

Using partial a priori knowledge of relative transfer functions to design an MVDR beamformer for a binaural hearing assistive device with external microphones

Randall ALI⁽¹⁾, Toon vanWATERSCHOOT⁽²⁾, Marc MOONEN⁽³⁾

⁽¹⁾KU Leuven, Dept. of Electrical Engineering (ESAT-STADIUS), Leuven, Belgium, randall.ali@esat.kuleuven.be

⁽²⁾KU Leuven, Dept. of Electrical Engineering (ESAT-ETC/STADIUS), Leuven, Belgium, tvanwate@esat.kuleuven.be

⁽³⁾KU Leuven, Dept. of Electrical Engineering (ESAT-STADIUS), Leuven, Belgium, marc.moonen@esat.kuleuven.be

Abstract

This paper considers a binaural hearing assistive device (HAD) equipped with a separate local microphone array (LMA) for the left and right ear, as well as external microphones (XMs) that may be located within the vicinity of this HAD. For such a system, a binaural minimum variance distortionless response (BMVDR) beamformer may be used for noise reduction, and for the preservation of the relevant binaural speech cues, provided that a reliable estimate of the left and right ear relative transfer function (RTF) vectors pertaining to all the microphones can be obtained. In this paper, an alternative approach is considered, which makes use of available partial a priori knowledge of these RTF vectors, i.e., known separate left and right ear RTF vectors for the respective LMAs on the binaural HAD. The procedure for this approach will be discussed, which requires the estimation of an appropriate scaling between the left and right ear RTF vectors, and the missing part of these RTF vectors pertaining to the XMs. An experiment involving a dummy head, two behind-the-ear dummy hearing aids, and XMs is also performed in order to evaluate the benefit of the proposed approach.

Keywords: Binaural MVDR, External Microphones, Hearing Assistive Device

1 INTRODUCTION

In noisy environments, speech intelligibility is inevitably degraded for individuals that suffer with a hearing impairment and hence hearing assistive devices (HADs) such as hearing aids (HAs) or cochlear implants (CIs) must perform speech enhancement tasks. In addition to the fundamental task of noise reduction, preservation of the binaural cues, i.e., the interaural time differences (ITDs) and interaural level differences (ILDs) is also important to maintain the spatial perception of the auditory scene.

For a binaural HAD equipped with a separate local microphone array (LMA) for the left and right ear, and a communication link between them, the binaural minimum variance distortionless response beamformer (BMVDR) (1) is known to exhibit substantial noise reduction and to preserve the ITD and ILD of a target speaker¹. In recent work (2, 3), such a binaural HAD has also been supplemented with an external microphone (XM) (e.g. a wearable microphone or the microphone on a mobile device) and it was demonstrated that the XM could contribute to additional noise reduction and preserve the relevant binaural cues. For the successful operation of the BMVDR in this case, an estimate of the entire vector of transfer functions from the target signal at a left ear reference microphone to all the other microphones, i.e., the left ear relative transfer function (RTF) vector, and a corresponding right ear RTF vector is required. However, obtaining such estimates becomes increasingly challenging in adverse acoustic conditions.

Therefore in this paper, generalising the system to include more than one XM, an alternative approach is considered, which makes use of available partial a priori knowledge of these RTF vectors, i.e., known separate left and right ear RTF vectors for the respective LMAs on the binaural HAD (4). In such a case, it is only the

¹Although, the BMVDR beamformer preserves the binaural cues for the target speaker, it distorts the binaural cues for the noise. However, in (1), several remedies have been proposed, and hence this work will focus only on the preservation of the binaural cues for the target speaker.

estimation of an appropriate scaling between the left and right ear RTF vectors, and the missing part of these RTF vectors pertaining to the XMs that need to be estimated.

The paper is organised as follows. In Section 2, the data model and notation are described. In Section 3, the state of the art procedure for estimating the entire RTF vector is reviewed. In Section 4, the proposed procedure that makes use of the partial a priori knowledge of the RTF vector is discussed. In Section 5, the proposed procedure is evaluated using recorded audio data, and conclusions are drawn in Section 6.

2 DATA MODEL

The scenario as depicted in Figure 1 is considered, in which a user of a binaural HAD is listening to one target speaker of interest in a noisy, reverberant environment. The binaural HAD consists of an LMA with M_a microphones for the left ear and an LMA with M_a microphones for the right ear. Additionally, there are M_e XMs randomly placed within the room ($M_e = 2$ in Fig.1). In the short-time Fourier transform (STFT) domain,

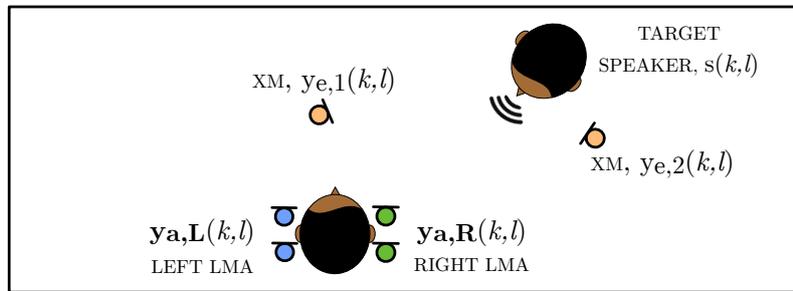


Figure 1. Scenario with a user of a binaural HAD having access to XMs, listening to the target speaker.

