

Distributed combined acoustic echo cancellation and noise reduction in wireless acoustic sensor and actuator networks

Santiago Ruiz, Toon van Waterschoot and Marc Moonen

Abstract—The paper presents distributed algorithms for combined acoustic echo cancellation (AEC) and noise reduction (NR) in a wireless acoustic sensor and actuator network (WASAN) where each node may have multiple microphones and multiple loudspeakers, and where the desired signal is a speech signal. A centralized integrated AEC and NR algorithm, i.e., multichannel Wiener filter (MWF), is used as starting point where echo signals are viewed as background noise signals and loudspeaker signals are used as additional input signals to the algorithm. By including prior knowledge (PK), namely that the loudspeaker signals do not contain any desired signal component, an alternative centralized cascade algorithm (PK-MWF) is obtained with an AEC stage first followed by an MWF-based NR stage. Distributed algorithms can then be obtained from the MWF and PK-MWF algorithm, i.e., the GEVD-DANSE and PK-GEVD-DANSE algorithm, respectively. In the former, each node performs a reduced dimensional integrated AEC and NR algorithm and broadcasts only 1 fused signal (instead of all its signals) to the other nodes. In the PK-GEVD-DANSE algorithm, each node performs a reduced dimensional cascade AEC and NR algorithm and broadcasts only 2 fused signals (instead of all its signals) to the other nodes. The distributed algorithms achieve the same performance as the corresponding centralized integrated (MWF) and cascade (PK-MWF) algorithm. It is observed, however, that the communication cost in the PK-GEVD-DANSE algorithm can be reduced, where each node then broadcasts only 1 fused signal (instead of 2 signals) to the other nodes, which finally results in an algorithm with a low communication cost as well as a low computational complexity in each node.

Index Terms—Distributed signal processing, wireless acoustic sensor and actuator networks, acoustic echo cancellation, noise reduction, multichannel Wiener filter

I. INTRODUCTION

MANY speech and audio signal processing applications, such as teleconferencing/telepresence, in-car communication and ambient intelligence, suffer from acoustic echoes and background noise which corrupt the desired audio signal. Acoustic echo cancellation (AEC) and noise reduction (NR)

This research work was carried out at the ESAT Laboratory of KU Leuven, in the frame of Research Council KU Leuven Project C3-19-00221 ‘‘Cooperative Signal Processing Solutions for IoT-based Multi-User Speech Communication Systems’’, VLAIO O&O Project nr. HBC.2020.2197 ‘SPIC: Signal Processing and Integrated circuits for Communications’, Fonds de la Recherche Scientifique - FNRS and the Fonds Wetenschappelijk Onderzoek - Vlaanderen under EOS Project no 30452698 ‘(MUSE-WINET) MUlti-SERVICE WIREless NETWORK’ and the European Research Council under the European Union’s Horizon 2020 research and innovation program / ERC Consolidator Grant: SONORA (no. 773268). This paper reflects only the authors’ views and the Union is not liable for any use that may be made of the contained information. The scientific responsibility is assumed by its authors.

techniques can be used to enhance the desired signal while reducing undesired signal components [1]–[4].

Solutions to combined AEC and NR have been presented in the literature, which fundamentally can be divided into cascade and integrated approaches [3]–[8]. A cascade approach consists of an AEC stage and an NR stage which can be combined in two ways, i.e., a multichannel AEC stage followed by a multichannel NR stage, or a single-channel AEC stage preceded by a multichannel NR stage. The order of the stages has performance implications on the combined system. The first combination requires an AEC that is robust against noise in the microphone signals. In the second combination the AEC stage receives a noise-reduced signal which may contain a far-end signal component, therefore the AEC stage should be able to track changes in the acoustic environment as well as in the NR filters. Integrated approaches aim to solve the problem by combining the AEC and NR tasks in a single optimization process [5], [7], [9].

Recently, a multichannel Kalman-based Wiener filter for speaker interference reduction was proposed in [10]. The filter is based on a multichannel AEC stage followed by a NR stage using a multichannel Wiener filter. The proposed method was developed and implemented for a specific set-up with three speakers. Combined AEC and NR was implemented using a Kalman filter for a single-channel scenario in [11]. The use of deep neural networks to solve combined AEC and NR has also gained significant attention [12]–[14]. Although these methods usually outperform model-based methods, their main drawback is their dependency on training sets, which limits their practical deployment in mobile devices [12].

Existing solutions are all based on centralized processing, which is usually prohibitive in a wireless acoustic sensor and actuator network (WASAN) in terms of complexity and communication cost [15]. Distributed algorithms have been developed to overcome this, such as, e.g., the distributed delay-and-sum beamformer for NR based on randomized gossiping presented in [16], which was extended to a distributed MVDR beamformer based on message passing in [17]. Both algorithms do not have a topology constraint and provide good performance at the expense of a high communication cost [16]. The distributed adaptive node-specific signal estimation (DANSE) algorithm as developed in [18], performs distributed NR, i.e., optimally enhances the desired signal component in the local microphone signals of each node. It achieves a performance as if all microphone signals in the network were available to each and every node, while each node is

still sharing only a fused version of its microphone signals with the other nodes. A combination of a neural network and beamforming was used in [19] for a real-time multi-channel speech enhancement algorithm, where a spectral mask estimation is performed via the deep neural network together with spatial filtering. All these distributed algorithms only consider NR.

In this paper, distributed algorithms for combined AEC and NR are presented, where each node may have multiple microphones and multiple loudspeakers, and where the desired signal is a speech signal. In a WASAN with K nodes, node $k \in \mathcal{K} = \{1, \dots, K\}$ contains m_k microphones and l_k loudspeakers. The loudspeakers play given (far-end) signals, and generate echo signals in the microphones (also in other nodes). Node k then has access to an $n_k = m_k + P \cdot l_k$ vector signal, where $P - 1$ will be defined as the order of the interframe filtering in the AEC stage in Section II. The total number of microphones and loudspeakers in the WASAN are denoted, respectively, by $M = \sum_{k=1}^K m_k$ and $L = \sum_{k=1}^K l_k$, and similarly, $N = \sum_{k=1}^K n_k$. Centralized, non-cooperative and distributed algorithms can be used for combined AEC and NR, where the following should be considered: A centralized cascade algorithm has an AEC stage with PL AEC filter input signals, and a NR stage with M channels. A non-cooperative cascade algorithm for node k (i.e. node k working in isolation) has an AEC stage with Pl_k AEC filter input signals, and a NR stage with m_k channels. A distributed algorithm aims to reduce computational complexity by performing local operations in each node and exchanging data with other nodes.

In [20] distributed combined AEC and NR was considered in a WASAN. Essentially, a centralized integrated algorithm, i.e., the multichannel Wiener filter (MWF), is first turned into an alternative centralized cascade algorithm by introducing prior knowledge (PK). In the MWF algorithm no distinction is made between loudspeaker and microphone signals, which means echo signals are viewed as additional background noise signals and loudspeaker signals are used as additional input signals to the algorithm. By including PK, namely that the loudspeaker signals do not contain any desired signal component, the MWF algorithm is turned into the PK-MWF algorithm, leading to the alternative centralized cascade algorithm, with an AEC stage first followed by an MWF-based NR stage. The resulting algorithm has a lower computational complexity and allows to substitute alternative algorithms in the AEC stage.

Both the MWF and PK-MWF algorithm can be turned into a distributed algorithm, namely the generalized eigenvalue decomposition (GEVD)-based DANSE (GEVD-DANSE) [18] and the PK-GEVD-DANSE [21]. In the GEVD-DANSE algorithm, each node in the network performs a reduced dimensional (dimension $n_k + K - 1$) integrated AEC and NR algorithm and broadcasts only 1 fused signal (instead of n_k signals) to the other nodes, and yet each node achieves the same performance as the centralized integrated algorithm, i.e., as if all loudspeaker and microphone signals were broadcast in the network. In the PK-GEVD-DANSE algorithm, each node in the network performs a reduced dimensional (dimension $n_k + 2(K - 1)$) cascade AEC and NR algorithm and broadcasts

only 2 fused signals (instead of n_k signals) to the other nodes, and yet each node again achieves the same performance as the centralized cascade algorithm.

The PK-GEVD-DANSE algorithm performs AEC and NR in each node based on sharing not only fused microphone and loudspeaker signals between the nodes, which act as desired signal references, but also fused loudspeaker signals, which act as noise references. In this paper, however, it will be shown that in an AEC context (unlike in the general PK-GEVD-DANSE context) there is no need for sharing noise references between the nodes, reducing the communication cost in the PK-GEVD-DANSE algorithm. Each node then effectively performs a reduced dimensional (dimension $n_k + K - 1$) cascade AEC and NR algorithm and broadcasts only 1 fused signal (instead of 2 signals) to the other nodes. It will be shown that, this PK-GEVD-DANSE algorithm again achieves a performance as if all signals were available to each and every node. Implementations of the PK-GEVD-DANSE algorithm are presented using the normalized least squares (NLMS) algorithm and QR decomposition based recursive least squares (QRD-RLS) algorithm in the AEC stage. Furthermore, monitoring of the loudspeaker activity by means of a voice activity detector (VAD) is proposed.

The paper is organized as follows. The data model is presented in Section II. The formulations for the centralized integrated and cascade algorithm are provided in Sections III and IV. The distributed integrated and cascade algorithm are described in Sections V and VI. Section VII describes the NLMS- and QRD-RLS-based algorithm in the AEC stage of the PK-GEVD-DANSE algorithm. Simulations are shown in Section VIII, and finally Section IX concludes the paper.

II. PROBLEM FORMULATION AND NOTATION

Consider a fully connected WASAN with K nodes (see Fig. 1), where node $k \in \mathcal{K} = \{1, \dots, K\}$ contains m_k microphones and l_k loudspeakers, and hence has access to the short-time Fourier transform (STFT) domain $n_k \times 1$ signal vector $\mathbf{y}_k(\kappa, l) = \begin{bmatrix} \mathbf{x}_k(\kappa, l) \\ \mathbf{u}_k(\kappa, l) \end{bmatrix}$, where κ is the frequency bin index, l the frame index (for brevity κ and l will be omitted in the following, except for a few cases where l has to be included explicitly) and, $n_k = m_k + Pl_k$. Vector \mathbf{u}_k contains l_k local loudspeaker signals sampled at the current and previous $P - 1$ frames, i.e.,

$$\mathbf{u}_k(l) = \begin{bmatrix} u_1(l) \\ \vdots \\ u_1(l - P + 1) \\ \vdots \\ u_{l_k}(l) \\ \vdots \\ u_{l_k}(l - P + 1) \end{bmatrix}. \quad (1)$$

Vector \mathbf{x}_k contains m_k local microphone signals sampled only at the current frame and is modeled as

$$\mathbf{x}_k = \mathbf{s}_k + \mathbf{n}_k = \mathbf{a}_k s + \mathbf{n}_k. \quad (2)$$

