

RESEARCH

Open Access



An integrated MVDR beamformer for speech enhancement using a local microphone array and external microphones

Randall Ali* , Toon van Waterschoot and Marc Moonen

Abstract

An integrated version of the minimum variance distortionless response (MVDR) beamformer for speech enhancement using a microphone array has been recently developed, which merges the benefits of imposing constraints defined from both a relative transfer function (RTF) vector based on a priori knowledge and an RTF vector based on a data-dependent estimate. In this paper, the integrated MVDR beamformer is extended for use with a microphone configuration where a microphone array, local to a speech processing device, has access to the signals from multiple external microphones (XMs) randomly located in the acoustic environment. The integrated MVDR beamformer is reformulated as a quadratically constrained quadratic program (QCQP) with two constraints, one of which is related to the maximum tolerable speech distortion for the imposition of the a priori RTF vector and the other related to the maximum tolerable speech distortion for the imposition of the data-dependent RTF vector. An analysis of how these maximum tolerable speech distortions affect the behaviour of the QCQP is presented, followed by the discussion of a general tuning framework. The integrated MVDR beamformer is then evaluated with audio recordings from behind-the-ear hearing aid microphones and three XMs for a single desired speech source in a noisy environment. In comparison to relying solely on an a priori RTF vector or a data-dependent RTF vector, the results demonstrate that the integrated MVDR beamformer can be tuned to yield different enhanced speech signals, which may be more suitable for improving speech intelligibility despite changes in the desired speech source position and imperfectly estimated spatial correlation matrices.

Keywords: Speech enhancement, Beamforming, Minimum variance distortionless response (MVDR) beamformer, External microphones

1 Introduction

Speech processing devices such as a hearing aid, a cochlear implant, or a mobile telephone are commonly equipped with an array of microphones to capture the acoustic environment. The received microphone signals are often a mixture of a desired speech signal plus some undesired noise (any combination of interfering speakers, background noises, and reverberation). As the quality

and intelligibility of the desired speech signal is susceptible to considerable degradation in the presence of such noise, the task of suppressing this noise and extracting the desired speech signal, known as speech enhancement, is of critical importance and has been the subject of extensive research [1–3].

While successful speech enhancement strategies have been developed with microphone arrays, in some applications, due to physical space constraints, the spatial variation between the observed microphone signals may not be sufficient to yield an acceptable degree of speech enhancement. Consequently, the potential of using more ad hoc microphone configurations consisting of randomly

*Correspondence: randall.ali@esat.kuleuven.be

KU Leuven, Department of Electrical Engineering (ESAT), STADIUS Center for Dynamical Systems, Signal Processing, and Data Analytics, Kasteelpark Arenberg 10, 3001 Leuven, Belgium

placed microphones to increase the spatial sampling of the acoustic environment has developed interest [4–12]. In this paper, a specific ad hoc microphone configuration is considered, where a microphone array located on some speech processing device, hereafter referred to as a local microphone array (LMA), is linked with multiple remote or external microphones (XMs) in a centralised processing framework, i.e. all microphone signals are sent to a fusion centre for processing. The terminology of a *local* microphone array is introduced since the microphone array is considered to be confined or fixed within some area of the acoustic environment relative to the XMs which are subject to movement.

When there is a single desired speech source, speech enhancement can be accomplished by using the minimum variance distortionless response (MVDR) beamformer [13, 14]. One of the important quantities required for computing the MVDR beamformer is a vector of acoustic transfer functions from the desired speech source to all of the microphones. More commonly, however, a vector of relative transfer functions (RTFs) is used instead, which is a normalised version of the acoustic transfer function vector with respect to some reference microphone [15]. In practice, for an LMA, this RTF vector may be measured a priori or based on assumptions regarding microphone characteristics, position, speaker location, and room acoustics (e.g. no reverberation). For instance, in assistive hearing devices, it is sometimes assumed that the desired speech source location is known and this knowledge can be subsequently used to define an a priori RTF vector [16–19]. Alternatively, it may be estimated in an online fashion from the observed microphone data [20, 21] so that it is a fully data-dependent estimate.

The situation under consideration throughout this paper is one in which there is an available a priori RTF vector pertaining only to the LMA that may or may not be sufficiently accurate with respect to the true RTF vector. In cases where the a priori RTF vector is not sufficiently accurate, then incorporating the use of a data-dependent RTF vector can be viewed as an opportunity for an improved performance provided that the data-dependent RTF vector is a better estimate of the true RTF vector. On the other hand, when acoustic conditions are adverse enough to significantly affect the accuracy of the data-dependent RTF vector, then relying on the a priori RTF vector can be viewed as a fall back or contingency strategy.