the microphone signals at one frequency, k , and one time frame, l , can be stacked into a vector and represented as follows:

$$\mathbf{y}(k,l) = \underbrace{\mathbf{a}(k,l)}_{\mathbf{x}(k,l)} s(k,l) + \mathbf{n}(k,l) \implies \begin{bmatrix} \mathbf{y}_{a,L}(k,l) \\ \mathbf{y}_{a,R}(k,l) \\ \mathbf{y}_e(k,l) \end{bmatrix} = \begin{bmatrix} \mathbf{a}_{a,L}(k,l) \\ \mathbf{a}_{a,R}(k,l) \\ \mathbf{a}_e(k,l) \end{bmatrix} s(k,l) + \begin{bmatrix} \mathbf{n}_{a,L}(k,l) \\ \mathbf{n}_{a,R}(k,l) \\ \mathbf{n}_e(k,l) \end{bmatrix} \quad (1)$$

where² $\mathbf{a}_{a,L} = [a_{a1,L}, a_{a2,L}, \dots, a_{aM_a,L}]^T$, $\mathbf{a}_{a,R} = [a_{a1,R}, a_{a2,R}, \dots, a_{aM_a,R}]^T$, and $\mathbf{a}_e = [a_{e1}, a_{e2}, \dots, a_{eM_e}]^T$ are the acoustic transfer functions (ATFs) from the target speaker to the microphones on the left LMA, the right LMA, and the XMs respectively. Furthermore, s is the target speaker, and $\mathbf{n}_{a,L}$, $\mathbf{n}_{a,R}$, and \mathbf{n}_e are the noise contributions similarly defined as $\mathbf{a}_{a,L}$, $\mathbf{a}_{a,R}$, and \mathbf{a}_e respectively. Without loss of generality, the first microphone in each of the LMAs is also chosen as the reference microphone:

$$y_{a1,L} = \mathbf{e}_L^T \mathbf{y} = s_{a1,L} + n_{a1,L} \quad y_{a1,R} = \mathbf{e}_R^T \mathbf{y} = s_{a1,R} + n_{a1,R} \quad (2)$$

where $s_{a1,L} = a_{a1,L}s$, $s_{a1,R} = a_{a1,R}s$, which are the speech components that need to be estimated, and \mathbf{e}_L and \mathbf{e}_R are all-zero vectors except for a one in the left and right LMA reference microphone position respectively. In order to perform the estimation, it is firstly convenient to re-define eq. (1) in terms of a relative transfer function (RTF) vector, as opposed to the ATF vector, \mathbf{a} . The RTF vector is simply the ATF vector normalised to a reference microphone. Therefore in the binaural context, a separate RTF vector can be defined for the left ear and another for the right ear. Hence, eq. (1) can be expressed as follows:

$$\mathbf{y} = \mathbf{h}_L s_{a1,L} + \mathbf{n} \quad \mathbf{y} = \mathbf{h}_R s_{a1,R} + \mathbf{n} \quad (3)$$

²The dependence on (k,l) is dropped for notational convenience.

where \mathbf{h}_L and \mathbf{h}_R are the RTF vectors defined as:

$$\mathbf{h}_L = \frac{1}{a_{a1,L}} \begin{bmatrix} \mathbf{a}_{a,L} \\ \mathbf{a}_{a,R} \\ \mathbf{a}_e \end{bmatrix} = \begin{bmatrix} \mathbf{h}_{a,L} \\ \varphi \mathbf{h}_{a,R} \\ \mathbf{h}_{e,L} \end{bmatrix} \quad \mathbf{h}_R = \frac{1}{a_{a1,R}} \begin{bmatrix} \mathbf{a}_{a,L} \\ \mathbf{a}_{a,R} \\ \mathbf{a}_e \end{bmatrix} = \begin{bmatrix} \frac{1}{\varphi} \mathbf{h}_{a,L} \\ \mathbf{h}_{a,R} \\ \frac{1}{\varphi} \mathbf{h}_{e,L} \end{bmatrix} \quad (4)$$

where $\mathbf{h}_{a,L} = \frac{\mathbf{a}_{a,L}}{a_{a1,L}}$ and $\mathbf{h}_{a,R} = \frac{\mathbf{a}_{a,R}}{a_{a1,R}}$ are the individual RTF vectors corresponding to each of the left and right LMAs and the complex scaling, $\varphi = \frac{a_{a1,R}}{a_{a1,L}}$. It should be noted that the part of the RTF vector pertaining to the XMs in \mathbf{h}_R is a scaled version of that in \mathbf{h}_L , where $\mathbf{h}_{e,L} = \frac{\mathbf{a}_e}{a_{a1,L}}$. In fact, it can be seen that $\mathbf{h}_L = \varphi \mathbf{h}_R$, which means that the RTF vectors are parallel.

The speech-plus-noise spatial correlation matrix, \mathbf{R}_{yy} , the noise-only correlation matrix, \mathbf{R}_{nn} , and the speech-only correlation matrix, \mathbf{R}_{xx} , all $\in \mathbb{C}^{(2M_a+M_e) \times (2M_a+M_e)}$, are given respectively as:

$$\mathbf{R}_{yy} = \mathbb{E}\{\mathbf{y}\mathbf{y}^H\}; \quad \mathbf{R}_{nn} = \mathbb{E}\{\mathbf{n}\mathbf{n}^H\}; \quad \mathbf{R}_{xx} = \mathbb{E}\{\mathbf{x}\mathbf{x}^H\} \quad (5)$$

where $\mathbb{E}\{\cdot\}$ is the expectation operator, H is the Hermitian transpose, and \mathbf{R}_{xx} is a rank-1 correlation matrix:

$$\mathbf{R}_{xx} = \mathbb{E}\{\mathbf{x}\mathbf{x}^H\} = \sigma_{s_{a1,L}}^2 \mathbf{h}_L \mathbf{h}_L^H = \sigma_{s_{a1,R}}^2 \mathbf{h}_R \mathbf{h}_R^H \quad (6)$$

where $\sigma_{s_{a1,L}}^2 = \mathbb{E}\{|s_{a1,L}|^2\}$ and $\sigma_{s_{a1,R}}^2 = \mathbb{E}\{|s_{a1,R}|^2\}$ are the speech powers in the reference left and right microphone respectively. It is also assumed that the speech components are uncorrelated with the noise components, and hence $\mathbf{R}_{yy} = \mathbf{R}_{xx} + \mathbf{R}_{nn}$. A perfect communication link is additionally assumed among the left and right LMAs in the binaural HAD, and the XMs, with no bandwidth constraints and synchronous sampling.

