

# OPTIMAL BINAURAL LCMV BEAMFORMING IN COMPLEX ACOUSTIC SCENARIOS: THEORETICAL AND PRACTICAL INSIGHTS

Nico Gößling<sup>1</sup>, Daniel Marquardt<sup>1,2</sup>, Ivo Merks<sup>2</sup>, Tao Zhang<sup>2</sup>, Simon Doclo<sup>1</sup>

<sup>1</sup>University of Oldenburg, Department of Medical Physics and Acoustics and Cluster of Excellence Hearing4All, Oldenburg, Germany

<sup>2</sup>Starkey Hearing Technologies, Eden Prairie, MN, 55344, USA

## ABSTRACT

Binaural beamforming algorithms for head-mounted assistive listening devices are crucial to improve speech quality and speech intelligibility in noisy environments, while maintaining the spatial impression of the acoustic scene. While the well-known BMVDR beamformer is able to preserve the binaural cues of one desired source, the BLCMV beamformer uses additional constraints to also preserve the binaural cues of interfering sources. In this paper, we provide theoretical and practical insights on how to optimally set the interference scaling parameters in the BLCMV beamformer for an arbitrary number of interfering sources. In addition, since in practice only a limited temporal observation interval is available to estimate all required beamformer quantities, we provide an experimental evaluation in a complex acoustic scenario using measured impulse responses from hearing aids in a cafeteria for different observation intervals. The results show that even rather short observation intervals are sufficient to achieve a decent noise reduction performance and that a proposed threshold on the optimal interference scaling parameters leads to smaller binaural cue errors in practice.

**Index Terms**— Hearing aids, binaural cues, noise reduction, beamforming, BLCMV, RTF

## 1. INTRODUCTION

For head-mounted assistive listening devices (e.g., hearing aids, cochlear implants), algorithms that use the microphone signals from both the left and the right hearing device are effective techniques to improve speech intelligibility, as the spatial information captured by all microphones can be exploited [1, 2]. Besides reducing undesired sources and limiting speech distortion, another important objective of binaural speech enhancement algorithms is the preservation of the listener's perception of the acoustical scene, in order to exploit the binaural hearing advantage [3] and to reduce confusions due to a mismatch between acoustical and visual information. To achieve binaural noise reduction with binaural cue preservation, two main concepts have been developed. In the first concept, a common real-valued spectro-temporal gain is applied to the reference microphone signals in the left and the right hearing device [4–10], ensuring perfect preservation of the instantaneous binaural cues but inevitably introducing speech distortion. The second concept, which is considered in this paper, is to apply a complex-valued filter to all available microphone signals on the left and the right hearing device using binaural extensions of spatial filtering techniques [11–19].

While the well-known binaural minimum variance distortionless response

This work was supported by the Collaborative Research Centre 1330 Hearing Acoustics and the Cluster of Excellence 1077 Hearing4all, funded by the German Research Foundation (DFG), by the joint Lower Saxony-Israeli Project ATHENA, and a research gift from Starkey Hearing Technologies.

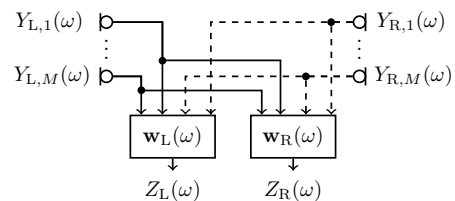


Fig. 1: Binaural hearing device configuration.

(BMVDR) beamformer [1] preserves the binaural cues (i.e., the interaural level difference (ILD) and interaural time difference (ITD)) of one desired source, the binaural linearly constrained minimum variance (BLCMV) beamformer [15] is also able to preserve the binaural cues of interfering sources. This is achievable by imposing interference scaling constraints for these sources. It should be noted that the BMVDR and BLCMV beamformers require an estimate of the correlation matrix that should be minimized and an estimate of the *relative transfer functions* (RTFs) of the desired and interfering sources. The performance of these beamformers may significantly deteriorate in case of estimation errors. Such estimation errors occur if only short temporal observation intervals for estimation can be used, e.g., due to dynamic spatial scenarios such as moving sources or head movement.

In this paper, we first derive optimal values for the interference scaling parameters in the BLCMV beamformer based on the BMVDR beamformer with RTF preservation (BMVDR-RTF) [13, 14] for an arbitrary number of interfering sources. Secondly, since these values are optimal in the sense of noise reduction but not robust against RTF estimation errors in practice, we propose to apply an upper and lower threshold on them. We evaluate the performance of the BMVDR beamformer and the BLCMV beamformer using the two different interference scaling parameters and measured impulse responses from hearing aids in a cafeteria [20] for several temporal observation intervals. The results show that even rather short temporal observation intervals lead to sufficient noise reduction performance and that the imposed threshold on the optimal interference scaling parameters can significantly reduce binaural cue errors.