It would therefore be seemingly advantageous to use both an a priori and a data-dependent RTF vector in practice. Such an approach has recently been investigated for an LMA only and resulted in an integrated version of the MVDR beamformer [22]. As opposed to imposing either the a priori RTF vector or the data-dependent RTF vector as a hard constraint, they were both softened into an unconstrained optimisation problem. It was

demonstrated that the resulting integrated MVDR beamformer is a convex combination of an MVDR beamformer that uses the a priori RTF vector, an MVDR beamformer that uses the data-dependent RTF vector, a linearly constrained minimum variance (LCMV) beamformer that uses both the a priori and data-dependent RTF vector, and an all-zero vector, each with real-valued weightings, revealing the versatile nature of such an integrated beamformer.

This paper therefore re-examines the integrated MVDR beamformer for the ad hoc microphone configuration consisting of an LMA located on some speech processing device linked with multiple XMs. Specifically, the integrated MVDR beamformer is reformulated from an alternative perspective, namely that of a quadratically constrained quadratic program (QCQP). This QCQP will consist of two constraints, one of which is related to the maximum tolerable speech distortion for the imposition of the a priori RTF vector and the other related to the maximum tolerable speech distortion for the imposition of the data-dependent RTF vector. With respect to the procedures for obtaining the RTF vectors, it is straightforward to obtain a data-dependent RTF vector; however, the notion of an a priori RTF vector when XMs are used with an LMA is a bit more ambiguous. In particular, since only partial a priori knowledge is usually available for the part of the RTF vector pertaining to the LMA, the other part pertaining to the XMs will have to be a data-dependent estimate and hence a procedure based on partial a priori knowledge [9] would be necessary. As a result, an integrated MVDR beamformer for a microphone configuration with an LMA and XMs will merge an a priori RTF vector that is based on partial a priori knowledge and a fully data-dependent one.

With the a priori and the data-dependent RTF vector for the LMA and XMs estimated, it will become evident that the optimal filter from the integrated MVDR beamformer, formulated as a QCQP, is identical to that which was derived from [22], where the Lagrangian multipliers associated with the QCQP are equivalent to the tuning parameters that have been considered in [22]. The additional insight of the QCQP formulation is that these tuning parameters or Lagrangian multipliers can be related to a maximum tolerable speech distortion for the imposition of the a priori or the data-dependent RTF vector. An analysis of this relationship is provided, which facilitates the tuning of the integrated MVDR beamformer from the more intuitive perspective of the maximum tolerable speech distortions as opposed to the combination of filters as in [22]. A general tuning framework will then be discussed along with the suggestion of some particular tuning strategies.

The integrated MVDR beamformer is then evaluated with audio recordings from behind-the-ear hearing aid

microphones (the LMA) and three XMs for a single desired speech source in a re-created cocktail party scenario. The results demonstrate that the integrated MVDR beamformer can be tuned to yield different enhanced speech signals, which can find a compromise between relying solely on an a priori RTF vector or a data-dependent RTF vector, and hence may be more suitable for improving speech intelligibility despite changes in the desired speech source position and imperfectly estimated spatial correlation matrices.

The paper is organised as follows. In Section 2, the data model is defined. In Section 3, the MVDR beamformer as applied to an LMA with XMs is discussed along with the procedures for obtaining the a priori RTF vector based on partial a priori knowledge and the data-dependent RTF vector. Section 4 reformulates the integrated MVDR beamformer as a QCQP and provides an analysis on the effect of the maximum tolerable speech distortions due to the imposition of the a priori RTF vector and the data-dependent RTF vector. In Section 5, a general tuning framework is presented, as well as some suggested tuning strategies. In Section 6, the integrated MVDR approach is analysed and evaluated with both simulated data as well as experimental data involving the use of behind-the-ear hearing aid microphones and three XMs. Conclusions are then drawn in Section 7.