The estimate of the speech component in the reference microphone of the left and right LMAs, i.e. the estimate of $s_{a1,L}$ and $s_{a1,R}$ is then obtained through the linear filtering of the microphone signals, with the complex-valued filters, \mathbf{w}_L and \mathbf{w}_R respectively:

$$\hat{s}_{a1,L} = \mathbf{w}_L^H \mathbf{y} \quad \hat{s}_{a1,R} = \mathbf{w}_R^H \mathbf{y} \quad (7)$$

The BMVDR beamformer filters, \mathbf{w}_L and \mathbf{w}_R , are then given by:

$$\mathbf{w}_L = \frac{\mathbf{R}_{nn}^{-1} \mathbf{h}_L}{\mathbf{h}_L^H \mathbf{R}_{nn}^{-1} \mathbf{h}_L} \quad \mathbf{w}_R = \frac{\mathbf{R}_{nn}^{-1} \mathbf{h}_R}{\mathbf{h}_R^H \mathbf{R}_{nn}^{-1} \mathbf{h}_R} \quad (8)$$

Consequently, in order to compute these filters, estimates are required for \mathbf{R}_{nn} , and the RTFs, \mathbf{h}_L , and \mathbf{h}_R . Typically, \mathbf{R}_{nn} can be estimated during periods of noise only with recursive averaging (3). Hence this paper focuses on the estimation of \mathbf{h}_L and \mathbf{h}_R .

3 ESTIMATING THE ENTIRE RTF VECTOR

Given $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$, which are estimates of \mathbf{R}_{yy} and \mathbf{R}_{nn} respectively, a generalised eigenvalue decomposition (GEVD) (5) or what is equivalently known as covariance whitening (3) can be used to estimate \mathbf{h}_L and \mathbf{h}_R . A spatial pre-whitening operation can be firstly defined from $\hat{\mathbf{R}}_{nn}$ using the Cholesky decomposition:

$$\hat{\mathbf{R}}_{nn} = \hat{\mathbf{R}}_{nn}^{1/2} \hat{\mathbf{R}}_{nn}^{H/2} \quad (9)$$

where $\hat{\mathbf{R}}_{nn}^{1/2}$ is a lower triangular matrix. Spatial pre-whitening is then performed by pre-multiplying the signal vector of interest by $\hat{\mathbf{R}}_{nn}^{-1/2}$. For an autocorrelation matrix, spatial pre-whitening is performed by pre-multiplying it by $\hat{\mathbf{R}}_{nn}^{-1/2}$ and post-multiplying it by $\hat{\mathbf{R}}_{nn}^{H/2}$. Using the definition of \mathbf{R}_{xx} from eq. (6) and that $\mathbf{R}_{yy} = \mathbf{R}_{xx} + \mathbf{R}_{nn}$, the following optimisation problem can be considered to estimate \mathbf{h}_L (and \mathbf{h}_R by an appropriate scaling):

$$\min_{\sigma_{s_{a1,L}}^2, \mathbf{h}_L} \|\hat{\mathbf{R}}_{nn}^{-1/2} ((\hat{\mathbf{R}}_{yy} - \hat{\mathbf{R}}_{nn}) - \sigma_{s_{a1,L}}^2 \mathbf{h}_L \mathbf{h}_L^H) \hat{\mathbf{R}}_{nn}^{H/2}\|_F^2 \quad (10)$$

where $\|\cdot\|_F$ is the frobenius norm. The solution to eq. (10) then follows from an eigenvalue decomposition (EVD) of $\hat{\mathbf{R}}_{\text{nn}}^{-1/2} \hat{\mathbf{R}}_{\text{yy}} \hat{\mathbf{R}}_{\text{nn}}^{-H/2}$ or equivalently, GEVD of the matrix pencil $\{\hat{\mathbf{R}}_{\text{yy}}, \hat{\mathbf{R}}_{\text{nn}}\}$:

$$\hat{\mathbf{R}}_{\text{nn}}^{-1} \hat{\mathbf{R}}_{\text{yy}} = \mathbf{U} \mathbf{\Sigma} \mathbf{U}^{-1} \quad (11)$$

where $\mathbf{\Sigma}$ is a diagonal matrix of the generalised eigenvalues arranged in descending order, and \mathbf{U} is an invertible matrix containing the corresponding generalised eigenvectors. The GEVD is also equivalent to a joint diagonalisation of $\hat{\mathbf{R}}_{\text{yy}}$ and $\hat{\mathbf{R}}_{\text{nn}}$:

$$\hat{\mathbf{R}}_{\text{yy}} = \mathbf{Q} \mathbf{\Sigma}_y \mathbf{Q}^H \quad \hat{\mathbf{R}}_{\text{nn}} = \mathbf{Q} \mathbf{\Sigma}_n \mathbf{Q}^H \quad (12)$$

where $\mathbf{\Sigma}_y$ and $\mathbf{\Sigma}_n$ are diagonal matrices, and $\mathbf{Q} = \mathbf{U}^{-H}$ is an invertible matrix. A rank-1 approximation to $(\hat{\mathbf{R}}_{\text{yy}} - \hat{\mathbf{R}}_{\text{nn}}) = \mathbf{Q}(\mathbf{\Sigma}_y - \mathbf{\Sigma}_n)\mathbf{Q}^H$ yields an estimate for \mathbf{R}_{xx} , $\hat{\mathbf{R}}_{\text{xx}} = \mathbf{Q} \mathbf{e}_1 \mathbf{e}_1^T (\mathbf{\Sigma}_y - \mathbf{\Sigma}_n) \mathbf{e}_1 \mathbf{e}_1^T \mathbf{Q}^H$, where $\mathbf{e}_1 \in \mathbb{C}^{2M_a + M_e}$ is an all-zero vector except for a one as the first element (and it is noted that $\mathbf{e}_1 = \mathbf{e}_L$). It can be shown (5) that this corresponds to the rank-1 approximation sought from eq. (10) so that the estimates to \mathbf{h}_L and \mathbf{h}_R then follow as:

$$\hat{\mathbf{h}}_L = \frac{\mathbf{Q} \mathbf{e}_1}{\mathbf{e}_L^T \mathbf{Q} \mathbf{e}_1} \quad \hat{\mathbf{h}}_R = \frac{\mathbf{Q} \mathbf{e}_1}{\mathbf{e}_R^T \mathbf{Q} \mathbf{e}_1} \quad (13)$$