Here, s is the desired speech source signal (also known as the dry signal), \mathbf{a}_k contains the acoustic transfer functions from the desired speech source position to the local microphones, \mathbf{s}_k is the desired speech component and \mathbf{n}_k is the noise component in the microphone signals of node k , modeled as

$$\mathbf{n}_k = \mathbf{G}_{kk} \mathbf{u}_k + \sum_{q \neq k} \mathbf{G}_{kq} \mathbf{u}_q + \mathbf{b}_k \quad (3)$$

where \mathbf{G}_{kk} is an $m_k \times Pl_k$ matrix representing the local echo paths from the local loudspeakers to the local microphones, \mathbf{G}_{kq} is an $m_k \times Pl_q$ matrix representing the echo paths from the loudspeakers in node q to the microphones in node k and finally \mathbf{u}_q contains the loudspeaker signals from node q . The background noise is assumed to be stationary with correlation matrix

$$\bar{\mathbf{R}}_{\mathbf{b}_k \mathbf{b}_k} = E\{\mathbf{b}_k \mathbf{b}_k^H\} \quad (4)$$

where $(\cdot)^H$ denotes the conjugate transpose operator and $E\{\cdot\}$ is the expected value operator. The following vectors are also defined,

$$\tilde{\mathbf{s}}_k = [\mathbf{s}_k^H \mathbf{0}_{1 \times Pl_k}]^H \quad (5)$$

$$\tilde{\mathbf{n}}_k = [\mathbf{n}_k^H \mathbf{u}_k^H]^H \quad (6)$$

$$\tilde{\mathbf{a}}_k = [\mathbf{a}_k^H \mathbf{0}_{1 \times Pl_k}]^H \quad (7)$$

$$\tilde{\mathbf{b}}_k = [\mathbf{b}_k^H \mathbf{0}_{1 \times Pl_k}]^H, \quad (8)$$

where $\mathbf{0}_{1 \times Pl_k}$ is a Pl_k -dimensional all-zero vector, and so that

$$\mathbf{y}_k = \tilde{\mathbf{s}}_k + \tilde{\mathbf{n}}_k = \tilde{\mathbf{a}}_k s + \tilde{\mathbf{n}}_k. \quad (9)$$

The N -dimensional vectors ($N = \sum_{k=1}^K n_k$), \mathbf{y} , \mathbf{s} , \mathbf{n} , \mathbf{a} and \mathbf{b} are the stacked versions of \mathbf{y}_k , $\tilde{\mathbf{s}}_k$, $\tilde{\mathbf{n}}_k$, $\tilde{\mathbf{a}}_k$ and $\tilde{\mathbf{b}}_k$ respectively, such that the signal vector \mathbf{y} can be characterized as follows

$$\mathbf{y} = \mathbf{s} + \mathbf{n} = \mathbf{a} \mathbf{s} + \mathbf{n}. \quad (10)$$

Assuming that the desired speech source signal and background noise are uncorrelated, and uncorrelated with the loudspeaker signals, correlation matrices can be defined as follows

$$\bar{\mathbf{R}}_{\mathbf{y} \mathbf{y}} = E\{\mathbf{y} \mathbf{y}^H\} = E\{\mathbf{s} \mathbf{s}^H\} + E\{\mathbf{n} \mathbf{n}^H\} = \bar{\mathbf{R}}_{\mathbf{s} \mathbf{s}} + \bar{\mathbf{R}}_{\mathbf{n} \mathbf{n}} \quad (11)$$

$$\bar{\mathbf{R}}_{\mathbf{s} \mathbf{s}} = \mathbf{a} \phi_s \mathbf{a}^H \quad (12)$$

$$\bar{\mathbf{R}}_{\mathbf{n} \mathbf{n}} = \mathbf{G} \Phi_{\mathbf{u}} \mathbf{G}^H + \bar{\mathbf{R}}_{\mathbf{b} \mathbf{b}} \quad (13)$$

$$\begin{aligned} \bar{\mathbf{R}}_{\mathbf{b} \mathbf{b}} &= E\{\mathbf{b} \mathbf{b}^H\} \\ &= \text{blockdiag}\{\bar{\mathbf{R}}_{\mathbf{b}_1 \mathbf{b}_1}, \mathbf{0}, \bar{\mathbf{R}}_{\mathbf{b}_2 \mathbf{b}_2}, \mathbf{0}, \dots, \bar{\mathbf{R}}_{\mathbf{b}_K \mathbf{b}_K}, \mathbf{0}\} \end{aligned} \quad (14)$$

where ϕ_s is the power spectral density (PSD) of the desired speech source signal, $\Phi_{\mathbf{u}} = E\{\mathbf{u} \mathbf{u}^H\}$ a $PL \times PL$ matrix representing the PSD of the loudspeaker signals ($L = \sum_{k=1}^K l_k$) with the PL -dimensional vector \mathbf{u} the stacked version of \mathbf{u}_k and

$$\tilde{\mathbf{G}}_{kk} = [\mathbf{G}_{kk}^H \mathbf{I}_{Pl_k \times Pl_k}]^H, \quad (15)$$

$$\tilde{\mathbf{G}}_{kq} = [\mathbf{G}_{kq}^H \mathbf{0}_{Pl_q \times Pl_k}]^H, \quad (16)$$

$$\mathbf{G} = \begin{bmatrix} \tilde{\mathbf{G}}_{11} & \dots & \tilde{\mathbf{G}}_{1K} \\ \vdots & \ddots & \vdots \\ \tilde{\mathbf{G}}_{K1} & \dots & \tilde{\mathbf{G}}_{KK} \end{bmatrix}. \quad (17)$$

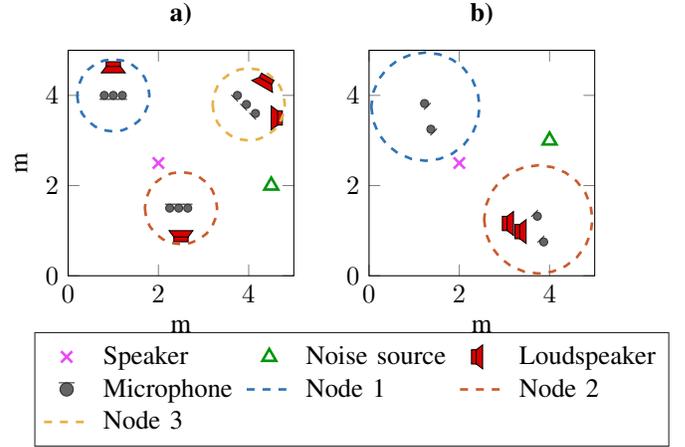


Fig. 1: Two example scenarios for a WASAN with a single target speaker and a single noise source: **a)** Three nodes each with 3 microphones and 1 or 2 loudspeakers. **b)** Two nodes each with 2 microphones. One node with a stereo loudspeaker signal.

Given that loudspeaker signals are generally non-stationary, e.g., speech and/or music signals, $\Phi_{\mathbf{u}}(l) \neq \Phi_{\mathbf{u}}(l')$ for $l \neq l'$. It is first assumed that $\Phi_{\mathbf{u}}(l) = \Phi_{\mathbf{u}}(l')$, $\forall l$, so that the noise \mathbf{n} is stationary, as required in the MWF algorithm in Section III. However this assumption will be revisited in Section III-A.

III. CENTRALIZED INTEGRATED AEC AND NR (MWF)

The node-specific combined AEC and NR task for node k is to estimate the desired signal d_k , defined here as the speech component in the first local microphone, i.e., $d_k = [1 \ 0] \mathbf{s}_k = \mathbf{e}_{d_k}^H \mathbf{s}$, where $\mathbf{0}$ is an all-zero vector with matching dimensions and $\mathbf{e}_{d_k}^H$ is a vector that selects the desired speech component in \mathbf{s} . The minimization of the mean squared error (MSE) between the desired signal and the filtered microphone and loudspeaker signals defines an optimal filter for node k ,

$$\bar{\mathbf{w}}_k = \arg \min_{\mathbf{w}_k} E\{|d_k - \mathbf{w}_k^H \mathbf{y}\|^2\}. \quad (18)$$

The node-specific signal estimate is then obtained as $\hat{d}_k = \bar{\mathbf{w}}_k^H \mathbf{y}$. The solution to this is the well-known MWF [22], [23], given by

$$\bar{\mathbf{w}}_k = \bar{\mathbf{R}}_{\mathbf{y} \mathbf{y}}^{-1} \bar{\mathbf{R}}_{\mathbf{y} d_k} = \bar{\mathbf{R}}_{\mathbf{y} \mathbf{y}}^{-1} \bar{\mathbf{R}}_{\mathbf{y} \mathbf{s}} \mathbf{e}_{d_k} = \bar{\mathbf{R}}_{\mathbf{y} \mathbf{y}}^{-1} \bar{\mathbf{R}}_{\mathbf{s} \mathbf{s}} \mathbf{e}_{d_k} \quad (19)$$

where $\bar{\mathbf{R}}_{\mathbf{y} d_k} = E\{\mathbf{y} d_k^H\}$ and $\bar{\mathbf{R}}_{\mathbf{y} \mathbf{s}} = E\{\mathbf{y} \mathbf{s}^H\}$. The final expression in (19) is obtained based on the assumption that \mathbf{s} and \mathbf{n} are uncorrelated (Section II).

In practice, by using a voice activity detector (VAD), $\bar{\mathbf{R}}_{\mathbf{y} \mathbf{y}}$ and $\bar{\mathbf{R}}_{\mathbf{n} \mathbf{n}}$ are first estimated during *speech-plus-noise* periods where the desired speech signal, loudspeaker signals and background noise are active, and *noise-only* periods where there is no activity of the desired speech signal and the other signals are active, respectively [24], i.e.,

$$\begin{aligned} \text{if VAD}(l) = 1: & \hat{\mathbf{R}}_{\mathbf{y} \mathbf{y}}(l) = \beta \hat{\mathbf{R}}_{\mathbf{y} \mathbf{y}}(l-1) + (1-\beta) \mathbf{y}(l) \mathbf{y}^H(l) \\ \text{if VAD}(l) = 0: & \hat{\mathbf{R}}_{\mathbf{n} \mathbf{n}}(l) = \beta \hat{\mathbf{R}}_{\mathbf{n} \mathbf{n}}(l-1) + (1-\beta) \mathbf{y}(l) \mathbf{y}^H(l) \end{aligned} \quad (20)$$

where $\hat{\mathbf{R}}_{yy}(l)$, $\hat{\mathbf{R}}_{nn}(l)$, $\mathbf{y}(l)$ represent $\hat{\mathbf{R}}_{yy}$, $\hat{\mathbf{R}}_{nn}$ and \mathbf{y} at frame l , respectively. The forgetting factor $0 < \beta < 1$ can be chosen depending on the variation of the statistics of the signals, i.e., if the statistics change slowly then β should be chosen close to 1 to obtain long-term estimates that mainly capture the spatial coherence between the microphone signals. For the time being, it is assumed that the loudspeaker signals and background noise are stationary (Section II), so that their contribution in $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ is the same. The following criterion will then be used to estimate $\hat{\mathbf{R}}_{ss}$ [21], [22],

$$\hat{\mathbf{R}}_{ss} = \arg \min_{\substack{\text{rank}(\mathbf{R}_{ss})=1 \\ \mathbf{R}_{ss} \succeq 0}} \left\| \hat{\mathbf{R}}_{nn}^{-1/2} \left(\hat{\mathbf{R}}_{yy} - \hat{\mathbf{R}}_{nn} - \mathbf{R}_{ss} \right) \hat{\mathbf{R}}_{nn}^{-H/2} \right\|_F^2 \quad (21)$$

where $\|\cdot\|_F$ denotes the Frobenius norm. Spatial pre-whitening is applied by pre- and post-multiplying by $\hat{\mathbf{R}}_{nn}^{-1/2}$ and $\hat{\mathbf{R}}_{nn}^{-H/2}$, respectively. The solution to (21) is based on a generalized eigenvalue decomposition (GEVD) of the $(N \times N)$ matrix pencil $\{\hat{\mathbf{R}}_{yy}, \hat{\mathbf{R}}_{nn}\}$ [22], [25]