## 2. CONFIGURATION AND NOTATION

Consider the binaural hearing device configuration in Fig. 1, consisting of a microphone array with  $M$  microphones on the left and the right hearing device. For an acoustic scenario with one desired source,  $P$  interfering sources and incoherent background noise, the  $m$ -th microphone signal of the left hearing device  $Y_{L,m}(\omega)$  can be written in the frequency-domain

as

$$Y_{L,m}(\omega) = X_{L,m}(\omega) + \sum_{p=1}^P U_{L,p,m}(\omega) + N_{L,m}(\omega), \quad (1)$$

with  $X_{L,m}(\omega)$  the desired speech component,  $U_{L,p,m}(\omega)$  the  $p$ -th interference component and  $N_{L,m}(\omega)$  the background noise component (e.g., diffuse noise) in the  $m$ -th microphone signal. The  $m$ -th microphone signal of the right hearing device  $Y_{R,m}(\omega)$  is defined similarly. For conciseness we will omit the frequency variable  $\omega$  in the remainder of the paper. We define the  $2M$ -dimensional stacked signal vector  $\mathbf{y}$  as

$$\mathbf{y} = [Y_{L,1} \dots Y_{L,M} Y_{R,1} \dots Y_{R,M}]^T, \quad (2)$$

where  $(\cdot)^T$  denotes the transpose and which can be written as

$$\mathbf{y} = \mathbf{x} + \mathbf{v}, \quad \mathbf{v} = \sum_{p=1}^P \mathbf{u}_p + \mathbf{n}, \quad (3)$$

where  $\mathbf{x}$ ,  $\mathbf{u}_p$  and  $\mathbf{n}$  are defined similarly as in (2) and  $\mathbf{v}$  denotes the overall undesired component, i.e., interference plus background noise components. For the coherent desired source  $S_x$  and the coherent interfering sources  $S_{u,p}$ , with  $p \in \{1, \dots, P\}$ , the vectors  $\mathbf{x}$  and  $\mathbf{u}_p$  can be written as

$$\mathbf{x} = S_x \mathbf{a}, \quad \mathbf{u}_p = S_{u,p} \mathbf{b}_p, \quad (4)$$

with  $\mathbf{a}$  and  $\mathbf{b}_p$  the *acoustic transfer functions* (ATFs) between all microphones and the desired and the  $p$ -th interfering source, respectively. Without loss of generality, we choose the first microphones on the left and the right hearing device as reference microphones, i.e.,

$$Y_L = \mathbf{e}_L^T \mathbf{y}, \quad Y_R = \mathbf{e}_R^T \mathbf{y}, \quad (5)$$

where  $\mathbf{e}_L$  and  $\mathbf{e}_R$  are  $2M$ -dimensional vectors with one element equal to 1 and the other elements equal to 0, i.e.,  $\mathbf{e}_L[1] = 1$  and  $\mathbf{e}_R[M+1] = 1$ . The correlation matrices of the background noise component, the desired speech component, the  $p$ -th interference component and all interference components are defined as

$$\mathbf{R}_n = \mathcal{E}\{\mathbf{nn}^H\}, \quad \mathbf{R}_x = \mathcal{E}\{\mathbf{xx}^H\} = \Phi_x \mathbf{a} \mathbf{a}^H, \quad (6)$$

$$\mathbf{R}_{u,p} = \mathcal{E}\{\mathbf{u}_p \mathbf{u}_p^H\} = \Phi_{u,p} \mathbf{b}_p \mathbf{b}_p^H, \quad \mathbf{R}_u = \sum_{p=1}^P \mathbf{R}_{u,p}, \quad (7)$$

where  $\mathcal{E}\{\cdot\}$  denotes the expectation operator,  $(\cdot)^H$  denotes the conjugate transpose and  $\Phi_x$  and  $\Phi_{u,p}$  denote the *power spectral density* (PSD) of the desired source and the  $p$ -th interfering source, respectively. Assuming statistical independence between the components in (1), the correlation matrix of the microphone signals  $\mathbf{R}_y$  can be written as

$$\mathbf{R}_y = \mathbf{R}_x + \mathbf{R}_u + \mathbf{R}_n = \mathbf{R}_x + \mathbf{R}_v, \quad (8)$$

with  $\mathbf{R}_v$  the correlation matrix of the overall undesired component. The output signal at the left hearing device  $Z_L$  is obtained by filtering the microphone signals with the  $2M$ -dimensional filter  $\mathbf{w}_L$ , i.e.,

$$Z_L = \mathbf{w}_L^H \mathbf{y} = \mathbf{w}_L^H \mathbf{x} + \sum_{p=1}^P \mathbf{w}_L^H \mathbf{u}_p + \mathbf{w}_L^H \mathbf{n}, \quad (9)$$

The output signal at the right hearing aid  $Z_R$  is similarly defined. Furthermore, we define the  $4M$ -dimensional filter vector  $\mathbf{w}$  as

$$\mathbf{w} = \begin{bmatrix} \mathbf{w}_L \\ \mathbf{w}_R \end{bmatrix}. \quad (10)$$