Finally, a substitution of $\hat{\mathbf{R}}_{\text{nn}}$ from eq. (12) and $\hat{\mathbf{h}}_L$ and $\hat{\mathbf{h}}_R$ from eq. (13) into eq. (8) results in the corresponding BMVDR filters:

$$\hat{\mathbf{w}}_L = \mathbf{U} \mathbf{e}_1 \mathbf{e}_1^T \mathbf{Q}^H \mathbf{e}_L \quad \hat{\mathbf{w}}_R = \mathbf{U} \mathbf{e}_1 \mathbf{e}_1^T \mathbf{Q}^H \mathbf{e}_R \quad (14)$$

4 USING PARTIAL A PRIORI KNOWLEDGE OF THE RTF VECTOR

As opposed to estimating the entire RTF vectors, \mathbf{h}_L , and \mathbf{h}_R , an alternative procedure may be followed if there is a priori knowledge of the RTF vectors for the separate left and right LMA, i.e., if a suitable approximation to $\mathbf{h}_{a,L}$ and $\mathbf{h}_{a,R}$ is available. For instance, such an approximation may be the measured RTF vectors for the separate left and right LMA in an anechoic room or RTF vectors from an existing binaural noise reduction system that uses only the LMAs. Denoting this approximation to $\mathbf{h}_{a,L}$ and $\mathbf{h}_{a,R}$ as $\tilde{\mathbf{h}}_{a,L}$ and $\tilde{\mathbf{h}}_{a,R}$ respectively, and recalling the definitions from eq. (4), an alternative optimisation problem to eq. (10) can be considered:

$$\min_{\sigma_{\text{sa},L}^2, \varphi, \mathbf{h}_{e,L}} \|\hat{\mathbf{R}}_{\text{nn}}^{-1/2} ((\hat{\mathbf{R}}_{\text{yy}} - \hat{\mathbf{R}}_{\text{nn}}) - \sigma_{\text{sa},L}^2 \begin{bmatrix} \tilde{\mathbf{h}}_{a,L} \\ \varphi \mathbf{h}_{a,R} \\ \mathbf{h}_{e,L} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{h}}_{a,L}^H & \varphi^* \tilde{\mathbf{h}}_{a,R}^H & \mathbf{h}_{e,L}^H \end{bmatrix}) \hat{\mathbf{R}}_{\text{nn}}^{-H/2}\|_F^2 \quad (15)$$

where now it is only the scaling, φ , and the RTF vector for the XMs, $\mathbf{h}_{e,L}$, which need to be found as opposed to the entire \mathbf{h}_L as in eq. (10). As will be discussed in the following, the solution can be realised in the block scheme of Figure 2, which consists of compressing the left and right LMA signals, an orthogonalisation operation, and finally a GEVD on a lower dimensional ($\mathbb{C}^{(M_e+2) \times (M_e+2)}$) matrix pencil. In order to solve eq. (15), the following blocking matrices, $\mathbf{C}_a \in \mathbb{C}^{2M_a \times (2M_a-2)}$, $\mathbf{C}_{a,L} \in \mathbb{C}^{M_a \times (M_a-1)}$, $\mathbf{C}_{a,R} \in \mathbb{C}^{M_a \times (M_a-1)}$, fixed beamformers, $\mathbf{F}_a \in \mathbb{C}^{2M_a \times 2}$, $\mathbf{f}_{a,L} \in \mathbb{C}^{M_a}$, $\mathbf{f}_{a,R} \in \mathbb{C}^{M_a}$, and transformation matrix, $\mathbf{T} \in \mathbb{C}^{(2M_a+M_e) \times (2M_a+M_e)}$ are firstly defined:

$$\mathbf{C}_a = \begin{bmatrix} \mathbf{C}_{a,L} & \mathbf{0} \\ \mathbf{0} & \mathbf{C}_{a,R} \end{bmatrix} \quad \mathbf{F}_a = \begin{bmatrix} \mathbf{f}_{a,L} & \mathbf{0} \\ \mathbf{0} & \mathbf{f}_{a,R} \end{bmatrix} \quad \mathbf{T} = \left[\begin{array}{c|c|c} \mathbf{C}_a & \mathbf{F}_a & \mathbf{0} \\ \hline \mathbf{0} & \mathbf{0} & \mathbf{I}_{M_e} \end{array} \right] \quad (16)$$

$$\mathbf{C}_{a,L}^H \tilde{\mathbf{h}}_{a,L} = \mathbf{0}; \mathbf{C}_{a,R}^H \tilde{\mathbf{h}}_{a,R} = \mathbf{0} \quad \mathbf{f}_{a,L}^H \tilde{\mathbf{h}}_{a,L} = 1; \mathbf{f}_{a,R}^H \tilde{\mathbf{h}}_{a,R} = 1$$

where $\mathbf{I}_{M_e} \in \mathbb{C}^{M_e \times M_e}$ is an identity matrix. The first two blocks in Fig. 2 apply the transformation, \mathbf{T}^H , to \mathbf{y} to yield a set of blocking matrix signals, $\mathbf{C}_a^H \mathbf{y}_a \in \mathbb{C}^{2M_a-2}$, two compressed signals, $\mathbf{f}_{a,L}^H \mathbf{y}_{a,L}$ and $\mathbf{f}_{a,R}^H \mathbf{y}_{a,R}$ resulting from the left and right fixed beamformers respectively, and the unaltered set of XM signals, \mathbf{y}_e . An alternative spatial pre-whitening operation can then be defined by applying the transformation to $\hat{\mathbf{R}}_{\text{nn}}$:

$$\mathbf{T}^H \hat{\mathbf{R}}_{\text{nn}} \mathbf{T}^H = \mathbf{L} \mathbf{L}^H \quad (17)$$