$$\begin{aligned} \hat{\mathbf{R}}_{yy} &= \hat{\mathbf{Q}} \hat{\Sigma}_{yy} \hat{\mathbf{Q}}^H \\ \hat{\mathbf{R}}_{nn} &= \hat{\mathbf{Q}} \hat{\Sigma}_{nn} \hat{\mathbf{Q}}^H \end{aligned} \quad (22)$$

where $\hat{\Sigma}_{yy}$ and $\hat{\Sigma}_{nn}$ are diagonal matrices and $\hat{\mathbf{Q}}$ is an invertible matrix. The speech correlation matrix estimate $\hat{\mathbf{R}}_{ss}$ is then [22]

$$\hat{\mathbf{R}}_{ss} = \hat{\mathbf{Q}} \text{diag}\{\hat{\sigma}_{y_1} - \hat{\sigma}_{n_1}, 0, \dots, 0\} \hat{\mathbf{Q}}^H \quad (23)$$

where $\hat{\sigma}_{y_1}$ and $\hat{\sigma}_{n_1}$ are the first diagonal element of $\hat{\Sigma}_{yy}$ and $\hat{\Sigma}_{nn}$, respectively, corresponding to the largest ratio $\hat{\sigma}_{y_1}/\hat{\sigma}_{n_1}$. Using (23) and $\hat{\mathbf{R}}_{yy}$ (cfr. (22)) in (19), the MWF estimate $\hat{\mathbf{w}}_k$ can be expressed as

$$\hat{\mathbf{w}}_k = \hat{\mathbf{Q}}^{-H} \text{diag} \left\{ 1 - \frac{\hat{\sigma}_{n_1}}{\hat{\sigma}_{y_1}}, 0, \dots, 0 \right\} \hat{\mathbf{Q}}^H \mathbf{e}_{d_k}. \quad (24)$$

The node-specific signal estimate is then obtained as $\hat{d}_k = \hat{\mathbf{w}}_k^H \mathbf{y}$. In this integrated algorithm, the MWF estimate depends on the loudspeaker signal statistics without exploiting the prior knowledge that there is no desired speech component in these loudspeaker signals. As a consequence, the combined AEC and NR fundamentally consists of a single NR stage in which acoustic echo is treated similarly to background noise.

A. Non-stationarity of loudspeaker signals and MWF assumptions

As mentioned in Section II, the loudspeaker signals are generally non-stationary, i.e., $\Phi_{\mathbf{u}}(l) \neq \Phi_{\mathbf{u}}(l')$ for $l \neq l'$. As a consequence their contribution in the *speech-plus-noise* and *noise-only* correlation matrices, $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$, respectively, may be different. This violates the basic stationarity assumption in the MWF algorithm described above. However, it is observed that this non-stationarity does not change significantly the GEVD of the matrix pencil $\{\hat{\mathbf{R}}_{yy}, \hat{\mathbf{R}}_{nn}\}$ because of the specific structure of $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ corresponding to the fact that the loudspeaker signals do not contain any desired speech

and background noise component. In particular, this will lead to the following structure in $\hat{\mathbf{Q}}$, $\hat{\Sigma}_{yy}$ and $\hat{\Sigma}_{nn}$ in (22)

$$\hat{\mathbf{Q}} = \left[\begin{array}{c|c|c} \hat{\mathbf{q}}_1 & \hat{\mathbf{Q}}_1 & \hat{\mathbf{q}}_2 \dots \hat{\mathbf{q}}_M \\ \hline N \times 1 & N \times PL & \end{array} \right] \quad (25)$$

$$\hat{\Sigma}_{yy} = \begin{bmatrix} \hat{\sigma}_{y_1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \hat{\Sigma}_{yy,1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \hat{\Sigma}_{yy,2} \\ & & & \ddots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \hat{\Sigma}_{yy,M-1} \end{bmatrix} \quad (26)$$

$$\hat{\Sigma}_{nn} = \begin{bmatrix} \hat{\sigma}_{n_1} & \mathbf{0} & \mathbf{0} \\ \mathbf{0} & \hat{\Sigma}_{nn,1} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} & \hat{\Sigma}_{nn,2} \\ & & & \ddots \\ \mathbf{0} & \mathbf{0} & \mathbf{0} & \hat{\Sigma}_{nn,M-1} \end{bmatrix} \quad (27)$$

where $(\hat{\mathbf{q}}_1 \dots \hat{\mathbf{q}}_M)$ are column vectors uniquely defined by the desired speech component and background noise ($M = \sum_{k=1}^K m_k$), hence containing zeros in the positions corresponding to the loudspeaker signals, and $\hat{\mathbf{Q}}_1$ contains PL columns which are uniquely defined by the loudspeaker signals and echo paths. The non-stationarity of the loudspeaker signals does not modify $(\hat{\mathbf{q}}_1 \dots \hat{\mathbf{q}}_M)$, $\hat{\sigma}_{y_1}, \hat{\sigma}_{n_1}$, $\hat{\Sigma}_{yy,2}$ and $\hat{\Sigma}_{nn,2}$. It also does not modify the column space spanned by $\hat{\mathbf{Q}}_1$. As a result, the first column of $\hat{\mathbf{Q}}^H$ in (24) is not modified, as well as all other relevant quantities in (24). Therefore, the MWF estimate in (24) is also not modified. Note that it is assumed here that the GEVLs corresponding to $\hat{\mathbf{Q}}_1$ are smaller than the GEVL corresponding to $\hat{\mathbf{q}}_1$, i.e. to the desired speech signal, so the latter continues to be the largest GEVL. For the unlikely scenario that a GEVL corresponding to $\hat{\mathbf{Q}}_1$ becomes the largest GEVL, $\hat{\mathbf{q}}_1$ may be monitored (based on its zeros structure) and tracked, so that the correct GEVL is still chosen.

IV. CENTRALIZED CASCADE AEC AND NR (PK-MWF)

Exploiting the prior knowledge that $\hat{\mathbf{R}}_{ss}$ has a specific zero structure (cfr. definition of \mathbf{s} and $\tilde{\mathbf{s}}_k$), the criterion in (21) can be redefined as

$$\hat{\mathbf{R}}_{ss} = \arg \min_{\substack{\text{rank}(\mathbf{R}_{ss})=1 \\ \mathbf{B}^H \mathbf{R}_{ss} \mathbf{B} = 0 \\ \mathbf{R}_{ss} \succeq 0}} \left\| \hat{\mathbf{R}}_{nn}^{-1/2} \left(\hat{\mathbf{R}}_{yy} - \hat{\mathbf{R}}_{nn} - \mathbf{R}_{ss} \right) \hat{\mathbf{R}}_{nn}^{-H/2} \right\|_F^2 \quad (28)$$

where \mathbf{B} is an $N \times PL$ block diagonal matrix

$$\mathbf{B} = \begin{bmatrix} \mathbf{B}_1 & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0} & \mathbf{B}_2 & \dots & \mathbf{0} \\ \vdots & \vdots & \ddots & \vdots \\ \mathbf{0} & \mathbf{0} & \dots & \mathbf{B}_K \end{bmatrix} \quad (29)$$

with the k^{th} diagonal block \mathbf{B}_k equal to

$$\mathbf{B}_k = \begin{bmatrix} \mathbf{0}_{m_k \times Pl_k} \\ \mathbf{I}_{Pl_k} \end{bmatrix}, \quad (30)$$

where \mathbf{I}_{Pl_k} is a $Pl_k \times Pl_k$ identity matrix. In the combined AEC and NR context \mathbf{B} is a selection matrix that selects the loudspeaker signals. In [21] it is shown that the inclusion of the constraint $\mathbf{B}^H \mathbf{R}_{ss} \mathbf{B} = 0$ leads to the reduced dimensional $(M \times M)$ matrix pencil $\{\mathbf{R}_{yy}^{\text{red}}, \mathbf{R}_{nn}^{\text{red}}\}$ with GEVD

$$\begin{aligned} \hat{\mathbf{R}}_{yy}^{\text{red}} &= \hat{\mathbf{Q}}^{\text{red}} \hat{\Sigma}_{yy}^{\text{red}} (\hat{\mathbf{Q}}^{\text{red}})^H \\ \hat{\mathbf{R}}_{nn}^{\text{red}} &= \hat{\mathbf{Q}}^{\text{red}} \hat{\Sigma}_{nn}^{\text{red}} (\hat{\mathbf{Q}}^{\text{red}})^H \end{aligned} \quad (31)$$

where $\hat{\mathbf{R}}_{yy}^{\text{red}} = \hat{\mathbf{C}}^H \hat{\mathbf{R}}_{yy} \hat{\mathbf{C}}$, $\hat{\mathbf{R}}_{nn}^{\text{red}} = \hat{\mathbf{C}}^H \hat{\mathbf{R}}_{nn} \hat{\mathbf{C}}$, $\mathbf{y}^{\text{red}} = \hat{\mathbf{C}}^H \mathbf{y}$, and with $\hat{\mathbf{C}}$ an $N \times M$ matrix obtained from the linearly-constrained minimum variance (LCMV) beamformer optimization criterion

$$\hat{\mathbf{C}} = \underset{\text{s.t. } \mathbf{H}^H \mathbf{C} = \mathbf{I}_M}{\arg \min} \left\| \text{trace}\{\mathbf{C}^H \hat{\mathbf{R}}_{nn} \mathbf{C}\} \right\| \quad (32)$$

where \mathbf{H} is a $N \times M$ block diagonal matrix

$$\mathbf{H} = \begin{bmatrix} \mathbf{H}_1 & 0 & \dots & 0 \\ 0 & \mathbf{H}_2 & \dots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \dots & \mathbf{H}_K \end{bmatrix} \quad (33)$$

with the k^{th} diagonal block equal to

$$\mathbf{H}_k = \begin{bmatrix} \mathbf{I}_{m_k} \\ \mathbf{0}_{Pl_k \times m_k} \end{bmatrix}, \quad (34)$$

such that $\mathbf{H}^H \mathbf{H} = \mathbf{I}_M$ and $\mathbf{B}^H \mathbf{H} = \mathbf{0}$. Hence $\hat{\mathbf{C}}$ can be defined based on a generalized sidelobe canceller (GSC) implementation as [21], [26]

$$\hat{\mathbf{C}} = \mathbf{H} - \mathbf{B} \hat{\mathbf{F}} \quad (35)$$

$$\hat{\mathbf{F}} = (\mathbf{B}^H \hat{\mathbf{R}}_{nn} \mathbf{B})^{-1} \mathbf{B}^H \hat{\mathbf{R}}_{nn} \mathbf{H} \quad (36)$$

where the filter $\hat{\mathbf{F}}$ operates on the loudspeaker signals ($\mathbf{B}^H \mathbf{y}$) and effectively serves as an AEC filter cancelling the echo components in the so-called fixed beamformer outputs corresponding to \mathbf{H} , i.e., the microphone signals ($\mathbf{H}^H \mathbf{y}$). The inclusion of the prior knowledge thus leads to a cascade algorithm where AEC is performed first and then NR. The AEC filter $\hat{\mathbf{F}}$ can also be implemented adaptively via an NLMS or QRD-RLS algorithm as will be explained in Section VII.