The RTF vectors of the desired and the interfering sources are defined by relating the ATF vectors to the ATF of the reference microphone on

the left and the right hearing device, i.e.,

$$\mathbf{a}_L = \frac{\mathbf{a}}{A_L}, \quad \mathbf{a}_R = \frac{\mathbf{a}}{A_R}, \quad \mathbf{b}_{L,p} = \frac{\mathbf{b}_p}{B_{L,p}}, \quad \mathbf{b}_{R,p} = \frac{\mathbf{b}_p}{B_{R,p}}. \quad (11)$$

The  $2M \times P$ -dimensional matrices  $\mathbf{B}_L$  and  $\mathbf{B}_R$  containing the RTF vectors of all interfering sources are defined as

$$\mathbf{B}_L = [\mathbf{b}_{L,1}, \dots, \mathbf{b}_{L,P}], \quad \mathbf{B}_R = [\mathbf{b}_{R,1}, \dots, \mathbf{b}_{R,P}]. \quad (12)$$

The binaural input and output *signal-to-noise ratio* (SNR) is defined as the ratio of the average input and output PSDs of the desired speech component and the background noise component, i.e.,

$$\text{SNR}^i = \frac{\mathbf{e}_L^T \mathbf{R}_x \mathbf{e}_L + \mathbf{e}_R^T \mathbf{R}_x \mathbf{e}_R}{\mathbf{e}_L^T \mathbf{R}_n \mathbf{e}_L + \mathbf{e}_R^T \mathbf{R}_n \mathbf{e}_R}, \quad \text{SNR}^o = \frac{\mathbf{w}_L^H \mathbf{R}_x \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_x \mathbf{w}_R}{\mathbf{w}_L^H \mathbf{R}_n \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_n \mathbf{w}_R}. \quad (13)$$

The binaural input and output *signal-to-interference ratio* (SIR) is defined as the ratio of the average input and output PSDs of the desired speech component and the interference components, i.e.,

$$\text{SIR}^i = \frac{\mathbf{e}_L^T \mathbf{R}_x \mathbf{e}_L + \mathbf{e}_R^T \mathbf{R}_x \mathbf{e}_R}{\mathbf{e}_L^T \mathbf{R}_u \mathbf{e}_L + \mathbf{e}_R^T \mathbf{R}_u \mathbf{e}_R}, \quad \text{SIR}^o = \frac{\mathbf{w}_L^H \mathbf{R}_x \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_x \mathbf{w}_R}{\mathbf{w}_L^H \mathbf{R}_u \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_u \mathbf{w}_R}. \quad (14)$$

The binaural input and output *signal-to-interference-plus-noise ratio* (SINR) is defined as the ratio of the average input and output PSDs of the desired speech component and the overall undesired component, i.e.,

$$\text{SINR}^i = \frac{\mathbf{e}_L^T \mathbf{R}_x \mathbf{e}_L + \mathbf{e}_R^T \mathbf{R}_x \mathbf{e}_R}{\mathbf{e}_L^T \mathbf{R}_v \mathbf{e}_L + \mathbf{e}_R^T \mathbf{R}_v \mathbf{e}_R}, \quad \text{SINR}^o = \frac{\mathbf{w}_L^H \mathbf{R}_x \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_x \mathbf{w}_R}{\mathbf{w}_L^H \mathbf{R}_v \mathbf{w}_L + \mathbf{w}_R^H \mathbf{R}_v \mathbf{w}_R}. \quad (15)$$

### 3. BINAURAL NOISE REDUCTION ALGORITHMS

In Section 3.1 and 3.2 we briefly review the BMVDR beamformer [1, 2, 12] and the BLCMV beamformer [15]. Based on the optimality of the BMVDR-RTF beamformer [14] in optimizing the SINR (or SNR) while preserving the binaural cues of all sources, in Section 3.3 we derive optimal values for the interference scaling parameters in the BLCMV beamformer in the case of an arbitrary number of interfering sources. Furthermore, in order to achieve a robust binaural cue preservation performance in case of estimation errors of the correlation matrices and the RTF vectors (Section 3.4), we propose to threshold these interference scaling parameters.

#### 3.1. BMVDR beamformer

The BMVDR beamformer aims at minimizing the output PSD in both hearing devices, while preserving the desired speech component in the reference microphone signals. The corresponding constrained optimization problem is given by

$$\min_{\mathbf{w}} \mathbf{w}^H \tilde{\mathbf{R}} \mathbf{w} \quad \text{subject to} \quad \mathbf{w}^H \mathbf{C} = \mathbf{g}, \quad (16)$$

with

$$\tilde{\mathbf{R}} = \begin{bmatrix} \mathbf{R} & \mathbf{0}_{2M \times 2M} \\ \mathbf{0}_{2M \times 2M} & \mathbf{R} \end{bmatrix}, \quad (17)$$

with  $\mathbf{R}$  either equal to the correlation matrix  $\mathbf{R}_y$  of the microphone signals, the correlation matrix  $\mathbf{R}_v$  of the overall undesired component or the correlation matrix  $\mathbf{R}_n$  of the background noise component. The constraint set in (16) is given by