2 Data model

2.1 Unprocessed signals

A microphone configuration consisting of an LMA of M_a microphones plus M_e XMs is considered with one desired speech source in a noisy, reverberant¹ environment. In the short-time Fourier transform (STFT) domain, the observed vector of microphone signals at frequency bin k and time frame l is represented as:

$$\mathbf{y}(k, l) = \underbrace{\mathbf{h}(k, l) s_{a,1}(k, l)}_{\mathbf{x}(k, l)} + \mathbf{n}(k, l) \quad (1)$$

where (dropping the dependency on k and l for brevity) $\mathbf{y} = [\mathbf{y}_a^T \mathbf{y}_e^T]^T$, $\mathbf{y}_a = [y_{a,1} \ y_{a,2} \ \dots \ y_{a,M_a}]^T$ is a vector of the LMA signals, $\mathbf{y}_e = [y_{e,1} \ y_{e,2} \ \dots \ y_{e,M_e}]^T$ is a vector of the XM signals, \mathbf{x} is the speech contribution, represented by $s_{a,1}$, the desired speech signal in the first (reference) microphone of the LMA, filtered with $\mathbf{h} = [\mathbf{h}_a^T \ \mathbf{h}_e^T]^T$, \mathbf{h}_a is the RTF vector for the LMA (where the first component of \mathbf{h}_a is equal to 1 since the first microphone is used as the reference), and \mathbf{h}_e is the RTF vector for the XM signals. Finally, $\mathbf{n} = [\mathbf{n}_a^T \ \mathbf{n}_e^T]^T$ represents the noise contribution. Variables with the subscript “a” refer to the LMA and variables with the subscript “e” refer to the XMs.

¹Reverberation is not explicitly included in the signal model as dereverberation is not addressed in this paper. This paper primarily focuses on noise reduction, although some dereverberation will be achieved as a fortunate by-product of beamforming.

The $(M_a + M_e) \times (M_a + M_e)$ spatial correlation matrix for the speech-plus-noise, noise-only, and speech-only signals is defined respectively as:

$$\mathbf{R}_{yy} = \mathbb{E} \{ \mathbf{y} \mathbf{y}^H \} \quad (2)$$

$$\mathbf{R}_{nn} = \mathbb{E} \{ \mathbf{n} \mathbf{n}^H \} \quad (3)$$

$$\mathbf{R}_{xx} = \mathbb{E} \{ \mathbf{x} \mathbf{x}^H \} \quad (4)$$

where $\mathbb{E}\{\cdot\}$ is the expectation operator and $\{\cdot\}^H$ is the Hermitian transpose. With the assumption of a single desired speech source from (1), \mathbf{R}_{xx} can be represented as a rank-1 correlation matrix as follows:

$$\mathbf{R}_{xx} = \sigma_{s_{a,1}}^2 \mathbf{h} \mathbf{h}^H \quad (5)$$

where $\sigma_{s_{a,1}}^2 = \mathbb{E} \{ |s_{a,1}|^2 \}$ is the desired speech power spectral density in the first microphone of the LMA. It is further assumed that the desired speech signal is uncorrelated with the noise signal, and hence $\mathbf{R}_{yy} = \mathbf{R}_{xx} + \mathbf{R}_{nn}$. The speech-plus-noise, noise-only, and speech-only correlation matrix can also be defined solely for the LMA signals respectively as $\mathbf{R}_{y_a y_a} = \mathbb{E} \{ \mathbf{y}_a \mathbf{y}_a^H \}$, $\mathbf{R}_{n_a n_a} = \mathbb{E} \{ \mathbf{n}_a \mathbf{n}_a^H \}$, and $\mathbf{R}_{x_a x_a} = \mathbb{E} \{ \mathbf{x}_a \mathbf{x}_a^H \}$, with $\mathbf{R}_{x_a x_a}$ also having the same rank-1 structure as in (5). It is assumed that all signal correlations can be estimated as if all signals were available in a centralised processor, i.e. a perfect communication link is assumed between the LMA and the XMs with no bandwidth constraints as well as synchronous sampling rates.

The estimate of the desired speech signal in the first microphone of the LMA, z_1 , is then obtained through a linear filtering of the microphone signals, such that:

$$z_1 = \mathbf{w}^H \mathbf{y} \quad (6)$$

where $\mathbf{w} = [\mathbf{w}_a^T \ \mathbf{w}_e^T]^T$ is a complex-valued filter.

2.2 Pre-whitened-transformed domain

As a pre-processing stage, the unprocessed microphone signals can be firstly transformed with the available a priori RTF vector for the LMA signals and then spatially pre-whitened using the resulting transformed noise-only correlation matrix, yielding a vector of pre-whitened-transformed (PWT) microphone signals. As discussed in [9] and subsequently reviewed in Section 3.1, these pre-processing steps essentially compress the M_a LMA signals into one signal. This signal is then used with the pre-processed M_e XM signals to obtain an estimate for the missing part of the RTF vector pertaining to the XMs when there is an available a priori RTF vector for the LMA. Therefore, PWT microphone signals will be adopted for convenience throughout this paper.