The prior knowledge speech correlation matrix estimate $\hat{\mathbf{R}}_{ss}$, i.e., the solution to (28), is then given as [21], [22],

$$\hat{\mathbf{R}}_{ss} = \mathbf{H} \hat{\mathbf{Q}}^{\text{red}} \text{diag}\{\hat{\sigma}_{y_1} - \hat{\sigma}_{n_1}, 0, \dots, 0\} (\hat{\mathbf{Q}}^{\text{red}})^H \mathbf{H}^H, \quad (37)$$

where $\hat{\sigma}_{y_1}$ and $\hat{\sigma}_{n_1}$ are the first diagonal element of $\hat{\Sigma}_{yy}^{\text{red}}$ and $\hat{\Sigma}_{nn}^{\text{red}}$, respectively, corresponding to the largest ratio $\hat{\sigma}_{y_i} / \hat{\sigma}_{n_i}$. Using this expression and the reduced dimensional $\hat{\mathbf{R}}_{yy}^{\text{red}}$ (cfr. (31)), the PK-MWF estimate $\hat{\mathbf{w}}_k$ can finally be expressed as [21]

$$\hat{\mathbf{w}}_k = \hat{\mathbf{C}} (\hat{\mathbf{Q}}^{\text{red}})^{-H} \text{diag}\left\{1 - \frac{\hat{\sigma}_{n_1}}{\hat{\sigma}_{y_1}}, 0, \dots, 0\right\} (\hat{\mathbf{Q}}^{\text{red}})^H \mathbf{H}^H \mathbf{e}_{d_k}. \quad (38)$$

The non-stationarity of the loudspeaker signals in this case does not affect the NR stage, as the joint-diagonalization is performed on the reduced dimensional ($M \times M$) matrix pencil $\{\hat{\mathbf{R}}_{yy}^{\text{red}}, \hat{\mathbf{R}}_{nn}^{\text{red}}\}$, therefore $\hat{\mathbf{Q}}^{\text{red}}$ will only have M columns defined by the desired speech components and background noise, and the echo signals are effectively removed by the AEC stage.

V. DISTRIBUTED INTEGRATED AEC AND NR (GEVD-DANSE)

The integrated AEC and NR algorithm of Section III can be implemented in a distributed fashion by means of the GEVD-DANSE algorithm [18] where each node instead of broadcasting n_k microphone and loudspeaker signals, broadcasts only 1 fused signal to the other nodes. Each node performs local operations, corresponding to a reduced dimensional version (dimension $n_k + (K - 1)$ in node k) of the MWF-based integrated AEC and NR algorithm of Section III (dimension N), based on n_k local microphone and loudspeaker signals and $(K - 1)$ fused signals received from the other nodes. The fused signal broadcast by node k is

$$z_k = \hat{\mathbf{p}}_k^H \mathbf{y}_k \quad (39)$$

where $\hat{\mathbf{p}}_k$ is an n_k -dimensional fusion vector. Then each node has access to a signal vector $\tilde{\mathbf{y}}_k = [\mathbf{y}_k^H \mathbf{z}_{-k}^H]^H$, where the subscript $-k$ refers to the concatenation of the fused signals of nodes other than k , so that $\mathbf{z}_{-k} = [z_1^* \dots z_{k-1}^* z_{k+1}^* \dots z_K^*]^H$, where $*$ represents the complex conjugate. The local filter $\hat{\mathbf{w}}_k$ is defined as

$$\hat{\mathbf{w}}_k = \hat{\mathbf{Q}}_k^{-H} \text{diag}\left\{1 - \frac{\hat{\sigma}_{n_1}}{\hat{\sigma}_{y_1}}, 0, \dots, 0\right\} \hat{\mathbf{Q}}_k^H [1 \ \mathbf{0}]^H \quad (40)$$

with the GEVD of the $(n_k + K - 1) \times (n_k + K - 1)$ matrix pencil $\{\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}, \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}\}$ given as

$$\begin{aligned} \hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} &= \hat{\mathbf{Q}}_k \hat{\Sigma}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} \hat{\mathbf{Q}}_k^H \\ \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} &= \hat{\mathbf{Q}}_k \hat{\Sigma}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} \hat{\mathbf{Q}}_k^H \end{aligned} \quad (41)$$

where $\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}$ is an estimate of $\bar{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} = E\{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^H\}$, $\hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}$ is an estimate of $\bar{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} = E\{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k^H\}$ and $\tilde{\mathbf{n}}_k$ corresponds to $\tilde{\mathbf{y}}_k$ in *noise-only* periods. The fusion vector is finally defined as

$$\hat{\mathbf{p}}_k = [\mathbf{I}_{n_k} \ \mathbf{0}] \hat{\mathbf{w}}_k. \quad (42)$$

In each time frame the nodes broadcast fused signals (39) using their current fusion vectors. One node then updates its fusion vector by means of (40)-(42). When the nodes update sequentially in a round-robin fashion (e.g. one node updates per time frame) the local signal estimates $\hat{d}_k = \hat{\mathbf{w}}_k^H \tilde{\mathbf{y}}_k$ have been shown to converge in each node to the centralized signal estimates obtained with (24) [18]. It has also been shown that when the nodes update simultaneously a relaxation factor ($\alpha_{r,S}$) is needed to avoid limit cycles. With this each filter is updated as a convex combination of its previous and newly computed version in (40) [18], [27].

VI. DISTRIBUTED CASCADE AEC AND NR (PK-GEVD-DANSE)

The cascade AEC and NR algorithm of Section IV can be implemented in a distributed fashion by means of the PK-GEVD-DANSE algorithm [21] where each node broadcasts 2 fused signals, i.e., a desired signal reference and a noise reference. In the context of combined AEC and NR, the second fused signal will be a fused loudspeaker signal. Each node then performs local operations, effectively corresponding to a reduced dimensional version (dimension $n_k + 2(K - 1)$ in node

k) of the PK-MWF-based cascade AEC and NR algorithm of Section IV (dimension N), based on n_k local microphone and loudspeaker signals and $2(K-1)$ fused signals received from the other nodes. The first fused signal broadcast by node k is given by (39) with the n_k -dimensional fusion vector $\hat{\mathbf{p}}_k$ to be redefined. The second fused signal broadcast by node k is

$$\tilde{z}_k = \hat{\lambda}_k^H \mathbf{B}_k^H \mathbf{y}_k = \hat{\lambda}_k^H \mathbf{u}_k \quad (43)$$

where $\hat{\lambda}_k$ is a Pl_k -dimensional fusion vector. Then each node has access to a signal vector $\tilde{\mathbf{y}}_k = [\mathbf{y}_k^H \mathbf{z}_{-k}^H \tilde{\mathbf{z}}_k^H]^H$, where \mathbf{z}_{-k} is defined as in Section V and $\tilde{\mathbf{z}}_k = [\tilde{z}_1^* \cdots \tilde{z}_{k-1}^* \tilde{z}_{k+1}^* \cdots \tilde{z}_K^*]^H$. A modification must be introduced in \mathbf{H}_k and \mathbf{B}_k to account for the extra signals broadcast from the other nodes, hence

$$\check{\mathbf{H}}_k = \left[\begin{array}{c|c} \mathbf{I}_{m_k} & \mathbf{0} \\ \hline \mathbf{0}_{Pl_k \times m_k} & \mathbf{I}_{K-1} \\ \hline \mathbf{0} & \mathbf{0}_{(K-1) \times (K-1)} \end{array} \right], \quad \check{\mathbf{B}}_k = \left[\begin{array}{c|c} \mathbf{0}_{m_k \times Pl_k} & \mathbf{0} \\ \hline \mathbf{I}_{Pl_k} & \mathbf{0}_{(K-1) \times (K-1)} \\ \hline \mathbf{0} & \mathbf{I}_{K-1} \end{array} \right] \quad (44)$$

where $\check{\mathbf{H}}_k$ is an $(n_k + 2(K-1)) \times (m_k + K-1)$ matrix and $\check{\mathbf{B}}_k$ is an $(n_k + 2(K-1)) \times (Pl_k + K-1)$ matrix. Then equations (35) and (36) become respectively

$$\hat{\mathbf{C}}_k = \check{\mathbf{H}}_k - \check{\mathbf{B}}_k \hat{\mathbf{F}}_k \quad (45)$$

$$\hat{\mathbf{F}}_k = (\check{\mathbf{B}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} \check{\mathbf{B}}_k)^{-1} \check{\mathbf{B}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} \check{\mathbf{H}}_k \quad (46)$$

where $\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}$ is an estimate of $\tilde{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} = E\{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k^H\}$, $\hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}$ is an estimate of $\tilde{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} = E\{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k^H\}$ and $\tilde{\mathbf{n}}_k$ corresponds to $\tilde{\mathbf{y}}_k$ in *noise-only* periods. The fusion vectors are defined as in (42) and as [21]

$$\hat{\lambda}_k = [\mathbf{I}_{Pl_k} \mathbf{0}] (\check{\mathbf{B}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} \check{\mathbf{B}}_k)^{-1} \check{\mathbf{B}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} \hat{\mathbf{w}}_k \quad (47)$$

where the local filter $\hat{\mathbf{w}}_k$ is defined as

$$\hat{\mathbf{w}}_k = \hat{\mathbf{C}}_k (\hat{\mathbf{Q}}_k^{\text{red}})^{-H} \text{diag} \left\{ 1 - \frac{\hat{\sigma}_{n_1}}{\hat{\sigma}_{y_1}}, 0, \dots, 0 \right\} (\hat{\mathbf{Q}}_k^{\text{red}})^H \check{\mathbf{H}}_k^H \mathbf{e}_{d_k} \quad (48)$$

with the GEVD of the reduced dimensional $(m_k + K-1) \times (m_k + K-1)$ pencil $\{\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}^{\text{red}}, \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}^{\text{red}}\}$ given as

$$\begin{aligned} \hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}^{\text{red}} &= \hat{\mathbf{Q}}_k^{\text{red}} \hat{\Sigma}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}^{\text{red}} (\hat{\mathbf{Q}}_k^{\text{red}})^H \\ \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}^{\text{red}} &= \hat{\mathbf{Q}}_k^{\text{red}} \hat{\Sigma}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}^{\text{red}} (\hat{\mathbf{Q}}_k^{\text{red}})^H \end{aligned} \quad (49)$$

where $\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}^{\text{red}} = \hat{\mathbf{C}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} \hat{\mathbf{C}}_k$, $\hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}^{\text{red}} = \hat{\mathbf{C}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} \hat{\mathbf{C}}_k$ and $\tilde{\mathbf{y}}_k^{\text{red}} = \hat{\mathbf{C}}_k^H \tilde{\mathbf{y}}_k$. In each time frame the nodes broadcast fused signals (39) and (43) using their current fusion vectors. One node then updates its fusion vectors by means of (44)-(49). When the nodes update sequentially in a round-robin fashion (e.g. one node updates per time frame) the local signal estimates $\hat{d}_k = \hat{\mathbf{w}}_k^H \tilde{\mathbf{y}}_k$ have been shown to converge in each node to the centralized signal estimates obtained with (38) [21].