$$\mathbf{C} = \begin{bmatrix} \mathbf{a}_L & \mathbf{0}_{2M \times 1} \\ \mathbf{0}_{2M \times 1} & \mathbf{a}_R \end{bmatrix}, \quad \mathbf{g} = [1 \quad 1], \quad (18)$$

requiring the RTF vectors of the desired source. The solution to the optimization problem in (16) using the constraint set in (18) is equal to [1, 12, 21]

$$\mathbf{w}_{\text{MVDR,L}} = \frac{\mathbf{R}^{-1} \mathbf{a}_L}{\mathbf{a}_L^H \mathbf{R}^{-1} \mathbf{a}_L}, \quad \mathbf{w}_{\text{MVDR,R}} = \frac{\mathbf{R}^{-1} \mathbf{a}_R}{\mathbf{a}_R^H \mathbf{R}^{-1} \mathbf{a}_R}. \quad (19)$$

From a theoretical point of view, in the case of perfectly estimated quantities (i.e., correlation matrices and RTF vector), using  $\mathbf{R} = \mathbf{R}_y$  or  $\mathbf{R} = \mathbf{R}_v$  in (19) is optimal in the SINR sense, whereas using  $\mathbf{R} = \mathbf{R}_n$  in (19) is optimal in the SNR sense. While the BMVDR beamformer preserves the binaural cues of the desired source, its major drawback is the distortion of the binaural cues of the interfering sources (and background noise), such that all sources are perceived as coming from the direction of the desired source. In practice, it should also be realized that using  $\mathbf{R} = \mathbf{R}_y$  may lead to target cancellation in the case of RTF estimation errors of the desired source [21] and that  $\mathbf{R}_v$  is not straightforward to estimate.

### 3.2. BLCMV beamformer

In order to also take binaural cue preservation of the interfering sources into account as well as control the amount of interference suppression, it has been proposed in [15] to add interference scaling constraints to the BMVDR beamformer, leading to the BLCMV beamformer. This corresponds to the constrained optimization problem in (18) with the constraint set

$$\mathbf{C}_1 = \begin{bmatrix} \mathbf{a}_L & \mathbf{B}_L & \mathbf{0}_{2M \times 1} & \mathbf{0}_{2M \times P} \\ \mathbf{0}_{2M \times 1} & \mathbf{0}_{2M \times P} & \mathbf{a}_R & \mathbf{B}_R \end{bmatrix}, \quad \mathbf{g}_1 = [1 \ \delta_L \ 1 \ \delta_R], \quad (20)$$

requiring the RTF vectors of the desired source and all interfering sources. The  $P$ -dimensional vectors  $\delta_L = [\delta_{L,1} \dots \delta_{L,P}]$  and  $\delta_R = [\delta_{R,1} \dots \delta_{R,P}]$  contain the *interference scaling parameters*, which control the suppression and the binaural cue preservation of the  $P$  interfering sources. The BLCMV beamformer is given by

$$\mathbf{w}_{\text{LCMV}} = \tilde{\mathbf{R}}^{-1} \mathbf{C}_1 \left( \mathbf{C}_1^H \tilde{\mathbf{R}}^{-1} \mathbf{C}_1 \right)^{-1} \mathbf{g}_1^H. \quad (21)$$

Setting  $\delta_{L,p} = \delta_{R,p}$  ensures binaural cue preservation of the  $p$ -th interfering source, while the absolute values of  $\delta_{L,p}$  and  $\delta_{R,p}$  directly determine the SIR improvement for the  $p$ -th interfering source. From a theoretical point of view, in the case of perfectly estimated quantities (i.e., correlation matrices and RTF vectors), setting  $\delta_{L,p} = \delta_{R,p} = 0$  in the BLCMV beamformer is optimal in the SIR sense, but not necessarily in the SINR or SNR sense. Moreover, in contrast to the BMVDR beamformer, the choice of the correlation matrix  $\mathbf{R}$  has no impact on the SINR, SNR and SIR improvement and the binaural cue preservation as these are completely determined by the interference scaling parameters. In practice, in the case of estimation errors the choice of the correlation matrix  $\mathbf{R}$  will obviously have an influence on the performance of the BLCMV beamformer (cf. Section 4).

### 3.3. Interference scaling parameters

As an extension of the method presented in [22] for an arbitrary number of interfering sources, in this section we propose a method to determine the interference scaling parameters that maximize the SINR or the SNR while preserving the binaural cues of the interfering sources. To this end, we will use the BMVDR beamformer with RTF preservation [14], denoted as BMVDR-RTF beamformer, which is a special case of the BLCMV beamformer. In the BMVDR-RTF beamformer the constraints related to the interfering sources only control the binaural cue preservation

while the amount of desired interference suppression is not specified, i.e.,

$$\frac{\mathbf{w}_L^H \mathbf{b}_p}{\mathbf{w}_R^H \mathbf{b}_p} = \frac{B_{L,p}}{B_{R,p}} \Rightarrow \frac{\mathbf{w}_L^H \mathbf{b}_{L,p}}{\mathbf{w}_R^H \mathbf{b}_{R,p}} = 1, \quad (22)$$

leading to the constraint set

$$\mathbf{C}_2 = \begin{bmatrix} \mathbf{a}_L & \mathbf{B}_L & \mathbf{0}_{2M \times 1} \\ \mathbf{0}_{2M \times 1} & -\mathbf{B}_R & \mathbf{a}_R \end{bmatrix}, \quad \mathbf{g}_2 = [1 \ \mathbf{0}_{1 \times P} \ 1]. \quad (23)$$