To define the transformation operation, an $M_a \times (M_a - 1)$ blocking matrix $\tilde{\mathbf{C}}_a$, and an $M_a \times 1$ fixed beamformer, $\tilde{\mathbf{f}}_a$, are firstly defined such that:

$$\tilde{\mathbf{C}}_a^H \tilde{\mathbf{h}}_a = \mathbf{0}; \quad \tilde{\mathbf{f}}_a^H \tilde{\mathbf{h}}_a = 1 \quad (7)$$

where $\tilde{\mathbf{h}}_a$ is an available a priori RTF vector (which is some pre-determined estimate or approximation of \mathbf{h}_a), and the notation ($\tilde{\cdot}$) refers to quantities based on available a priori knowledge. Using $\tilde{\mathbf{C}}_a$ and $\tilde{\mathbf{f}}_a$, an $(M_a + M_e) \times (M_a + M_e)$ transformation matrix, $\tilde{\mathbf{Y}}$, can be defined as:

$$\tilde{\mathbf{Y}} = \begin{bmatrix} \tilde{\mathbf{Y}}_a & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{M_e} \end{bmatrix} = \begin{bmatrix} [\tilde{\mathbf{C}}_a \tilde{\mathbf{f}}_a] & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{M_e} \end{bmatrix} \quad (8)$$

where $\tilde{\mathbf{Y}}_a = [\tilde{\mathbf{C}}_a \tilde{\mathbf{f}}_a]$ and in general \mathbf{I}_ϑ denotes the $\vartheta \times \vartheta$ identity matrix. Consequently, the transformed speech-plus-noise signals and the transformed noise-only signals are defined respectively as:

$$\tilde{\mathbf{Y}}^H \mathbf{y} = \begin{bmatrix} \tilde{\mathbf{C}}_a^H \mathbf{y}_a \\ \tilde{\mathbf{f}}_a^H \mathbf{y}_a \\ \mathbf{y}_e \end{bmatrix}; \quad \tilde{\mathbf{Y}}^H \mathbf{n} = \begin{bmatrix} \tilde{\mathbf{C}}_a^H \mathbf{n}_a \\ \tilde{\mathbf{f}}_a^H \mathbf{n}_a \\ \mathbf{n}_e \end{bmatrix} \quad (9)$$

This transformation domain is simply the LMA signals that pass through a blocking matrix and a fixed beamformer as is done in the first stage of a typical generalised sidelobe canceller (i.e. the adaptive implementation of an MVDR beamformer) [23], along with the unprocessed XM signals.

A spatial pre-whitening operation can now be defined from the transformed noise-only correlation matrix by using the Cholesky decomposition:

$$\mathbb{E} \left\{ \left(\tilde{\mathbf{Y}}^H \mathbf{n} \right) \left(\tilde{\mathbf{Y}}^H \mathbf{n} \right)^H \right\} = \mathbf{L} \mathbf{L}^H \quad (10)$$

where \mathbf{L} is an $(M_a + M_e) \times (M_a + M_e)$ lower triangular matrix.

A transformed signal vector can then be pre-whitened by pre-multiplying it with \mathbf{L}^{-1} and will be denoted with an underbar ($\underline{\cdot}$). Hence, the signal model for the unprocessed microphone signals from (1) can be expressed in the PWT domain as²:

$$\underline{\mathbf{y}}(k, l) = \mathbf{L}^{-1}(k, l) \tilde{\mathbf{Y}}^H(k, l) \mathbf{y}(k, l) \quad (11)$$

$$= \underbrace{\mathbf{h}(k, l) s_{a,1}(k, l) + \mathbf{n}(k, l)}_{\underline{\mathbf{x}}(k, l)} \quad (12)$$

where $\underline{\mathbf{y}}$ consists of the PWT LMA and XM signals, i.e.