In an AEC context (unlike the general PK-GEVD-DANSE context), the above algorithm and its communication cost can

be reduced as follows. First, the loudspeaker signals ($\check{\mathbf{B}}_k^H \tilde{\mathbf{y}}_k$) do not contain any desired speech component, hence

$$\check{\mathbf{B}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k} \hat{\mathbf{w}}_k \simeq \check{\mathbf{B}}_k^H \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k} \hat{\mathbf{w}}_k = \mathbf{0}, \quad (50)$$

which in (47) leads to $\hat{\lambda}_k \simeq \mathbf{0}$. In (50) the " \simeq " is replaced by an equality if estimated correlation matrices are replaced by true statistical quantities. The " $=$ " in (50) is obtained by substituting (45), (46) and (48). Therefore the extra fused signal \tilde{z}_k does not need to be communicated to the other nodes. Therefore, $\check{\mathbf{H}}_k$ and $\check{\mathbf{B}}_k$ are reduced to

$$\check{\mathbf{H}}_k = \left[\begin{array}{c|c} \mathbf{I}_{m_k} & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{I}_{K-1} \end{array} \right], \quad \check{\mathbf{B}}_k = \left[\begin{array}{c} \mathbf{0} \\ \hline \mathbf{I}_{Pl_k} \\ \hline \mathbf{0} \end{array} \right] \quad (51)$$

with $\check{\mathbf{H}}_k$ an $(n_k + K-1) \times (m_k + K-1)$ matrix and $\check{\mathbf{B}}_k$ an $(n_k + K-1) \times Pl_k$ matrix. The vector $\tilde{\mathbf{y}}_k$ is then reduced to

$$\tilde{\mathbf{y}}_k = \begin{bmatrix} \mathbf{y}_k \\ \mathbf{z}_{-k} \end{bmatrix}. \quad (52)$$

This effectively leads to a distributed cascade algorithm where each node k shares only 1 fused signal with the other nodes. Each node performs local operations corresponding to a reduced dimensional version (dimension $n_k + K-1$ in node k) of the PK-MWF-based cascade AEC and NR algorithm of Section IV.

VII. PRACTICAL CONSIDERATIONS

Up until now it has been assumed that the loudspeaker signals are stationary and so always active, which allows to estimate the speech correlation matrix $\hat{\mathbf{R}}_{\text{ss}}$ based on (21) or (28). In practice this assumption is not a valid one when the loudspeakers play back speech or music signals. Therefore the VAD should be able to detect the activity of the desired speech signal in the presence of loudspeaker signals which may also contain speech signals, and other background noise signals. A cascade approach then allows to consider two aspects in the AEC stage. Firstly, different adaptive filtering algorithms can be used such as NLMS, RLS, etc., to estimate the AEC filters. Secondly, the activity of the loudspeaker signals can be monitored and then used to control the adaptive filters. This section describes adaptive implementations of the PK-GEVD-DANSE algorithm described in Section VI. An NLMS or QRD-RLS algorithm is used to update \mathbf{F}_k in the AEC stage. Then the procedure to deal with the loudspeaker activity detection is also presented.

A. NLMS-based PK-GEVD-DANSE

A first adaptive implementation of the PK-GEVD-DANSE algorithm (Section VI) can be obtained using the NLMS algorithm as [28]

$$\hat{\mathbf{F}}_k(l) = \hat{\mathbf{F}}_k(l-1) + \mu \mathbf{u}_k(l) \left(\left(\check{\mathbf{H}}_k - \check{\mathbf{B}}_k \hat{\mathbf{F}}_k(l-1) \right)^H \tilde{\mathbf{y}}_k(l) \right)^H \quad (53)$$

$$\mu = \frac{\mu_F}{\hat{\Phi}_{\mathbf{u}_k}(l) + \delta} \quad (54)$$

$$\hat{\Phi}_{\mathbf{u}_k}(l) = \mathbf{u}_k^H(l) \mathbf{u}_k(l) \quad (55)$$

where $\hat{\Phi}_{\mathbf{u}_k}(l)$ is an estimate of the PSD of the local loudspeaker signals, δ is a regularization term and μ_F is the step size. The filter $\hat{\mathbf{F}}_k(l)$ is updated whenever the desired speech signal is not active and the loudspeaker signals are active (i.e., no double-talk).

B. QRD-RLS-based PK-GEVD-DANSE

Alternatively $\hat{\mathbf{F}}_k$ can be computed using a QRD-RLS algorithm. It is assumed that

$$\mathbf{U}_k = \begin{bmatrix} (\check{\mathbf{B}}_k^H \check{\mathbf{y}}_k(l))^H \\ (\check{\mathbf{B}}_k^H \check{\mathbf{y}}_k(l-1))^H \\ \vdots \\ (\check{\mathbf{B}}_k^H \check{\mathbf{y}}_k(l-L_f+1))^H \end{bmatrix} \quad (56)$$

admits a QR-decomposition (QRD)

$$\mathbf{U}_k = \mathbf{Q}_k \begin{bmatrix} \mathcal{R}_k \\ \mathbf{0} \end{bmatrix} \quad (57)$$

where \mathbf{Q}_k is an orthogonal matrix, such that $\mathbf{Q}_k^H \mathbf{Q}_k = \mathbf{I}$, \mathcal{R}_k is an upper triangular matrix and $L_f > Pl_k$. The microphone signals are expressed as

$$\mathbf{Y}_k = \begin{bmatrix} (\check{\mathbf{H}}_k^H \check{\mathbf{y}}_k(l))^H \\ (\check{\mathbf{H}}_k^H \check{\mathbf{y}}_k(l-1))^H \\ \vdots \\ (\check{\mathbf{H}}_k^H \check{\mathbf{y}}_k(l-L_f+1))^H \end{bmatrix} \quad (58)$$

and based on this, the following optimization criterion can be defined

$$\begin{aligned} \hat{\mathbf{F}}_k &= \arg \min_{\mathbf{F}_k} \|\mathbf{Y}_k - \mathbf{U}_k \mathbf{F}_k\|_2^2 \\ &= \arg \min_{\mathbf{F}_k} \|\mathbf{Q}_k^H (\mathbf{Y}_k - \mathbf{U}_k \mathbf{F}_k)\|_2^2 \\ &= \arg \min_{\mathbf{F}_k} \left\| \begin{bmatrix} \mathcal{Z}_k \\ * \end{bmatrix} - \begin{bmatrix} \mathcal{R}_k \\ \mathbf{0} \end{bmatrix} \mathbf{F}_k \right\|_2^2 \end{aligned} \quad (59)$$

where the rows of \mathcal{Z}_k correspond to the first Pl_k rows of $\mathbf{Q}_k^H \mathbf{Y}_k$. The star $*$ represents entries to the matrix that are not of interest, however it has been shown that these can be used to compute least squares residuals [28]. The solution to (59) is given by

$$\hat{\mathbf{F}}_k = \mathcal{R}_k^{-1} \mathcal{Z}_k \quad (60)$$

which can be computed by backsubstitution. A QRD-updating can be used in frame l if \mathcal{R}_k is available from frame $l-1$, by means of the following recursions

$$\begin{bmatrix} \mathcal{R}_k(l) & \mathcal{Z}_k(l) \\ \mathbf{0} & * \end{bmatrix} = \mathcal{G}^H \begin{bmatrix} \gamma \mathcal{R}_k(l-1) & \gamma \mathcal{Z}_k(l-1) \\ (\check{\mathbf{B}}_k^H \check{\mathbf{y}}_k(l))^H & (\check{\mathbf{H}}_k^H \check{\mathbf{y}}_k(l))^H \end{bmatrix} \quad (61)$$

$$\hat{\mathbf{F}}_k(l) = \mathcal{R}_k(l)^{-1} \mathcal{Z}_k(l). \quad (62)$$

where γ is a forgetting factor and \mathcal{G} is an orthogonal transformation. The matrix on the right hand side of (61) is retriangularized via a series of complex Givens rotations [29]. The new quantities are then used to obtain $\hat{\mathbf{F}}_k(l)$ (cfr. (60)).

A complete description of the NLMS-based and QRD-RLS-based PK-GEVD-DANSE algorithm using simultaneous

updating is shown as Algorithm 1 and 2, respectively. The fusion vectors in line 12 and 13 in Algorithm 1 and 2 respectively, can be updated once every D frames, using the previous and current values, and a relaxation factor as in [27].

Algorithm 1: NLMS based PK-GEVD-DANSE

- 1 Construct $\check{\mathbf{H}}_k$ and $\check{\mathbf{B}}_k$ based on m_k , l_k and K (cfr. (51));
 - 2 Randomly initialize fusion vector $\hat{\mathbf{p}}_k \forall k \in \mathcal{K}$;
 - 3 Each node $k \in \mathcal{K}$ performs the following simultaneously;
 - 4 **for** $l = 1, 2, 3, \dots$ **do**
 - 5 Collect observations $\mathbf{y}_k(l)$
 - 6 Compute $z_k(l)$ (cfr. (39)) and broadcast to other nodes;
 - 7 Construct $\check{\mathbf{y}}_k(l)$ based on $\mathbf{z}_{-k}(l)$;
 - 8 Update $\hat{\mathbf{R}}_{\check{\mathbf{y}}_k \check{\mathbf{y}}_k}$ or $\hat{\mathbf{R}}_{\check{\mathbf{n}}_k \check{\mathbf{n}}_k}$ similar to (20) based on VAD;
 - 9 Compute $\Phi_{u_k}(l)$ and update $\hat{\mathbf{F}}_k(l)$, following (55) and (53), respectively;
 - 10 Update its local LCMV beamformer $\hat{\mathbf{C}}_k$ using (45) and estimate $\hat{\mathbf{R}}_{\check{\mathbf{y}}_k \check{\mathbf{y}}_k}^{\text{red}}$ and $\hat{\mathbf{R}}_{\check{\mathbf{n}}_k \check{\mathbf{n}}_k}^{\text{red}}$;
 - 11 Update $\hat{\mathbf{w}}_k$ using (48) by means of the GEVD of $\{\hat{\mathbf{R}}_{\check{\mathbf{y}}_k \check{\mathbf{y}}_k}^{\text{red}}, \hat{\mathbf{R}}_{\check{\mathbf{n}}_k \check{\mathbf{n}}_k}^{\text{red}}\}$;
 - 12 Compute fusion vector $\hat{\mathbf{p}}_k$ following (42);
 - 13 Compute estimated node-specific desired signal $\hat{d}_k(l) = \hat{\mathbf{w}}_k^H \check{\mathbf{y}}_k(l)$;
-

C. Loudspeaker activity detection

The loudspeaker signals activity at each node directly affects the adaptation and convergence of the adaptive filters in the AEC stage. For this, a VAD to monitor the activity of the local loudspeaker signals is required to control the updating of the *noise-only* correlation matrix and therefore the adaptation of $\hat{\mathbf{F}}_k$. Using two binary VADs, i.e., VAD_s for the desired speech signal and VAD_p for the local loudspeaker signals, brings four possible combined outcomes for which the operations are defined in the following sections.

1) $\text{VAD}_s = 1$ and $\text{VAD}_p = 0$: The desired speech signal is active and the local loudspeaker signals are inactive. Update *speech-plus-noise* correlation matrix. Do not update $\hat{\mathbf{F}}_k$. This stage allows to collect information about the desired speech signal in the correlation matrix.

2) $\text{VAD}_s = 0$ and $\text{VAD}_p = 1$: The desired speech signal is inactive and the local loudspeaker signals are active. Update *noise-only* correlation matrix and update $\hat{\mathbf{F}}_k$. This is similar to what a standard AEC filter does, where the information about the loudspeaker signals is used to update the AEC filters when the desired speech signal is not active.

3) $\text{VAD}_s = 0$ and $\text{VAD}_p = 0$: Both the desired speech signal and local loudspeaker signals are inactive. Update *noise-only* correlation matrix and do not update $\hat{\mathbf{F}}_k$. This updates the background noise component in the *noise-only* correlation matrix, and given the absence of loudspeaker signals activity, no adaptation is performed of the AEC filters.