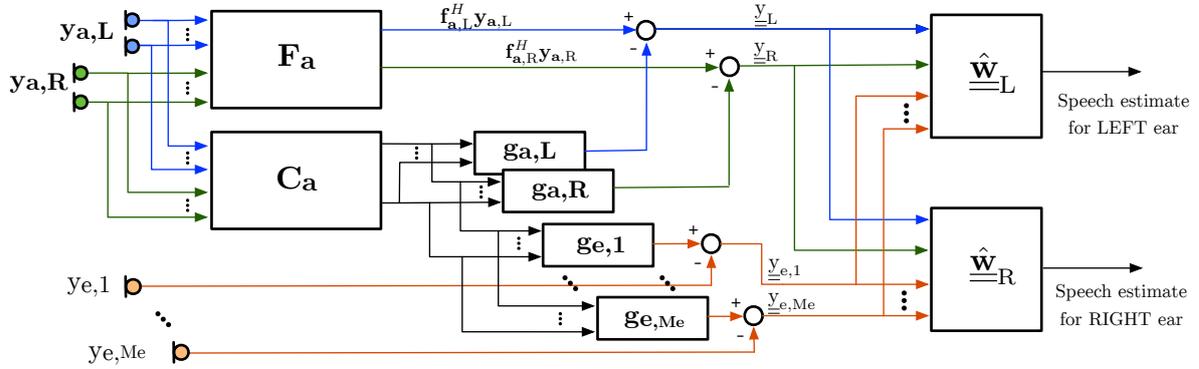


Figure 2. Block scheme for a BMVDR that uses partial a priori knowledge of the RTF vectors.

where \mathbf{L} is a lower triangular matrix. It can then be shown (6) that eq. (15) can be equivalently re-written as:

$$\min_{\sigma_{s_{a1,L}}^2, \varphi, \mathbf{h}_{e,L}} \|\mathbf{L}^{-1} \mathbf{T}^H ((\hat{\mathbf{R}}_{yy} - \hat{\mathbf{R}}_{nn}) - \sigma_{s_{a1,L}}^2 \begin{bmatrix} \tilde{\mathbf{h}}_{a,L} \\ \tilde{\boldsymbol{\varphi}}_{\mathbf{h}_{a,R}} \\ \mathbf{h}_{e,L} \end{bmatrix} \begin{bmatrix} \tilde{\mathbf{h}}_{a,L}^H & \tilde{\boldsymbol{\varphi}}_{\mathbf{h}_{a,R}}^H & \mathbf{h}_{e,L}^H \end{bmatrix}) \mathbf{T} \mathbf{L}^{-H}\|_F^2 \quad (18)$$

from which the solution follows eventually from a GEVD of a matrix pencil consisting of the lower dimensional $\mathbb{C}^{(M_e+2) \times (M_e+2)}$ correlation matrices (6) (7) corresponding to the $M_e + 2$ compressed signals, $\underline{\mathbf{y}} = [\underline{y}_L \ \underline{y}_R \ \underline{y}_{e,1} \ \dots \ \underline{y}_{e,M_e}]^T \in \mathbb{C}^{M_e+2}$. These compressed signals can be computed by performing an orthogonalisation involving the previously transformed signals, or equivalently, with a GSC beamformer, \mathbf{A} :

$$\underline{\mathbf{y}} = \mathbf{A}^H \mathbf{y} = \mathbf{F}^H \mathbf{y} - \mathbf{G}^H \mathbf{C}_a^H \mathbf{E}^T \mathbf{y} \quad (19)$$

with $\mathbf{A} = \mathbf{F} - \mathbf{E} \mathbf{C}_a \mathbf{G}$, $\mathbf{F} \in \mathbb{C}^{(2M_a+M_e) \times (M_e+2)}$ defined as $\mathbf{F} = \begin{bmatrix} \mathbf{F}_a & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{M_e} \end{bmatrix}$, the selection matrix, $\mathbf{E} = [\mathbf{I}_{2M_a} \mid \mathbf{0}_{(2M_a \times M_e)}]^T$, and $\mathbf{G} \in \mathbb{C}^{(2M_a-2) \times (M_e+2)}$, which can be computed as follows:

$$\mathbf{G} = (\mathbf{C}_a^H \mathbf{E}^T \hat{\mathbf{R}}_{nn} \mathbf{E} \mathbf{C}_a)^{-1} \mathbf{C}_a^H \mathbf{E}^T \hat{\mathbf{R}}_{nn} \mathbf{F} = [\mathbf{g}_{a,L} \ \mathbf{g}_{a,R} \ \mathbf{g}_{e,1} \ \dots \ \mathbf{g}_{e,M_e}] \quad (20)$$

is comprised of the \mathbb{C}^{M_e+2} orthogonalisation filters in each column. As depicted in Fig. 2, $\mathbf{g}_{a,L}$ and $\mathbf{g}_{a,R}$ will orthogonalise the noise components of the signals $\mathbf{f}_{a,L}^H \mathbf{y}_{a,L}$ and $\mathbf{f}_{a,R}^H \mathbf{y}_{a,R}$ onto the noise components of the blocking matrix signals, and $[\mathbf{g}_{e,1} \ \dots \ \mathbf{g}_{e,M_e}]$ will orthogonalise the noise components of each of the respective XMs onto the noise components of the blocking matrix signals.