The BMVDR-RTF beamformer is given by [14]

$$\mathbf{w}_{\text{RTF}} = \tilde{\mathbf{R}}^{-1} \mathbf{C}_2 \left( \mathbf{C}_2^H \tilde{\mathbf{R}}^{-1} \mathbf{C}_2 \right)^{-1} \mathbf{g}_2^H, \quad (24)$$

and either maximizes the SINR ( $\mathbf{R} = \mathbf{R}_y$  or  $\mathbf{R} = \mathbf{R}_v$ ) or the SNR ( $\mathbf{R} = \mathbf{R}_n$ ), while preserving the binaural cues of all sources.

Hence, the optimal interference scaling parameters for the BLCMV beamformer (in the SINR or SNR sense) can be determined as

$$\delta_p^{\text{opt}} = \delta_{L,p} = \delta_{R,p} = \mathbf{w}_{\text{RTF,L}}^H \mathbf{b}_{L,p} = \mathbf{w}_{\text{RTF,R}}^H \mathbf{b}_{R,p} \quad (25)$$

However, using the optimal interference scaling parameters may lead to problems in practice due to estimation errors of the correlation matrices and RTF vectors. More in particular, in the case of SINR maximization, the corresponding interference scaling parameters may be rather small, leading to a decreased binaural cue preservation performance (cf. simulations in Section 4). On the other hand, in the case of SNR maximization, the corresponding interference scaling parameters may be rather large, depending on the position of the interfering source, leading to an unsatisfying SINR improvement. Hence, we propose to enforce an upper and lower threshold on the optimal interference scaling parameters, i.e.,

$$\delta_p^{\text{thr}} = \begin{cases} |\delta_p^{\text{opt}}|, & \text{if } \delta^{\min} < |\delta_p^{\text{opt}}| < \delta^{\max}, \\ \delta^{\min}, & \text{if } |\delta_p^{\text{opt}}| \leq \delta^{\min}, \\ \delta^{\max}, & \text{if } |\delta_p^{\text{opt}}| \geq \delta^{\max}. \end{cases} \quad (26)$$

The thresholds have been experimentally obtained as  $\delta^{\min} = 0.2$  and  $\delta^{\max} = 0.4$ , limiting the theoretically possible SIR improvement for each interfering source between 8 dB and 14 dB.

### 3.4. Estimation of correlation matrices and RTFs

All considered binaural beamformers require an estimate of the RTF vectors  $\mathbf{a}_L$  and  $\mathbf{a}_R$  of the desired source (cf. (11)). In addition, the BLCMV and BMVDR-RTF beamformers require an estimate of the RTF vectors  $\mathbf{b}_{L,p}$  and  $\mathbf{b}_{R,p}$  of each interfering source. In this paper, we will estimate these RTFs using the covariance whitening approach [23, 24], which is based on the generalized eigenvalue decomposition (GEVD) of the speech + noise correlation matrix  $\mathbf{R}_{\text{xn}} = \mathbf{R}_x + \mathbf{R}_n$  and the background noise correlation matrix  $\mathbf{R}_n$  or the GEVD of the interference + noise correlation matrix  $\mathbf{R}_{v,p} = \mathbf{R}_{u,p} + \mathbf{R}_n$  and  $\mathbf{R}_n$ . While  $\mathbf{R}_n$  can be estimated exploiting the assumed stationarity of the background noise, estimating  $\mathbf{R}_{\text{xn}}$  and  $\mathbf{R}_{v,p}$  from the available mixture is not straightforward. Due to limited source activity and possible spatial changes of the acoustic scenarios, the *temporal observation interval* that is available in practice for estimating these correlation matrices is typically limited. We assume that the correlation matrix  $\mathbf{R}_{\text{xn}}$  can be estimated from an observation interval consisting of  $T_L$  frames (corresponding to  $L$  seconds) where only the desired source and the background noise are active, i.e.,

$$\hat{\mathbf{R}}_{\text{xn}} = \frac{1}{T_L} \sum_{t=1}^{T_L} \left( \mathbf{x}(t) + \mathbf{n}(t) \right) \left( \mathbf{x}(t) + \mathbf{n}(t) \right)^H, \quad (27)$$

where  $t$  is the frame index. Similarly, we assume that the correlation matrix  $\mathbf{R}_{v,p}$  can be estimated from an observation interval of  $T_L$  frames where only the  $p$ -th interfering source and the background noise are active.