Algorithm 2: QRD-RLS based PK-GEVD-DANSE

- 1 Construct $\tilde{\mathbf{H}}_k$ and $\tilde{\mathbf{B}}_k$ based on m_k , l_k and K (cfr. (51));
- 2 Randomly initialize fusion vector $\hat{\mathbf{p}}_k \forall k \in \mathcal{K}$;
- 3 Initialize the matrices $\mathcal{R}_k = 10^{-6} \mathbf{I}_{P l_k}$ and $\mathcal{Z}_k = \mathbf{0}_{P l_k \times (m_k + K - 1)}$;
- 4 Each node $k \in \mathcal{K}$ performs the following simultaneously;
- 5 **for** $l = 1, 2, 3, \dots$ **do**
- 6 Collect observations $\mathbf{y}_k(l)$;
- 7 Compute $z_k(l)$ (cfr. (39)) and broadcast to other nodes;
- 8 Construct $\tilde{\mathbf{y}}_k(l)$ based on $\mathbf{z}_{-k}(l)$;
- 9 Update $\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}$ or $\hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}$ similar to (20) based on VAD;
- 10 Update $\mathcal{R}_k(l)$ and $\mathcal{Z}_k(l)$ based on (61) and compute $\hat{\mathbf{F}}_k(l)$ according to (62);
- 11 Update its local LCMV beamformer $\hat{\mathbf{C}}_k$ using (45) and estimate $\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}^{\text{red}}$ and $\hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}^{\text{red}}$;
- 12 Update $\hat{\mathbf{w}}_k$ using (48) by means of the GEVD of $\{\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}^{\text{red}}, \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}^{\text{red}}\}$;
- 13 Compute fusion vector $\hat{\mathbf{p}}_k$ following (42);
- 14 Compute estimated node-specific desired signal $\hat{d}_k(l) = \hat{\mathbf{w}}_k^H \tilde{\mathbf{y}}_k(l)$;

4) $VAD_s = 1$ and $VAD_p = 1$: Both the desired speech signal and local loudspeaker signals are active (double-talk). Update *speech-plus-noise* correlation matrix and do not update $\hat{\mathbf{F}}_k$. The effects of double-talk in the filter adaptation have been well documented in the literature, where it is shown that the filter adaptation can suffer from this, and even diverge.

VIII. SIMULATIONS

A. Simulation scenario

In order to assess the performance of the proposed algorithms under different situations, the scenarios shown in Fig. 1 are considered. The simulations aim at representing situations that may occur in real life, such as more than two nodes in the WASAN, stereo echo cancellation, long echo paths and reverberation times, double-talk and highly correlated signals being reproduced in two or more nodes. The simulations are grouped in batch processing, per-frame processing and adaptive processing simulations.

B. Batch-processing

This section outlines the batch simulations carried out using the MWF and PK-MWF algorithm presented in Sections III and IV respectively and the iterative GEVD-DANSE and PK-GEVD-DANSE algorithm described in Sections V and VI, respectively. First, a comparison of these algorithms using the normalized mean squared error (NMSE) metric is shown in Fig. 2 for a WASAN consisting of 4 nodes, with $\{6, 3, 8, 3\}$ microphones and $\{1, 3, 2, 3\}$ loudspeakers respectively. For comparison, the NMSE that each node would achieve working

in isolation using the PK-MWF algorithm is also shown. The desired, loudspeaker and noise signals were uncorrelated white noise signals. The loudspeaker and noise signals were continuously active while the desired signal had an ON/OFF behaviour. The NMSE was computed in the discrete time domain as

$$\text{NMSE} = 10 \log_{10} \left(\frac{1}{T} \sum_{t=1}^T \frac{(d_k(t) - \hat{d}_k(t))^2}{(d_k(t))^2} \right) \text{ dB} \quad (63)$$

where T is the sample duration of the signal and t the discrete time index. In frames where the desired signal $d_k(t)$ was zero, the power from previous frames was used for the normalization. It is observed in Fig. 2 that including the PK in both the centralized and distributed implementations reduces the NMSE in the estimation of the desired signal. In all nodes the distributed implementation outperforms each node working in isolation. This scenario deals with a large number of microphone and loudspeaker signals, and the activity patterns of the latter may not be very realistic. The scenario depicted

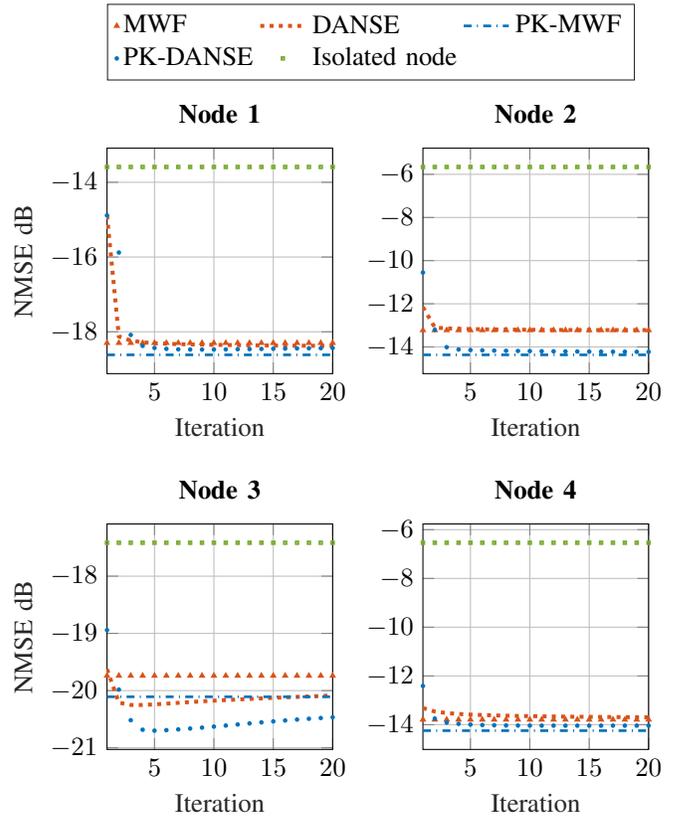


Fig. 2: NMSE for the MWF, GEVD-DANSE, PK-MWF and PK-GEVD-DANSE algorithms at each node for a WASAN with 4 nodes, with $\{6, 3, 8, 3\}$ microphones signals and $\{1, 3, 2, 3\}$ loudspeakers, respectively.

in Fig. 1a is now analyzed with no interframe filtering, i.e., $P = 1$. The performance of the algorithm was measured in terms of the echo return loss enhancement (ERLE), the signal-to-noise ratio (SNR) and the NMSE. The simulations were set up as follows. Firstly, the microphone and loudspeaker signals were simulated at each node using room impulse responses

of 500 samples long simulated with the randomized image method described in [30] and a sampling frequency of 16 kHz. The reflection coefficient of all surfaces in the room was set to 0.15 (for a reverberation time $T_{60} = 0.11$ s), and the random displacement of the image sources to 0.13 m. The inter-microphone distance of the arrays was set to 20 cm for all the nodes. The microphone signals were created such that the signal-to-echo-and-noise ratio at microphone 1 in node 2 was -5 dB. Then, the corresponding vector \mathbf{y}_k for each node was transformed to the STFT domain using a square-root Hann window of 512 samples using 50 % overlap. The correlation matrices in (22), (31) and $\{\hat{\mathbf{R}}_{\tilde{\mathbf{y}}_k \tilde{\mathbf{y}}_k}, \hat{\mathbf{R}}_{\tilde{\mathbf{n}}_k \tilde{\mathbf{n}}_k}\}$ in Section VI were computed by selecting the time frames where the desired speech signal was active and not active, respectively, based on an ideal VAD. An ideal VAD was used to isolate the influence of VAD errors. In practice VAD information may be shared among the nodes [31], using a speaker-selective VAD [32] and/or estimating the speech presence probability in a distributed fashion [33]. All nodes in Fig. 1a had a loudspeaker reproducing a speech signal, which were simultaneously active only when the desired speech signal was not active. The second loudspeaker in node 3 was reproducing a music signal which was continuously active. The desired speech signal was produced by a speaker located in the centre of the room. A continuously active localized noise source was also included, producing babble noise.

PK-GEVD-DANSE was run with simultaneous node updating with a relaxation factor $\alpha_{r,S} = 0.9$, to guarantee convergence as suggested in [24]. The ERLE was computed with non-overlapping windows of 1024 samples. The average ERLE (over the time frames) and SNR are shown in Fig. 3. Both metrics were computed for the first microphone in each node. The SNR was computed by filtering the noise component at each microphone signal with the filter obtained for each implementation. The ERLE and SNR when each node works in isolation (ISO) are also shown. PK-GEVD-DANSE is abbreviated to PK-DANSE in the legends for brevity. The NMSE for the three algorithms at the first microphone of each node is shown in Fig. 4.

It can be seen in all nodes that including the PK reduces the error in the estimation of the desired speech signal. In node 3 PK-GEVD-DANSE and PK-MWF outperform MWF in terms of ERLE and SNR. Notice that node 3 is the furthest away from the desired speech source location, it is very close to the noise source location and has two different loudspeaker signals. PK-GEVD-DANSE performs better, in all nodes, in terms of ERLE and SNR than a node working in isolation. The scenario depicted in Fig. 1b was simulated with echo paths of 4096 samples long and a $T_{60} \approx 1.1$ s. The use of previous frames is investigated for $P = 2, 4, 8$ and 16 in the same scenario. Fig. 5 shows the NMSE for PK-GEVD-DANSE with different values of P and the PK-MWF, with frame size of 512 samples. It is observed that for all P , except $P = 64$, PK-GEVD-DANSE outperforms PK-MWF ($P = 1$) at node 2, whereas in node 1 this happens for $P > 2$. In node 2, for $P > 16$ the NMSE was not reduced further, which could be related to the presence of the babble noise and the increasing number of filter coefficients to be estimated. Figure 6 shows

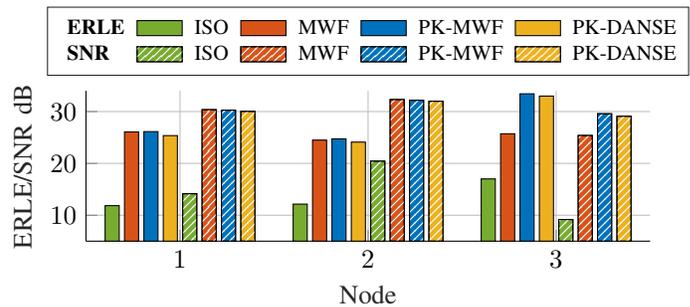


Fig. 3: Average ERLE and SNR computed at the first microphone of each node in Fig. 1a. The ERLE and SNR when the nodes work in isolation (ISO) are also shown.

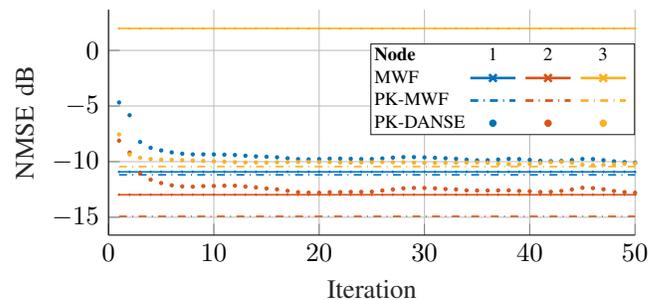


Fig. 4: NMSE at each of the nodes in Fig. 1a.

the NMSE for PK-GEVD-DANSE with different values of P and the PK-MWF, with frame size of 1024 samples. A similar behaviour is observed where the best result is obtained with a filter which effectively has the same number of taps (i.e. $1024 \times \frac{P}{2}$) as the simulated echo path i.e. 4096.