$\underline{\mathbf{y}} = \begin{bmatrix} \mathbf{y}_a^T & \mathbf{y}_e^T \end{bmatrix}^T$, $\underline{\mathbf{n}} = \mathbf{L}^{-1} \tilde{\mathbf{Y}}^H \mathbf{n}$, the PWT RTF vector $\underline{\mathbf{h}} = \mathbf{L}^{-1} \tilde{\mathbf{Y}}^H \mathbf{h}$, and the respective correlation matrices are:

$$\underline{\mathbf{R}}_{\mathbf{y}\mathbf{y}} = \mathbb{E} \left\{ \underline{\mathbf{y}} \underline{\mathbf{y}}^H \right\} = \mathbf{L}^{-1} \tilde{\mathbf{Y}}^H \mathbf{R}_{\mathbf{y}\mathbf{y}} \tilde{\mathbf{Y}} \mathbf{L}^{-H} \quad (13)$$

$$\underline{\mathbf{R}}_{\mathbf{nn}} = \mathbb{E} \left\{ \underline{\mathbf{n}} \underline{\mathbf{n}}^H \right\} = \mathbf{L}^{-1} \tilde{\mathbf{Y}}^H \mathbf{R}_{\mathbf{nn}} \tilde{\mathbf{Y}} \mathbf{L}^{-H} = \mathbf{I}_{(M_a + M_e)} \quad (14)$$

$$\underline{\mathbf{R}}_{\mathbf{xx}} = \mathbb{E} \left\{ \underline{\mathbf{x}} \underline{\mathbf{x}}^H \right\} = \sigma_{s_{a,1}}^2 \underline{\mathbf{h}} \underline{\mathbf{h}}^H \quad (15)$$

where the expression for $\underline{\mathbf{R}}_{\mathbf{nn}}$ is a direct consequence of (10). With the assumption of the desired speech source and noise being uncorrelated, it also holds that $\underline{\mathbf{R}}_{\mathbf{y}\mathbf{y}} = \underline{\mathbf{R}}_{\mathbf{xx}} + \underline{\mathbf{R}}_{\mathbf{nn}}$. In the PWT domain, the estimate of the desired speech signal in the first microphone of the LMA, z_1 , which is equivalent to (6), is then obtained through a linear filtering of the PWT microphone signals, such that:

$$z_1 = \underline{\mathbf{w}}^H \underline{\mathbf{y}} \quad (16)$$

where $\underline{\mathbf{w}} = \mathbf{L}^H \tilde{\mathbf{Y}}^{-1} \mathbf{w}$ is a complex-valued filter³.

3 MVDR with an LMA and XMs

The MVDR beamformer minimises the noise power spectral density after filtering (minimum variance), with a constraint that the desired speech signal should not be subject to any distortion (distortionless response), which is specified by an appropriate RTF vector for the MVDR beamformer. For the unprocessed microphone signals, the MVDR beamformer problem can be formulated as:

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimise}} && \mathbf{w}^H \mathbf{R}_{\mathbf{nn}} \mathbf{w} \\ & \text{s.t.} && \mathbf{w}^H \mathbf{h} = 1 \end{aligned} \quad (17)$$

The solution to (17) yields the optimal filter:

$$\mathbf{w}_{\text{mvdr}} = \frac{\mathbf{R}_{\mathbf{nn}}^{-1} \mathbf{h}}{\mathbf{h}^H \mathbf{R}_{\mathbf{nn}}^{-1} \mathbf{h}} \quad (18)$$

with the desired speech signal estimate, $z_1 = \mathbf{w}_{\text{mvdr}}^H \mathbf{y}$. In practice, both $\mathbf{R}_{\mathbf{nn}}$ and \mathbf{h} are unknown and hence must be estimated.

A data-dependent estimate can typically be obtained for $\mathbf{R}_{\mathbf{nn}}$, for instance by recursive averaging, with a voice activity detector [24] or a speech presence probability (SPP) estimator [25]. This data-dependent estimate will be denoted as $\hat{\mathbf{R}}_{\mathbf{nn}}$ and in general the notation ($\hat{\cdot}$) will refer to any data-dependent estimate.

In the PWT domain, it can be seen that using $\hat{\mathbf{R}}_{\mathbf{nn}}$ in (10) will result in an estimate for the pre-whitening operator as $\hat{\mathbf{L}}$ and hence from (14), $\hat{\mathbf{R}}_{\mathbf{nn}}$ can be expressed as:

$$\hat{\mathbf{R}}_{\mathbf{nn}} = \tilde{\mathbf{Y}}^{-H} \hat{\mathbf{L}} \hat{\mathbf{L}}^H \tilde{\mathbf{Y}}^{-1} \quad (19)$$

Replacing $\mathbf{R}_{\mathbf{nn}}$ in (17) with $\hat{\mathbf{R}}_{\mathbf{nn}}$ in (19) then results in the MVDR beamformer problem formulated in the PWT domain as:

²The dependence on k and l is included here as a reminder and for completeness in the signal model. It will be dropped again unless explicitly required.