A new signal model that resembles eq. (3) can in fact then be realised by substituting eq. (3) into eq. (19):

$$\underline{\mathbf{y}} = \underline{\mathbf{h}}_L s_{a1,L} + \underline{\mathbf{n}} \quad \underline{\mathbf{y}} = \underline{\mathbf{h}}_R s_{a1,R} + \underline{\mathbf{n}} \quad (21)$$

where $\underline{\mathbf{h}}_L = \mathbf{A}^H \mathbf{h}_L = [1 \ \varphi \ \mathbf{h}_{e,L}^T]^T$, $\underline{\mathbf{h}}_R = \mathbf{A}^H \mathbf{h}_R = [\frac{1}{\varphi} \ 1 \ \frac{1}{\varphi} \mathbf{h}_{e,L}^T]^T$, $\underline{\mathbf{n}} = \mathbf{A}^H \mathbf{n}$, and the associated correlation matrices, $\underline{\mathbf{R}}_{yy} = \mathbb{E}\{\underline{\mathbf{y}}\underline{\mathbf{y}}^H\}$ and $\underline{\mathbf{R}}_{nn} = \mathbb{E}\{\underline{\mathbf{n}}\underline{\mathbf{n}}^H\}$. These compressed signals, $\underline{\mathbf{y}}$, can now be used to design the BMVDR as opposed to \mathbf{y} . Estimates of $s_{a1,L}$ and $s_{a1,R}$ can be obtained by directly filtering $\underline{\mathbf{y}}$ with the BMVDR filters, $\underline{\mathbf{w}}_L \in \mathbb{C}^{M_e+2}$ and $\underline{\mathbf{w}}_R \in \mathbb{C}^{M_e+2}$ respectively (as in eq. (7) and depicted in Fig. 2):

$$\underline{\mathbf{w}}_L = \frac{\underline{\mathbf{R}}_{nn}^{-1} \underline{\mathbf{h}}_L}{\underline{\mathbf{h}}_L^H \underline{\mathbf{R}}_{nn}^{-1} \underline{\mathbf{h}}_L} \quad \underline{\mathbf{w}}_R = \frac{\underline{\mathbf{R}}_{nn}^{-1} \underline{\mathbf{h}}_R}{\underline{\mathbf{h}}_R^H \underline{\mathbf{R}}_{nn}^{-1} \underline{\mathbf{h}}_R} \quad (22)$$

Hence, similar to eq. (8), estimates are now required for $\underline{\mathbf{R}}_{nn}$, $\underline{\mathbf{h}}_L$, and $\underline{\mathbf{h}}_R$ in order to compute these filters. With the estimates of the correlation matrices, $\hat{\underline{\mathbf{R}}}_{yy} = \mathbf{A}^H \hat{\underline{\mathbf{R}}}_{yy} \mathbf{A}$ and $\hat{\underline{\mathbf{R}}}_{nn} = \mathbf{A}^H \hat{\underline{\mathbf{R}}}_{nn} \mathbf{A}$, the GEVD procedure from Section 3 then follows from the matrix pencil $\{\hat{\underline{\mathbf{R}}}_{yy}, \hat{\underline{\mathbf{R}}}_{nn}\}$ to estimate $\underline{\mathbf{h}}_L$ and $\underline{\mathbf{h}}_R$:

$$\hat{\underline{\mathbf{R}}}_{nn}^{-1} \hat{\underline{\mathbf{R}}}_{yy} = \underline{\mathbf{U}} \underline{\mathbf{\Sigma}} \underline{\mathbf{U}}^{-1}; \quad \hat{\underline{\mathbf{R}}}_{yy} = \underline{\mathbf{Q}} \underline{\mathbf{\Sigma}}_y \underline{\mathbf{Q}}^H; \quad \hat{\underline{\mathbf{R}}}_{nn} = \underline{\mathbf{Q}} \underline{\mathbf{\Sigma}}_n \underline{\mathbf{Q}}^H \quad (23)$$

where $\underline{\mathbf{\Sigma}}$ is a diagonal matrix of the generalised eigenvalues arranged in descending order, and $\underline{\mathbf{U}}$ contains the corresponding generalised eigenvectors, $\underline{\mathbf{\Sigma}}_y$, and $\underline{\mathbf{\Sigma}}_n$ are also diagonal matrices, and $\underline{\mathbf{Q}} = \underline{\mathbf{U}}^{-H}$ is an invertible matrix. The estimates for $\underline{\mathbf{h}}_L$ and $\underline{\mathbf{h}}_R$ then follow as:

$$\hat{\underline{\mathbf{h}}}_L = \frac{\underline{\mathbf{Q}} \underline{\mathbf{e}}_1}{\underline{\mathbf{e}}_1^T \underline{\mathbf{Q}} \underline{\mathbf{e}}_1} \quad \hat{\underline{\mathbf{h}}}_R = \frac{\underline{\mathbf{Q}} \underline{\mathbf{e}}_1}{\underline{\mathbf{e}}_1^T \underline{\mathbf{Q}} \underline{\mathbf{e}}_1} \quad (24)$$

where the selection vectors, $\underline{\mathbf{e}}_1 = \underline{\mathbf{e}}_L = [1, 0, \mathbf{0}_{(1 \times M_e)}]^T$ and $\underline{\mathbf{e}}_R = [0, 1, \mathbf{0}_{(1 \times M_e)}]^T$. Finally, the substitution of $\hat{\underline{\mathbf{R}}}_{nn}$ from eq. (23) and $\hat{\underline{\mathbf{h}}}_L$ and $\hat{\underline{\mathbf{h}}}_R$ from eq. (24) into eq. (22) results in the corresponding BMVDR filters:

$$\hat{\underline{\mathbf{w}}}_L = \underline{\mathbf{U}} \underline{\mathbf{e}}_1 \underline{\mathbf{e}}_1^T \underline{\mathbf{Q}}^H \underline{\mathbf{e}}_L \quad \hat{\underline{\mathbf{w}}}_R = \underline{\mathbf{U}} \underline{\mathbf{e}}_1 \underline{\mathbf{e}}_1^T \underline{\mathbf{Q}}^H \underline{\mathbf{e}}_R \quad (25)$$