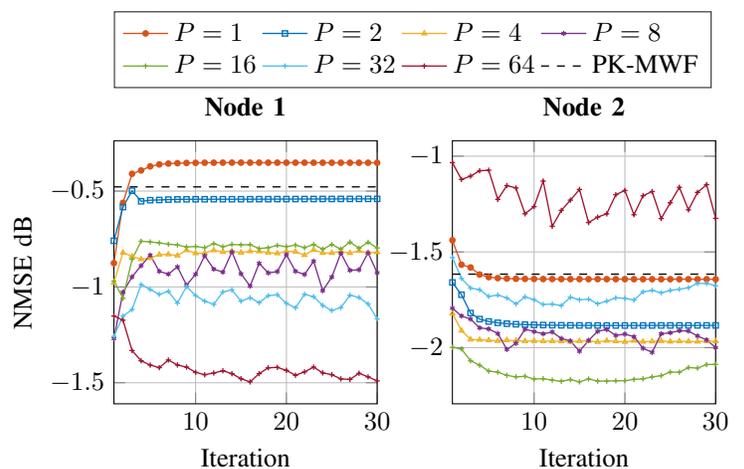


Fig. 5: NMSE for PK-GEVD-DANSE at each node in Fig. 1b with $P = 1, 2, 4, 8, 16$ and a frame size of 512 samples.

C. Per-frame processing

A per-frame processing approach is now used for PK-GEVD-DANSE, where the correlation matrices are updated based on (20), and its NMSE is shown in Fig. 7 and compared

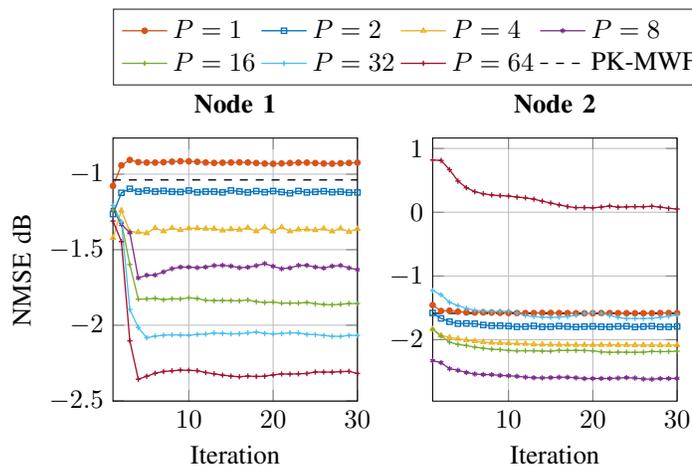


Fig. 6: NMSE for PK-GEVD-DANSE at each node in Fig. 1b with $P = 1, 2, 4, 8, 16$ and a frame size of 1024 samples.

to batch results for MWF and PK-MWF. The high NMSE values before frame 50 are due to the initial updating of the correlation matrices. The closed-form expression for \mathbf{F}_k was used and regularization was applied before inverting the matrix in (46). The forgetting factor for the correlation matrices was chosen such that data from 6 s ago is weighted with a factor of 0.1, based on the following expression

$$\beta = \exp\left(\frac{\ln(0.1)}{6 \frac{F_s}{R/2}}\right) \quad (64)$$

where F_s is the sampling frequency and R the frame size. A per-frame approach is also used for $P > 1$ and it is compared to the batch solution with the same frame size in Fig. 8.

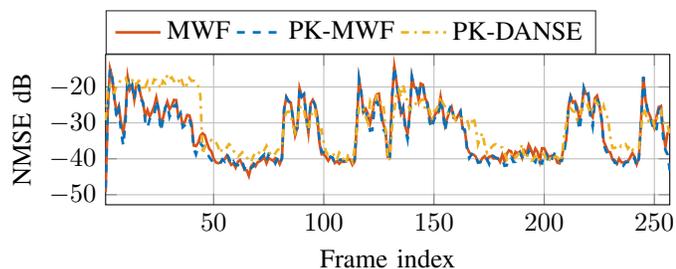


Fig. 7: NMSE at node 2 in Fig. 1b. Frame size of 4096 samples.

D. Adaptive processing

The adaptive simulations presented in this section show the results for the algorithms described in Sections VII-A and VII-B. The results are shown in terms of the NMSE for each frame index. The NMSE for the NLMS-based PK-MWF and PK-GEVD-DANSE are shown in Fig. 9 and 10, respectively. The scenario used is shown in Fig. 1b. The first frames are used to update the *speech-plus-noise* and *noise-only* correlation matrices. The PK-GEVD-DANSE takes a longer time to reduce the NMSE initially (around frame index 50)

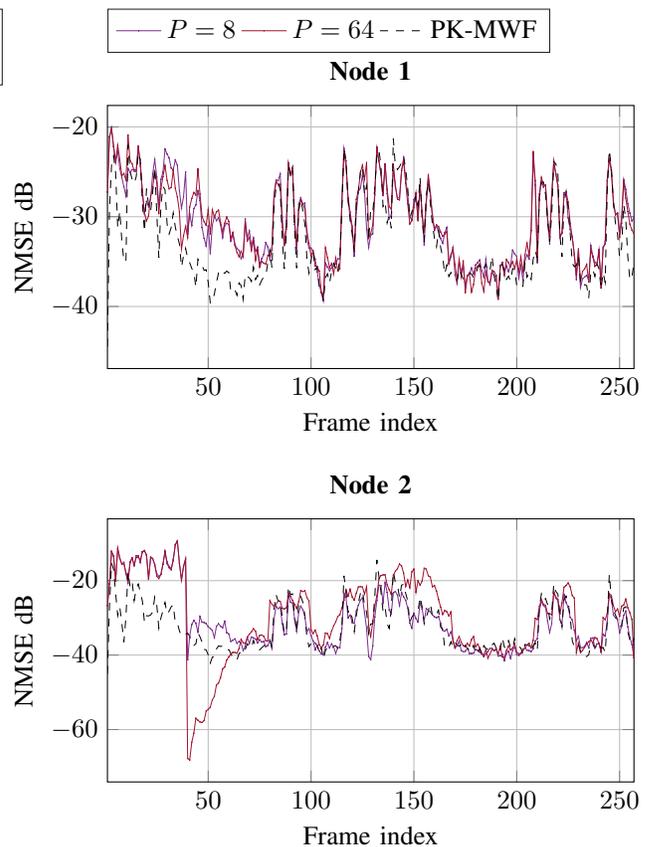


Fig. 8: NMSE for PK-GEVD-DANSE at each of the nodes in Fig. 1b with $P = 8, 64$ and a frame size of 512 samples. The PK-MWF batch solution is shown as reference with same frame size

but after that it reaches lower values than the centralized implementation. It should be noted that the fusion vectors are computed and modified more often than in a batch-implementation which leads to some incorrect entries in the correlation matrices. However once they reach a stable point (around frame index 175), the algorithm performs similarly to the centralized algorithm.

The QRD-based PK-MWF and PK-GEVD-DANSE are shown in Fig. 11 and 12 using the scenario in Fig. 1b. The upper triangular matrix was initialized with a soft-constrained initialization factor as $\delta \mathbf{I}$, where \mathbf{I} is an identity matrix with matching dimensions, according to the implementation. The scenario is time-invariant, i.e., the impulse responses do not change with time, hence the forgetting factor was set to 1 for the PK-MWF. For PK-GEVD-DANSE, the forgetting factor was set to 0.97 due to the fact that the fusion vector in each node is updated over time. Similar results to the NLMS-based PK-MWF are observed for the QRD-based PK-MWF. The results for the QRD-based PK-GEVD-DANSE algorithm are more stable than for the NLMS-based PK-GEVD-DANSE. The NMSE values are more consistent and slightly lower in the desired speech signal parts than for the NLMS-based PK-GEVD-DANSE. It should be noted that the tuning of the NLMS-based algorithms can take some time, which is not the case with the QRD-based algorithms, where only the forgetting

factor and initialization of the upper triangular matrix ($\delta\mathbf{I}$) are needed.

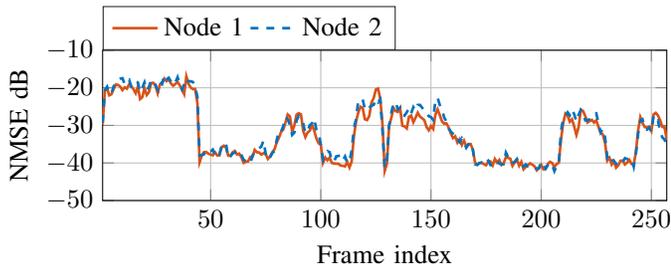


Fig. 9: NMSE at node 1 and 2 using the NLMS-based PK-MWF in scenario depicted in Fig. 1b using a step size $\mu_F = 0.02$, $\delta = 0.0154$ (the sum of the power of the loudspeaker signals) and a frame size of 4096 samples.

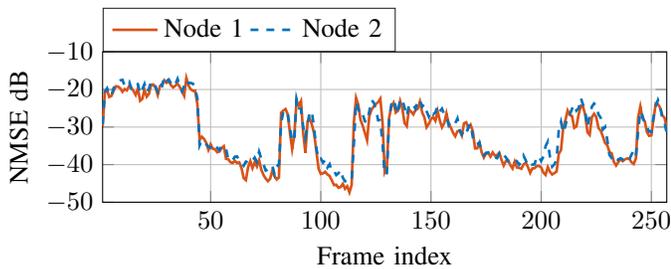


Fig. 10: NMSE at node 1 and 2 using the NLMS-based PK-GEVD-DANSE in scenario depicted in Fig. 1b using a step size $\mu_F = 0.002$, $\delta = 0.0154$ and a frame size of 4096 samples.

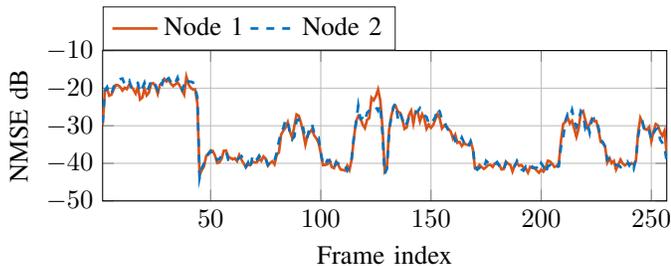


Fig. 11: NMSE at node 1 and 2 using the QRD-based PK-MWF in scenario depicted in Fig. 1b with forgetting factor $\gamma = 1$.

To keep a low algorithmic delay, the frame size was reduced to 512 to estimate the 4096 samples long echo paths in the scenario in Fig. 1b and $P = 8$. The resulting NMSE for both adaptive implementations is shown in Fig. 13. A similar NMSE to that obtained with the respective PK-MWF using a frame size of 4096 is observed.