#### 4. EXPERIMENTAL RESULTS

In this section, we experimentally investigate the effect of the temporal observation interval on the performance of the BMVDR beamformer ( $\mathbf{w}_{\text{MVDR}}$ ) and the BLCMV beamformer using either the optimal interference scaling parameters ( $\mathbf{w}_{\text{LCMV}}(\delta^{\text{opt}})$ ) or the proposed thresholded interference scaling parameters ( $\mathbf{w}_{\text{LCMV}}(\delta^{\text{thr}})$ ) (cf. Section 3.3).

We consider three different acoustic scenarios comprising of one desired source, one or two interfering sources and diffuse background noise (cf. Table 1 for source positions). The desired source was a male German speaker, the first interfering source was a male Dutch speaker and the second interfering source was a male English speaker. The desired speech and interference components were generated by convolving the desired and interfering source signals with measured impulse responses of binaural behind-the-ear hearing aids mounted on a dummy head in a cafeteria ( $T_{60} \approx 1250\text{ms}$ ) [20], with  $M = 2$  microphones per hearing aid. For background noise we used real ambient noise recorded in the same cafeteria with the same setup. The sampling frequency was 16 kHz. All signals start with 2 s of noise-only, followed by about 20 s of all sources being active. The broadband input SNR was set to 5 dB and the SIRs were set to 0 dB.

The noise correlation matrix  $\mathbf{R}_n$  was estimated using the 2 s noise-only segment. To estimate the correlation matrices  $\mathbf{R}_y$ ,  $\mathbf{R}_v$ ,  $\mathbf{R}_{\text{xin}}$  and  $\mathbf{R}_{v,p}$ , we considered different temporal observation intervals (starting at 2 s), whose length  $L$  ranged between 0.1 s and 3 s. To estimate the correlation matrices  $\mathbf{R}_v$ ,  $\mathbf{R}_{\text{xin}}$  and  $\mathbf{R}_{v,p}$  the algorithm had access to the respective mixtures. The RTF vectors of the desired source and the interfering source(s) were then calculated based on these estimated correlation matrices (cf. Section 3.4). Please note that shorter temporal observation intervals correspond to larger estimation errors.

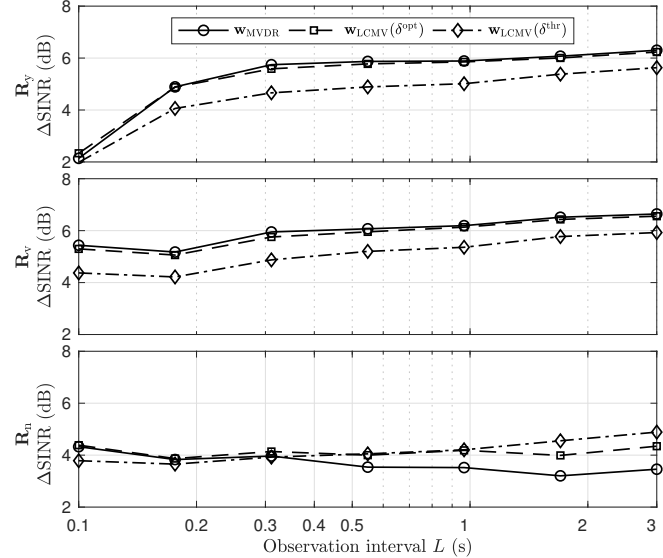
The microphone signals were processed using a weighted overlap-add framework with a block length of 256 with 50% overlap and a square-root Hann window. The BMVDR and BLCMV beamformers were calculated using three different correlation matrices, i.e.,  $\mathbf{R} = \mathbf{R}_y$  (maximizing SINR with possible target cancellation),  $\mathbf{R} = \mathbf{R}_v$  (maximizing SINR) and  $\mathbf{R} = \mathbf{R}_n$  (maximizing SNR). The filters were used as fixed filters over the whole signal.

As performance measures we used the binaural SINR improvement and the binaural cues errors, i.e., ILD and ITD errors, that we calculated using a model of binaural auditory processing [25]. All performance measures were averaged over all frequencies and all acoustic scenarios.