³Since the sequence of operations from \mathbf{w} to $\underline{\mathbf{w}}$ is not exactly that of a PWT signal vector, a slightly different notation is used for this quantity.

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimise}} \quad \mathbf{w}^H \mathbf{w} \\ & \text{s.t.} \quad \mathbf{w}^H \mathbf{h} = 1 \end{aligned} \quad (20)$$

where \mathbf{w} is redefined as $\mathbf{w} = \hat{\mathbf{L}}^H \tilde{\mathbf{Y}}^{-1} \mathbf{w}$ and \mathbf{h} is redefined as $\mathbf{h} = \hat{\mathbf{L}}^{-1} \tilde{\mathbf{Y}}^H \mathbf{h}$. The solution to (20) then yields the optimal filter in the PWT domain:

$$\mathbf{w}_{\text{mvdr}} = \frac{\mathbf{h}}{\mathbf{h}^H \mathbf{h}} \quad (21)$$

with the desired speech signal estimate, $z_1 = \mathbf{w}_{\text{mvdr}}^H \mathbf{y}$. As \mathbf{h} is still unknown, however, it means that \mathbf{h} is also unknown and an estimate for this component is still required. Using the same $\hat{\mathbf{R}}_{\text{nn}}$, two general approaches for the estimation of \mathbf{h} can be considered, either making use of an available a priori RTF vector pertaining to the LMA or making use of only the observable microphone data, i.e. a fully data-dependent estimate. The remainder of this section elaborates on these procedures.

3.1 Using an a priori RTF vector

For a microphone configuration consisting of only an LMA, it is not uncommon to use an a priori RTF vector, $\tilde{\mathbf{h}}_{\mathbf{a}}$, in place of the true RTF vector. As mentioned earlier, this may be measured a priori or based on several assumptions regarding the spatial scenario and acoustic environment. For the inclusion of XMs into the microphone configuration, however, the notion of an a priori RTF vector is not so straightforward as no immediate prior knowledge with respect to the XMs can be exploited since there are no restrictions on what type of XMs can be used or where they must be placed in the acoustic environment. Hence, an a priori RTF vector cannot be prescribed, as was the case for the LMA only. However, since a priori information would typically only be available for the LMA, an a priori RTF vector for a microphone configuration of an LMA with XMs can be defined as follows:

$$\tilde{\mathbf{h}} = \begin{bmatrix} \tilde{\mathbf{h}}_{\mathbf{a}}^T & \mathbf{h}_{\mathbf{e}}^T \end{bmatrix}^T \quad (22)$$

which consists partially of the a priori RTF vector pertaining to the LMA, $\tilde{\mathbf{h}}_{\mathbf{a}}$, and partially of the RTF vector pertaining to the XM, $\mathbf{h}_{\mathbf{e}}$, which is unknown and remains to be estimated. The estimate of $\mathbf{h}_{\mathbf{e}}$ will be denoted as $\hat{\mathbf{h}}_{\mathbf{e}}$ to emphasise that it is constrained by the a priori knowledge set by $\tilde{\mathbf{h}}_{\mathbf{a}}$ but estimated from the observed microphone data. In [9], a procedure involving the generalised eigenvalue decomposition (GEVD) was used for obtaining $\hat{\mathbf{h}}_{\mathbf{e}}$ which is subsequently reviewed and re-framed in the PWT domain.

In the PWT domain, using (13)–(15), a rank-1 matrix approximation problem can be firstly formulated to estimate the entire RTF vector [9]:

$$\underset{\sigma_{s_{a,1}}^2, \mathbf{h}}{\text{minimise}} \quad \left\| \left(\hat{\mathbf{R}}_{\text{yy}} - \hat{\mathbf{R}}_{\text{nn}} \right) - \sigma_{s_{a,1}}^2 \hat{\mathbf{L}}^{-1} \tilde{\mathbf{Y}}^H \mathbf{h} \mathbf{h}^H \tilde{\mathbf{Y}} \hat{\mathbf{L}}^{-H} \right\|_F^2 \quad (23)$$

where $\|\cdot\|_F$ is the Frobenius norm, and:

$$\hat{\mathbf{R}}_{\text{yy}} = \hat{\mathbf{L}}^{-1} \tilde{\mathbf{Y}}^H \hat{\mathbf{R}}_{\text{yy}} \tilde{\mathbf{Y}} \hat{\mathbf{L}}^{-H} \quad (24)$$