IX. CONCLUSION

It has been shown that the GEVD-DANSE algorithm from [18] can be adopted for distributed integrated AEC and NR, and that the PK-GEVD-DANSE algorithm from [21] can be adopted for distributed cascade AEC and NR in a WASAN. In

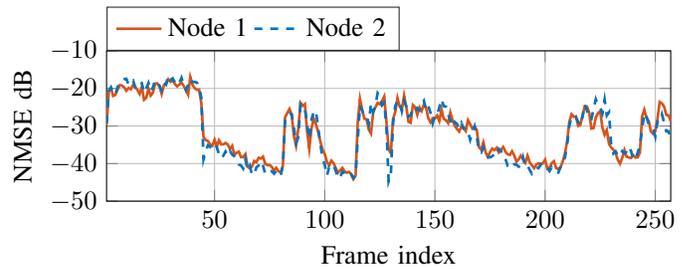


Fig. 12: NMSE at node 1 and 2 using the QRD-based PK-GEVD-DANSE in scenario depicted in Fig. 1b with forgetting factor $\gamma = 0.97$.

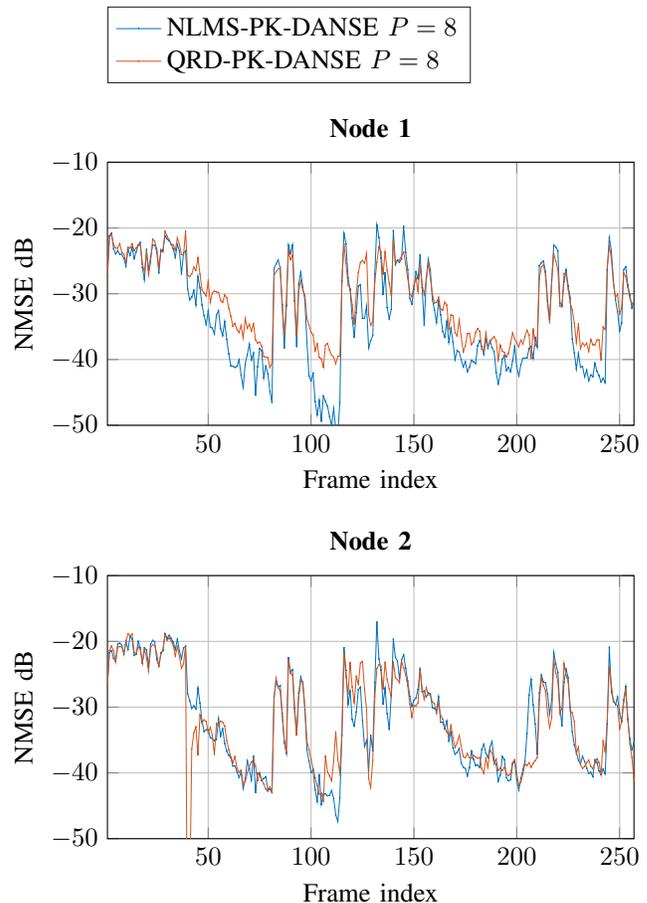


Fig. 13: NMSE for the NLMS-based and QRD-based PK-GEVD-DANSE algorithms at each of the nodes in Fig. 1b with $P = 8$ and a frame size of 512 samples. For the NLMS-based algorithm the step size $\mu_F = 0.001$, $\delta = 10^{-2}$. For the QRD-based algorithm $\gamma = 0.996$.

addition, the communication cost of the PK-GEVD-DANSE algorithm has been reduced resulting in each node in the network broadcasting only 1 fused signal (instead of 2 signals) to the other nodes. The performance of the algorithms has been verified in terms of AEC quantified with the ERLE, as well as in terms of NR quantified with the SNR. As in most cascaded approaches, it has been shown that the AEC stage can be implemented using an NLMS- or QRD-RLS-based algorithm.

REFERENCES

- [1] W. Herbordt, W. Kellermann, and S. Nakamura, ““joint optimization of acoustic echo cancellation and adaptive beamforming”,” in *Topics in acoustic echo and noise control*, E. Hansler and G. Schmidt, Eds. Springer, 2006, pp. 19–50.
- [2] J. Benesty, J. R. Jensen, M. G. Christensen, and J. Chen, *Speech enhancement: A signal subspace perspective*. Elsevier, 2014.
- [3] E. Bohmler, J. Freudenberger, and S. Stenzel, “Combined echo and noise reduction for distributed microphones,” in *Proc. 2011 Joint Workshop Hands-Free Speech Commun. and Microphone Arrays (HSCMA ’11)*. Edinburgh, United Kingdom, 2011, pp. 98–103.
- [4] S. Gustafsson, R. Martin, and P. Vary, “Combined acoustic echo control and noise reduction for hands-free telephony,” *Signal Process.*, vol. 64, no. 1, pp. 21–32, 1998.
- [5] G. Rombouts and M. Moonen, “An integrated approach to acoustic noise and echo cancellation,” *Signal Process.*, vol. 85, no. 4, pp. 849–871, 2005.
- [6] W. Herbordt and W. Kellermann, “Frequency-domain integration of acoustic echo cancellation and a generalized sidelobe canceller with improved robustness,” *European Trans. Telecommun.*, vol. 13, no. 2, pp. 123–132, 2002.
- [7] S. J. Park, C. G. Cho, C. Lee, and D. H. Youn, “Integrated echo and noise canceler for hands-free applications,” *IEEE Trans. Circuits Syst. II, Analog Digit. Signal Process.*, vol. 49, no. 3, pp. 188–195, 2002.
- [8] A. Cohen, A. Barnov, S. Markovich-Golan, and P. Kroon, “Joint beamforming and echo cancellation combining QRD based multichannel AEC and MVDR for reducing noise and non-linear echo,” in *Proc. 26th European Signal Process. Conf. (EUSIPCO ’18)*. Rome, Italy, 2018, pp. 6–10.
- [9] 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, 2006.
- [10] P. Meyer, S. Elshamy, and T. Fingscheidt, “A multichannel Kalman-based Wiener filter approach for speaker interference reduction in meetings,” in *Proc. 2020 IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP ’20)*. Barcelona, Spain (virtual conference), 2020, pp. 451–455.
- [11] K. Nathwani, “Joint acoustic echo and noise cancellation using spectral domain Kalman filtering in double-talk scenario,” in *Proc. 2018 Int. Workshop Acoustic Signal Enhancement (IWAENC ’18)*. Tokyo, Japan, 2018, pp. 1–330.
- [12] A. Fazel, M. El-Khamy, and J. Lee, “CAD-AEC: Context-aware deep acoustic echo cancellation,” in *Proc. 2020 IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP ’20)*. Barcelona, Spain (virtual conference), 2020, pp. 6919–6923.
- [13] H. Zhang, K. Tan, and D. Wang, “Deep learning for joint acoustic echo and noise cancellation with nonlinear distortions,” in *Proc. Interspeech ’19*. Graz, Austria, 2019, pp. 4255–4259.
- [14] H. Seo, M. Lee, and J.-H. Chang, “Integrated acoustic echo and background noise suppression based on stacked deep neural networks,” *Appl. Acoust.*, vol. 133, pp. 194–201, 2018.
- [15] I. F. Akyildiz, T. Melodia, and K. R. Chowdhury, “A survey on wireless multimedia sensor networks,” *Computer Networks*, vol. 51, no. 4, pp. 921–960, 2007.
- [16] Y. Zeng and R. C. Hendriks, “Distributed delay and sum beamformer for speech enhancement via randomized gossip,” *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 22, no. 1, 2013.
- [17] R. Heusdens, G. Zhang, R. C. Hendriks, Y. Zeng, and W. B. Kleijn, “Distributed MVDR beamforming for (wireless) microphone networks using message passing,” in *Proc. 2012 Int. Workshop Acoustic Signal Enhancement (IWAENC ’12)*. Aachen, Germany, 2012.
- [18] A. Bertrand and M. Moonen, “Distributed adaptive node-specific signal estimation in fully connected sensor networks—part II: Simultaneous and asynchronous node updating,” *IEEE Trans. Signal Process.*, vol. 58, no. 10, pp. 5292–5306, 2010.
- [19] E. Ceolini and S. C. Liu, “Combining deep neural networks and beamforming for real-time multi-channel speech enhancement using a wireless acoustic sensor network,” in *Proc. 2019 Int. Workshop Mach. Learning Signal Process. (MLSP ’19)*. Pittsburgh, PA, USA, 2019, pp. 1–6.
- [20] S. Ruiz, T. van Waterschoot, and M. Moonen, “Distributed combined acoustic echo cancellation and noise reduction using GEVD-based distributed adaptive node specific signal estimation with prior knowledge,” in *Proc. 28th European Signal Process. Conf. (EUSIPCO ’20)*. Amsterdam, The Netherlands, 2021.
- [21] R. Van Rompaey and M. Moonen, “Distributed adaptive node-specific signal estimation in a wireless sensor network with partial prior knowledge of the desired source steering vector,” in *Proc. 27th European Signal Process. Conf. (EUSIPCO ’19)*. A Coruna, Spain, 2019.
- [22] R. Serizel, M. Moonen, B. Van Dijk, and J. Wouters, “Low-rank approximation based multichannel Wiener filter algorithms for noise reduction with application in cochlear implants,” *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 22, no. 4, pp. 785–799, 2014.
- [23] J. Benesty, J. Chen, Y. A. Huang, and S. Doclo, “Study of the Wiener filter for noise reduction,” in *Speech Enhancement*. Springer, 2005, pp. 9–41.
- [24] A. Bertrand and M. Moonen, “Robust distributed noise reduction in hearing aids with external acoustic sensor nodes,” *EURASIP J. Adv. Signal Process.*, vol. 2009, p. 12, 2009.
- [25] F. Jabloun and B. Champagne, “Signal subspace techniques for speech enhancement,” in *Speech Enhancement*. Springer, 2005, pp. 135–159.
- [26] S. Doclo, S. Gannot, M. Moonen, and A. Spriet, “9. “acoustic beamforming for hearing aid applications”,” in *Handbook on array processing and sensor networks*, S. Haykin and K. R. Liu, Eds. Wiley Online Library, 2010, pp. 269–302.
- [27] J. Szurley, A. Bertrand, and M. Moonen, “Improved tracking performance for distributed node-specific signal enhancement in wireless acoustic sensor networks,” in *Proc. 2013 IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP ’13)*. Vancouver, Canada, 2013, pp. 336–340.
- [28] S. Haykin, *Adaptive filter theory*. Prentice-Hall, Inc., 1996.
- [29] G. Rombouts and M. Moonen, “QRD-based unconstrained optimal filtering for acoustic noise reduction,” *Signal Process.*, vol. 83, no. 9, pp. 1889–1904, 2003.
- [30] E. De Sena, N. Antonello, M. Moonen, and T. van Waterschoot, “On the modeling of rectangular geometries in room acoustic simulations,” *IEEE/ACM Trans. Audio Speech Lang. Process.*, vol. 23, no. 4, pp. 774–786, 2015.
- [31] V. Berisha, H. Kwon, and A. Spanias, “Real-time implementation of a distributed voice activity detector,” in *Proc. 4th IEEE Workshop Sensor Array Multichannel Process. (SAM ’06)*. Waltham, MA, USA, 2006, pp. 659–662.
- [32] S. Maraboina, D. Kolossa, P. Bora, and R. Orglmeister, “Multi-speaker voice activity detection using ICA and beamforming analysis,” in *Proc. 14th European Signal Process. Conf. (EUSIPCO ’06)*. Florence, Italy, 2006.
- [33] Y. Zhao, J. K. Nielsen, J. Chen, and M. G. Christensen, “Model-based distributed node clustering and multi-speaker speech presence probability estimation in wireless acoustic sensor networks,” *J. Acoust. Soc. Amer.*, vol. 147, no. 6, pp. 4189–4201, 2020.