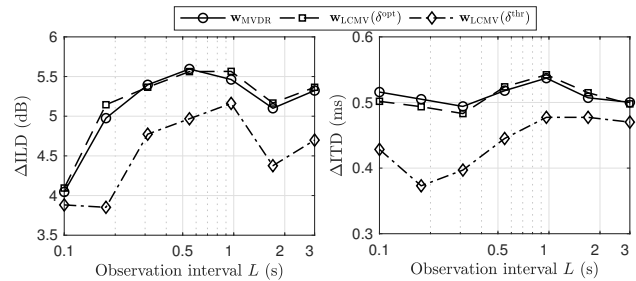
Figure 2 depicts the SINR improvement for different lengths of the temporal observation interval and for different correlation matrices, while Figure 3 depicts the binaural cue errors of the first interfering source for the same temporal observation intervals and  $\mathbf{R} = \mathbf{R}_v$ . First, it can be observed that when using  $\mathbf{R}_y$  or  $\mathbf{R}_v$  the SINR improvement is generally larger than when using  $\mathbf{R}_n$ . This is expected because using the noise correlation matrix  $\mathbf{R}_n$  is maximizing the SNR and not the SINR. Second, when using  $\mathbf{R}_y$  or  $\mathbf{R}_v$ , an apparent difference can be seen for small observation intervals below 200 ms. The small observation intervals lead to larger estimation errors for the correlation matrices and hence also for the RTF vectors, such that the drop in SINR improvement observed when using  $\mathbf{R}_y$  is probably attributed to target cancellation. For longer observation intervals and hence smaller estimation errors, the difference between using  $\mathbf{R}_y$  and  $\mathbf{R}_v$  is smaller. As expected, the SINR improvement of the BLCMV beamformer using the thresholded interference scaling parameters  $\delta^{\text{thr}}$  is smaller than for the BLCMV beamformer using the optimal interference scaling parameters  $\delta^{\text{opt}}$ . Although, looking at the binaural cue errors, using  $\delta^{\text{thr}}$  in the BLCMV beamformer leads to much better binaural cue preservation, while using  $\delta^{\text{opt}}$  leads to similar binaural cue errors as for the BMVDR beamformer. This difference is especially visible for the ITD error at small observation intervals and is also confirmed by

Scenario	1	2	3
Desired	$-35^\circ$	$0^\circ$	$0^\circ$
Interfering	$150^\circ$	$-35^\circ$	$-35^\circ, 150^\circ$

**Table 1:** Spatial scenarios ( $0^\circ$ : frontal direction.  $-90^\circ$ : left hand side.  $90^\circ$ : right hand side).



**Fig. 2:** SINR improvement for different temporal observation intervals for  $\mathbf{R} = \mathbf{R}_y$  (top),  $\mathbf{R} = \mathbf{R}_v$  (mid) and  $\mathbf{R} = \mathbf{R}_n$  (bottom).



**Fig. 3:** Binaural cue errors of the first interfering source ( $\mathbf{R} = \mathbf{R}_v$ ).

informal listening tests. Third, when using  $\mathbf{R}_n$ , the BLCMV beamformer outperforms the BMVDR beamformer for longer observation intervals above 300ms because of the additional constraints. Additionally, using  $\delta^{\text{thr}}$  in the BLCMV beamformer apparently leads to marginally better SINR improvement in this case. Because  $\mathbf{R}_v$  is in practice very hard to accurately estimate, it should be recommended to use  $\mathbf{R}_n$  when short observation intervals are required (e.g., in dynamic acoustic scenarios) and to use  $\delta^{\text{thr}}$  in the BLCMV beamformer to prevent binaural cue errors.

#### 5. CONCLUSIONS

In this paper, we proposed optimal values for the interference scaling parameters in the BLCMV beamformer for an arbitrary number of interfering sources based on the BMVDR-RTF beamformer. We showed how to set these parameters in practice such that a robust performance in the case of estimation errors can be achieved. Evaluation results in a complex acoustic scenario showed that even short temporal observation intervals for estimating the required correlation matrices and RTF vectors are sufficient to achieve a decent noise reduction performance and binaural cue preservation.