5 RESULTS AND DISCUSSION

In order to evaluate these algorithms, audio recordings were made in an L-shaped room of height 3.8 m, whose longer dimensions were approximately 6 m \times 5.5 m, with an estimated broadband reverberation time of 1.5 s. Similar to Fig. 1, a Neumann KU-100 dummy head was placed in a central location of the room and equipped with two behind-the-ear (BTE) hearing aids (HAs), each consisting of an LMA with two microphones spaced approximately 1.3 cm apart. Denoting 0° as the azimuth direction directly in front of the dummy head and positive angles as clockwise, two XMs (AKG-CK-97-O) were positioned such that one was at 60° and 1 m away from the dummy head, and the other was at 20 cm below the dummy head. A Genelec 8030C loudspeaker, 1 m from the dummy head, was used to generate the speech signal from a target male speaker (8). Four of the same loudspeakers, facing away from the dummy head toward opposite corners, were used to generate uncorrelated excerpts of babble noise to approximate a diffuse noise field. The speech and noises were recorded separately, but mixed together at an input signal-to-noise ratio (SNR) of 0 dB at the reference (frontal) microphone on the right LMA.

Four algorithms were then evaluated: (i) the BMVDR from Section 3 (BMVDR- $\hat{\mathbf{h}}$), (ii) the BMVDR from Section 4 (BMVDR- $\tilde{\mathbf{h}}$), (iii) the BMVDR from Section 3 but only using the left and right LMA (BMVDR- $\hat{\mathbf{h}}_a$), and (iv) the BMVDR as from Section 4 but only using the left and right LMA (BMVDR- $\tilde{\mathbf{h}}_a$). For the processing of the algorithms, the Weighted Overlap and Add (WOLA) method (9), with a Discrete Fourier Transform (DFT) size of 128 samples, 50% overlap, and a sampling frequency of 16 kHz was used. Using the speech presence probability (SPP) (10), periods for which the speech was active were extracted if the SPP > 0.6, and not active if the SPP < 0.4. $\hat{\underline{\mathbf{R}}}_{yy}$ and $\hat{\underline{\mathbf{R}}}_{nn}$ were then estimated accordingly via recursive averaging with an averaging time of 1 s. Finally, the a priori RTFs, $\tilde{\mathbf{h}}_{a,L}$ and $\tilde{\mathbf{h}}_{a,R}$ were computed from the DFT of 128 samples (including only the direct path component) of the measured impulse responses (IRs) from the target speaker position to each of the left and right LMAs. The metrics used to evaluate the following experiments were the change in speech intelligibility-weighted SNR improvement (11) (Δ SI-SNR) from the input SI-SNR at LM1, and change in binaural cue errors, Δ ITD and Δ ILD, using the auditory model from (12). As the room was quite reverberant, the reference binaural signal for computing Δ ITD and Δ ILD was the speech signal of the target speaker convolved with the time domain direct path IRs used to define $\tilde{\mathbf{h}}_{a,L}$ and $\tilde{\mathbf{h}}_{a,R}$. All metrics were computed in 2 s time frames with a 50% overlap.

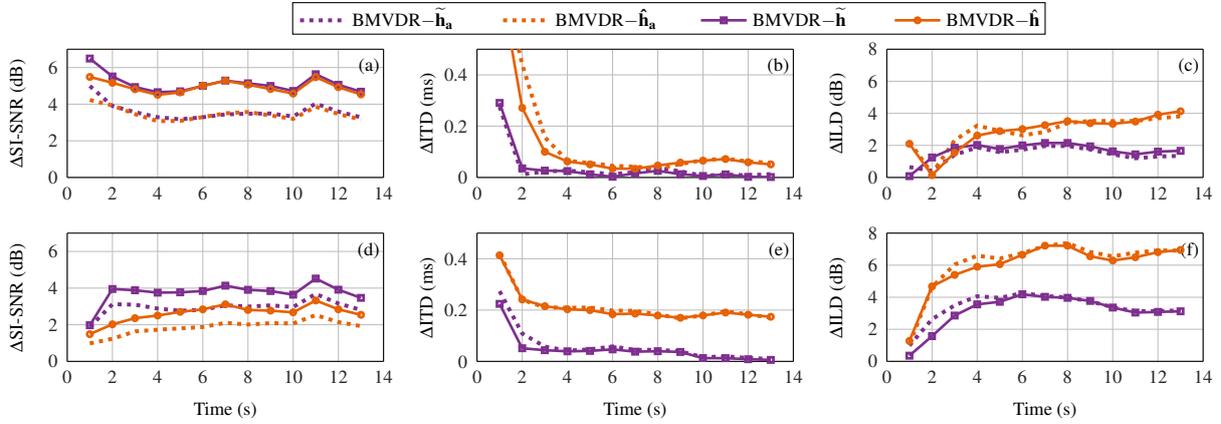


Figure 3. Performance of the various BMVDR algorithms for a target speaker located at 45° . The upper plots ((a),(b),(c)) are the results for accurately estimated correlation matrices, while the lower plots ((d),(e),(f)) are the results for less accurately estimated correlation matrices.

Figure 3 displays the results for an input signal of 13 s when the target speaker was at 45° . The upper plots ((a),(b),(c)) display the results when the SPP was computed on the speech-only signal as received by the reference microphone of the right LMA, and represents a case when $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ have been accurately estimated. In the lower plots ((d),(e),(f)), the SPP was computed using the actual noisy signal as received by the same microphone and represents a case when $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ have not been as accurately estimated.

