IEEE-sponsored workshops and conferences. He served as a Distinguished Lecturer of the IEEE Signal Processing Society (SPS) from 2007 to 2008. He was the Chair of the IEEE SPS Technical Committee for Audio and Acoustic Signal Processing from 2008 to 2010, a Member of the IEEE James L. Flanagan Award Committee from 2011 to 2014, a Member of the SPS Board of Governors (2013-2015), Vice President Technical Directions of the IEEE Signal Processing Society (2016-2018), a Member of the SPS Nominations Appointments Committee (2019-2022) and currently serves as a Member of the SPS Fellow Evaluation Committee and as Chair of the EURASIP Fellow Evaluation Committee. He was awarded the Julius von Haast Fellowship by the Royal Society of New Zealand in 2012 and the Group Technical Achievement Award of the European Association for Signal Processing (EURASIP) in 2015. In 2016, he was a Visiting Fellow at Australian National University, Canberra, Australia. He was elevated to EURASIP Fellow in 2021 and is an IEEE Life Fellow.



Shoji Makino (Life Fellow, IEEE) received the B.E., M.E., and Ph.D. degrees from Tohoku University, Sendai, Japan, in 1979, 1981, and 1993, respectively. He joined NTT in 1981 and the University of Tsukuba, Ibaraki, Japan, in 2009. He is currently a Professor with Waseda University, Kitakyushu, Japan. He has authored or coauthored of more than 400 articles in journals and conference proceedings and is responsible for more than 200 patents. His research interests include adaptive filtering technologies, blind source separation of convolutive mixtures

of speech, the realization of acoustic echo cancellation, and acoustic signal processing for speech and audio applications. He was the recipient of the IEEE Signal Processing Society Leo L. Beranek Meritorious Service Award in 2022, the ICA Unsupervised Learning Pioneer Award in 2006, the IEEE MLSP Competition Award in 2007, the IEEE SPS Best Paper Award in 2014, the Achievement Award for Science and Technology from the Japanese Government in 2015, the Hoko Award of the Hattori Hokokai Foundation in 2018, the Honorary Member Award of the IEICE in 2022, the Outstanding Contribution Award of the IEICE in 2018, the Technical Achievement Award of the IEICE in 2017 and 1997, the Outstanding Technological Development Award of the ASJ in 1995, and 8 best paper awards. He was a member of the IEEE Jack S. Kilby Signal Processing Medal Committee (2015–2018) and the James L. Flanagan Speech & Audio Processing Award Committee (2008–2011). He was on IEEE SPS Board of Governors (2018–2020), Technical Directions Board (2013–2014), Awards Board (2006–2008), Conference Board (2002–2004), and Fellow Evaluation Committee (2018–2020). He was a Keynote Speaker at ICA 2007, a Tutorial Speaker at ICASSP 2007, Interspeech 2011, and EMBC 2013. He was an Associate Editor for the IEEE TRANSACTIONS ON SPEECH AND AUDIO PROCESSING (2002–2005) and an Associate Editor for the EURASIP Journal on Advances in Signal Processing (2005–2012). He was the Guest Editor of the Special Issue of the IEEE Signal Processing Magazine (2013–2014). He was the Chair of SPS Audio and Acoustic Signal Processing Technical Committee (2013–2014) and the Chair of the Blind Signal Processing Technical Committee of the IEEE Circuits and Systems Society (2009–2010). He was the General Chair of IWAENC 2018, WASPAA 2007, IWAENC 2003, the Organizing Chair of ICA 2003, and is the designated Plenary Chair of ICASSP 2012. Dr. Makino is an IEEE SPS Distinguished Lecturer (2009–2010), an IEICE Fellow, a Board member of the ASJ, and a member of EURASIP.



Jingdong Chen (Fellow, IEEE) received the Ph.D. degree in pattern recognition and intelligence control from the Chinese Academy of Sciences, Beijing, China, in 1998. He is currently a Professor with Northwestern Polytechnical University, Xi'an, China. Prior to this position, he worked at Bell Laboratories, Murray Hill, New Jersey, WeVoice Inc., New Jersey, Griffith University, Brisbane, Australia, and Advanced Telecommunication Research Institute International (ATR), Kyoto, Japan, for more than a decade. His research interests include array

signal processing, adaptive signal processing, speech enhancement, adaptive noise/echo control, signal separation, speech communication, and artificial intelligence. He served as an Associate Editor for the IEEE Trans. Audio, Speech, Lang. Process. from 2008 to 2014, as a Technical Committee (TC) Member of the IEEE Signal Processing Society (SPS) TC on Audio and Electroacoustics from 2007 to 2009 and a member of the IEEE SPS TC on Audio and Acoustic Signal Processing from 2018 to 2021. He is currently serving as the Chair of IEEE Xi'an Section. He was the General Co-Chair of ACM WUWNET 2018 and IWAENC 2016, the Technical Program Chair of IEEE TENCON 2013, a Technical Program Co-Chair of IEEE WASPAA 2009, IEEE ChinaSIP 2014, IEEE ICSPCC 2014, and IEEE ICSPCC 2015, and helped organize many other conferences. Dr. Chen was the recipient of the 2008 Best Paper Award from the IEEE Signal Processing Society (with Benesty, Huang, and Doclo), the Best Paper Award from the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics in 2011 (with Benesty), the Bell Labs Role Model Teamwork Award twice, respectively, in 2009 and 2007, the NASA Tech Brief Award twice, respectively, in 2010 and 2009, and the Young Author Best Paper Award from the 5th National Conference on Man-Machine Speech Communications in 1998. He is a co-author of a paper for which C. Pan was the recipient of the IEEE R10 (Asia-Pacific Region) Distinguished Student Paper Award (First Prize) in 2016. He was also the recipient of the Japan Trust International Research Grant from the Japan Key Technology Center in 1998 and the Distinguished Young Scientists.