## 6. REFERENCES

- [1] S. Doclo, W. Kellermann, S. Makino, and S.E. Nordholm, "Multichannel Signal Enhancement Algorithms for Assisted Listening Devices: Exploiting spatial diversity using multiple microphones," *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 18–30, Mar. 2015.
- [2] S. Doclo, S. Gannot, D. Marquardt, and E. Hadad, "Binaural Speech Processing with Application to Hearing Devices," in *Audio Source Separation and Speech Enhancement*, chapter 18. Wiley, 2018.
- [3] A. W. Bronkhorst and R. Plomp, "The effect of head-induced interaural time and level differences on speech intelligibility in noise," *The Journal of the Acoustical Society of America*, vol. 83, no. 4, pp. 1508–1516, 1988.
- [4] T. Lotter and P. Vary, "Dual-channel speech enhancement by superdirective beamforming," *EURASIP Journal on Applied Signal Processing*, vol. 2006, pp. 1–14, 2006.
- [5] G. Grimm, V. Hohmann, and B. Kollmeier, "Increase and Subjective Evaluation of Feedback Stability in Hearing Aids by a Binaural Coherence-based Noise Reduction Scheme," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 7, pp. 1408–1419, Sep. 2009.
- [6] A. H. Kamkar-Parsi and M. Bouchard, "Improved Noise Power Spectrum Density Estimation for Binaural Hearing Aids Operating in a Diffuse Noise Field Environment," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 4, pp. 521–533, May 2009.
- [7] A. H. Kamkar-Parsi and M. Bouchard, "Instantaneous Binaural Target PSD Estimation for Hearing Aid Noise Reduction in Complex Acoustic Environments," *IEEE Transactions on Instrumentation and Measurement*, vol. 60, no. 4, pp. 1141–1154, Apr. 2011.
- [8] K. Reindl, Y. Zheng, A. Schwarz, S. Meier, R. Maas, A. Sehr, and W. Kellermann, "A stereophonic acoustic signal extraction scheme for noisy and reverberant environments," *Computer Speech and Language*, vol. 27, no. 3, pp. 726–745, 2013.
- [9] R. Baumgärtel, M. Krawczyk-Becker, D. Marquardt, C. Völker, H. Hu, T. Herzke, G. Coleman, K. Adiloğlu, S. M. A. Ernst, T. Gerkmann, S. Doclo, B. Kollmeier, V. Hohmann, and M. Dietz, "Comparing binaural signal processing strategies I: Instrumental evaluation," *Trends in Hearing*, vol. 19, pp. 1–16, 2015.
- [10] D. Marquardt and S. Doclo, "Noise power spectral density estimation for binaural noise reduction exploiting direction of arrival estimates," in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, New Paltz NY, USA, Oct. 2017, pp. 234–238.
- [11] R. Aichner, H. Buchner, M. Zourub, and W. Kellermann, "Multi-channel source separation preserving spatial information," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Honolulu HI, USA, Apr. 2007, pp. 5–8.
- [12] B. Cornelis, S. Doclo, T. Van den Bogaert, J. Wouters, and M. Moonen, "Theoretical analysis of binaural multi-microphone noise reduction techniques," *IEEE Transactions on Audio, Speech and Language Processing*, vol. 18, no. 2, pp. 342–355, Feb. 2010.
- [13] D. Marquardt, E. Hadad, S. Gannot, and S. Doclo, "Theoretical Analysis of Linearly Constrained Multi-channel Wiener Filtering Algorithms for Combined Noise Reduction and Binaural Cue Preservation in Binaural Hearing Aids," *IEEE/ACM Trans. on Audio, Speech, and Language Processing*, vol. 23, no. 12, pp. 2384–2397, Dec. 2015.
- [14] E. Hadad, D. Marquardt, S. Doclo, and S. Gannot, "Theoretical Analysis of Binaural Transfer Function MVDR Beamformers with Interference Cue Preservation Constraints," *IEEE/ACM Trans. Audio, Speech and Language Proc.*, vol. 23, no. 12, pp. 2449–2464, Dec. 2015.
- [15] E. Hadad, S. Doclo, and S. Gannot, "The Binaural LCMV Beamformer and its Performance Analysis," *IEEE/ACM Trans. on Audio, Speech, and Language Proc.*, vol. 24, no. 3, pp. 543–558, 2016.
- [16] E. Hadad, D. Marquardt, W. Pu, S. Gannot, S. Doclo, Z.-Q. Luo, I. Merks, and T. Zhang, "Comparison of two binaural beamforming approaches for hearing aids," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, New Orleans, USA, Mar. 2017, pp. 236–240.
- [17] A. I. Koutrouvelis, R. C. Hendriks, R. Heusdens, and J. Jensen, "Relaxed binaural LCMV beamforming," *IEEE/ACM Trans. on Audio, Speech and Language Processing*, vol. 25, no. 1, pp. 137–152, Jan. 2017.
- [18] W. Pu, J. Xiao, T. Zhang, and Z.-Q. Luo, "A penalized inequality-constrained minimum variance beamformer with applications in hearing aids," in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, New Paltz NY, USA, Oct. 2017, pp. 175–179.
- [19] D. Marquardt and S. Doclo, "Interaural Coherence Preservation for Binaural Noise Reduction Using Partial Noise Estimation and Spectral Postfiltering," *IEEE/ACM Trans. on Audio, Speech and Language Processing*, vol. 26, no. 7, pp. 1257–1270, 2018.
- [20] H. Kayser, S. Ewert, J. Annemüller, T. Rohdenburg, V. Hohmann, and B. Kollmeier, "Database of multichannel In-Ear and Behind-The-Ear Head-Related and Binaural Room Impulse Responses," *Eurasip Journal on Advances in Signal Processing*, vol. 2009, pp. 10 pages, 2009.
- [21] B. D. Van Veen and K. M. Buckley, "Beamforming: a versatile approach to spatial filtering," *IEEE ASSP Magazine*, vol. 5, no. 2, pp. 4–24, Apr. 1988.
- [22] D. Marquardt, E. Hadad, S. Gannot, and S. Doclo, "Optimal binaural LCMV beamformers for combined noise reduction and binaural cue preservation," in *Proc. International Workshop on Acoustic Signal Enhancement (IWAENC)*, Juan-les-Pins, France, Sep. 2014, pp. 288–292.
- [23] S. Markovich, S. Gannot, and I. Cohen, "Multichannel eigenspace beamforming in a reverberant noisy environment with multiple interfering speech signals," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 6, pp. 1071–1086, Aug 2009.
- [24] S. Markovich-Golan and S. Gannot, "Performance analysis of the covariance subtraction method for relative transfer function estimation and comparison to the covariance whitening method," in *Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, Brisbane, Australia, Apr. 2015, pp. 544–548.
- [25] M. Dietz, S. D. Ewert, and V. Hohmann, "Auditory model based direction estimation of concurrent speakers from binaural signals," *Speech Communication*, vol. 53, pp. 592–605, 2011.