Firstly, in terms of $\Delta\text{SI-SNR}$, it can be observed that the $\text{BMVDR}-\hat{\mathbf{h}}$ and $\text{BMVDR}-\tilde{\mathbf{h}}$ (i.e., algorithms that use the LMAs along with the XMs) always demonstrate an improvement over $\text{BMVDR}-\hat{\mathbf{h}}_a$ and $\text{BMVDR}-\tilde{\mathbf{h}}_a$ (i.e., algorithms that use the LMAs only). Furthermore, the ΔITD and ΔILD remain relatively constant between $\text{BMVDR}-\hat{\mathbf{h}}_a$ and $\text{BMVDR}-\tilde{\mathbf{h}}$, as well as between $\text{BMVDR}-\hat{\mathbf{h}}$ and $\text{BMVDR}-\tilde{\mathbf{h}}$, suggesting that the inclusion of the XMs does not contribute to any additional distortion in the binaural cues.

The impact of accurate estimation of $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ is also evident. In the upper plots ((a),(b),(c)), where $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ have been accurately estimated, both the $\text{BMVDR}-\tilde{\mathbf{h}}$ and $\text{BMVDR}-\hat{\mathbf{h}}$ perform quite similarly. In the lower plots ((d),(e),(f)), where $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$ have been estimated less accurately, the absolute values of all metrics now indicate a generally lower performance for all algorithms. However, in such a case, it is seen that the $\text{BMVDR}-\tilde{\mathbf{h}}$ offers an improved $\Delta\text{SI-SNR}$ and maintenance of the binaural cues in comparison with the $\text{BMVDR}-\hat{\mathbf{h}}$. Therefore, in cases of inaccurate estimation of $\hat{\mathbf{R}}_{yy}$ and $\hat{\mathbf{R}}_{nn}$, the approach proposed in this paper is seen to be beneficial. The resulting audio files from this experiment may be listened to for subjective evaluation (13).

6 CONCLUSIONS

An approach to designing a BMVDR beamformer for a binaural HAD consisting of separate LMAs for the left and right ear, as well as XMs has been developed. As opposed to the conventional approach of estimating entire RTF vectors, partial a priori knowledge of these RTF vectors, i.e., known separate left and right ear RTF vectors for the LMAs on the binaural HAD was incorporated so that only the computation of an appropriate scaling and the missing part of these RTF vectors pertaining to the XMs was required. An experiment with a dummy head, two behind-the-ear dummy hearing aids, and XMs indicated that the proposed approach is especially beneficial in cases where there may be inaccurate estimation of the relevant correlation matrices needed for the BMVDR beamformer.

ACKNOWLEDGEMENTS

This research work was carried out at the ESAT Laboratory of KU Leuven, in the frame of KU Leuven C2-16-00449 'Distributed Digital Signal Processing for Ad-hoc Wireless Local Area Audio Networking', and KU Leuven Internal Funds VES/19/004. The research leading to these results has received funding from 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.

REFERENCES

- [1] Doclo, S.; Gannot, S.; Marquardt, D.; Hadad, E. Ch. 18: Binaural Speech Processing with Application to Hearing Devices, Audio Source Separation and Speech Enhancement, Wiley, 1st Edition, 2018.
- [2] Szurley J.; Bertrand A.; van Dijk B.; Moonen M. Binaural noise cue preservation in a binaural noise reduction system with a remote microphone signal. *IEEE/ACM Trans. Audio Speech Lang. Process.*, 24(5), 2016, pp 952–966.
- [3] Gößling, N.; Doclo, S. RTF-Based Binaural MVDR Beamformer Exploiting an External Microphone in a Diffuse Noise Field. *Proc. ITG Conference on Speech Communication*, Oldenburg, Germany, 2018.
- [4] Koutrouvelis A.; Hendriks R.; Heusdens R.; Jensen J.; Guo M. Binaural beamforming using pre-determined relative acoustic transfer functions. *Proc. 25th European Signal Process. Conf. (EUSIPCO '17)*, Kos island, Greece, 2017.
- [5] Serizel R.; Moonen M.; van Dijk B.; Wouters J. Low-rank Approximation Based Multichannel Wiener Filter Algorithms for Noise Reduction with Application in Cochlear Implants. *IEEE/ACM Trans Audio, Speech, Lang. Process.*, 22(4), 2014, pp 785–99.
- [6] Ali R.; Bernardi G.; van Waterschoot T.; Moonen M. Methods of extending a generalised sidelobe canceller with external microphones. *IEEE/ACM Trans Audio, Speech, Lang. Process.*, to appear, 2019.
- [7] Van Rompaey, R.; Moonen, M. Distributed Adaptive Node-Specific Signal Estimation in Wireless Sensor Networks with Partial Prior Knowledge of the Desired Source Steering Vector. Internal Report 19-29, ESAT-STADIUS, KU Leuven (Leuven, Belgium). <ftp://ftp.esat.kuleuven.be/sista/rvanromp/19-29.pdf>.
- [8] Bang and Olufsen; Music for Archimedes; CD B&O 101
- [9] Crochiere, R. A weighted overlap-add method of short-time Fourier analysis/synthesis. *IEEE Trans Acoust., Speech, Signal Process.*, 28(1), 1980, pp 99–102.
- [10] Gerkmann T.; Hendriks R.; Noise power estimation based on the probability of speech presence. *Proc. 2011 IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA '11)*, New Paltz, NY, 2011.
- [11] Greenberg J; Peterson P; Zurek P. Intelligibility-weighted measures of speech-to-interference ratio and speech system performance. *J. Acoust. Soc. Amer.*, 94(5), 1993, pp 3009–3010.
- [12] Dietz M.; Ewert S.; Hohmann V. Auditory model based direction estimation of concurrent speakers from binaural signals. *Speech Communication*, 53(5), 2011, pp 592–605.
- [13] Ali R. <ftp://ftp.esat.kuleuven.be/pub/SISTA/rali/Reports/ICA2019/AudioBIN>, 2019.