$$\hat{\mathbf{R}}_{\text{nn}} = \hat{\mathbf{L}}^{-1} \tilde{\mathbf{Y}}^H \hat{\mathbf{R}}_{\text{nn}} \tilde{\mathbf{Y}} \hat{\mathbf{L}}^{-H} = \mathbf{I}_{(M_a + M_e)} \quad (25)$$

where $\hat{\mathbf{R}}_{\text{yy}}$ is the data-dependent estimate of \mathbf{R}_{yy} . From (22), an a priori RTF vector in the PWT domain can be defined as follows:

$$\tilde{\mathbf{h}} = \hat{\mathbf{L}}^{-1} \tilde{\mathbf{Y}}^H \begin{bmatrix} \tilde{\mathbf{h}}_{\mathbf{a}}^T & \mathbf{h}_{\mathbf{e}}^T \end{bmatrix}^T = \hat{\mathbf{L}}^{-1} \begin{bmatrix} \mathbf{0}^T & 1 & \mathbf{h}_{\mathbf{e}}^T \end{bmatrix}^T \quad (26)$$

where $\mathbf{0}$ is a vector of $(M_a - 1)$ zeros. Replacing \mathbf{h} with the a priori RTF vector from (22) then results in:

$$\underset{\sigma_{s_{a,1}}^2, \mathbf{h}_{\mathbf{e}}}{\text{minimise}} \quad \left\| \left(\hat{\mathbf{R}}_{\text{yy}} - \hat{\mathbf{R}}_{\text{nn}} \right) - \sigma_{s_{a,1}}^2 \tilde{\mathbf{h}} \tilde{\mathbf{h}}^H \right\|_F^2 \quad (27)$$

where now only an estimate is required for $\mathbf{h}_{\mathbf{e}}$, which in turn will define the a priori RTF vector. As discussed in [9], it can be observed that it is only the lower $(M_e + 1) \times (M_e + 1)$ blocks of $\hat{\mathbf{R}}_{\text{yy}}$ and $\hat{\mathbf{R}}_{\text{nn}}$ that are required for estimating $\mathbf{h}_{\mathbf{e}}$. Hence, (27) can be reduced to:

$$\underset{\sigma_{s_{a,1}}^2, \mathbf{h}_{\mathbf{e}}}{\text{minimise}} \quad \left\| \left(\hat{\mathbf{R}}_{\text{yy}} - \hat{\mathbf{R}}_{\text{nn}} \right) - \sigma_{s_{a,1}}^2 \mathbf{J}^T \tilde{\mathbf{h}} \tilde{\mathbf{h}}^H \mathbf{J} \right\|_F^2 \quad (28)$$

where $\mathbf{J} = \left[\mathbf{0}_{(M_e+1) \times (M_a-1)} \mid \mathbf{I}_{(M_e+1)} \right]^T$ is a selection matrix, $\hat{\mathbf{R}}_{\text{yy}} = \mathbf{J}^T \hat{\mathbf{R}}_{\text{yy}} \mathbf{J}$, and $\hat{\mathbf{R}}_{\text{nn}} = \mathbf{J}^T \hat{\mathbf{R}}_{\text{nn}} \mathbf{J} = \mathbf{I}_{M_e+1}$. The solution of (28) then follows from a GEVD of the matrix pencil $\left\{ \hat{\mathbf{R}}_{\text{yy}}, \hat{\mathbf{R}}_{\text{nn}} \right\}$ or equivalently from the eigenvalue decomposition (EVD) of $\hat{\mathbf{R}}_{\text{yy}}$ [26]:

$$\hat{\mathbf{R}}_{\text{yy}} = \hat{\mathbf{V}} \hat{\mathbf{\Gamma}} \hat{\mathbf{V}}^H \quad (29)$$

where $\hat{\mathbf{V}}$ is a $(M_e + 1) \times (M_e + 1)$ unitary matrix of eigenvectors and $\hat{\mathbf{\Gamma}}$ is a diagonal matrix with the associated eigenvalues in descending order. The estimate of $\mathbf{h}_{\mathbf{e}}$ then follows from the appropriate scaling of the principal eigenvector, $\hat{\mathbf{v}}_{\mathbf{p}}$:

$$\begin{bmatrix} \mathbf{0} \\ 1 \\ \hat{\mathbf{h}}_{\mathbf{e}} \end{bmatrix} = \frac{\hat{\mathbf{L}} \hat{\mathbf{V}}_{\mathbf{p}}}{\mathbf{e}_{M_a}^T \hat{\mathbf{L}} \hat{\mathbf{V}}_{\mathbf{p}}} = \frac{\hat{\mathbf{L}} \hat{\mathbf{V}}_{\mathbf{p}}}{\hat{l}_{M_a} \hat{v}_{p,1}} \quad (30)$$

where \mathbf{e}_{M_a} is an $(M_a + M_e)$ selection vector consisting of all zeros except for a one in the M_a th position, $\hat{v}_{p,1}$ is the first element of $\hat{\mathbf{v}}_{\mathbf{p}}$, and \hat{l}_{M_a} is the real-valued (M_a, M_a) th

element in $\hat{\mathbf{L}}$. Substitution of this expression into (26) finally yields the a priori RTF vector in the PWT domain as⁴:

$$\tilde{\mathbf{h}} = \frac{1}{\hat{l}_{M_a} \hat{v}_{p,1}} [\mathbf{0} \ \hat{\mathbf{v}}_p]^T \quad (31)$$

Finally, replacing \mathbf{h} in (21) with $\tilde{\mathbf{h}}$ from (31) results in the MVDR beamformer based on a priori knowledge pertaining to the LMA:

$$\tilde{\mathbf{w}}_{\text{mvdr}} = \hat{l}_{M_a} \hat{v}_{p,1}^* [\mathbf{0} \ \hat{\mathbf{v}}_p]^T \quad (32)$$

which will be referred to as MVDR-AP. The corresponding speech estimate is then computed using (16):

$$\tilde{z}_1 = \hat{l}_{M_a} \hat{v}_{p,1} \hat{\mathbf{v}}_p^H \begin{bmatrix} \mathbf{y}_{a,M_a} \\ \mathbf{y}_e \end{bmatrix} \quad (33)$$

As a consequence of incorporating the a priori information into the rank-1 speech model, it can be seen that it is only necessary to filter the last $(M_e + 1)$ elements of \mathbf{y} , i.e. \mathbf{y}_{a,M_a} and \mathbf{y}_e , with the lower order, $(M_e + 1)$ filter defined by $\hat{l}_{M_a} \hat{v}_{p,1}^* \hat{\mathbf{v}}_p$.

3.2 Using a data-dependent RTF vector

In the PWT domain, it is (23) that needs to be solved in order to obtain a fully data-dependent estimate of the RTF vector pertaining to the LMA and the XMs. The solution to (23) follows from a GEVD of the matrix pencil $\{\hat{\mathbf{R}}_{yy}, \hat{\mathbf{R}}_{nn}\}$ or equivalently from the EVD of $\hat{\mathbf{R}}_{yy}$:

$$\hat{\mathbf{R}}_{yy} = \hat{\mathbf{Q}} \hat{\mathbf{\Lambda}} \hat{\mathbf{Q}}^H \quad (34)$$

where $\hat{\mathbf{Q}}$ is an $(M_a + M_e) \times (M_a + M_e)$ unitary matrix of eigenvectors and $\hat{\mathbf{\Lambda}}$ is a diagonal matrix with the associated eigenvalues in descending order. The estimated RTF vector is then given by the principal (first in this case) eigenvector, $\hat{\mathbf{q}}_p$:

$$\hat{\mathbf{h}} = \frac{\tilde{\mathbf{\Upsilon}}^{-H} \hat{\mathbf{L}} \hat{\mathbf{q}}_p}{\hat{\eta}_q} \quad (35)$$

where $\hat{\eta}_q = \mathbf{e}_1^T \tilde{\mathbf{\Upsilon}}^{-H} \hat{\mathbf{L}} \hat{\mathbf{q}}_p$ and \mathbf{e}_1 is an $(M_a + M_e)$ selection vector with a one as the first element and zeros everywhere else. In the PWT domain, this data-dependent RTF vector then becomes:

$$\hat{\mathbf{h}} = \hat{\mathbf{L}}^{-1} \tilde{\mathbf{\Upsilon}}^H \hat{\mathbf{h}} = \frac{\hat{\mathbf{q}}_p}{\hat{\eta}_q} \quad (36)$$

Replacing \mathbf{h} in (21) with $\hat{\mathbf{h}}$ from (36) results in the MVDR beamformer that makes use of a data-dependent RTF vector:

⁴It is acknowledged that there is a slight abuse of notation here as the estimate for $\tilde{\mathbf{h}}$ should be denoted as $\hat{\tilde{\mathbf{h}}}$. However, in favour of legibility and to stress that the estimation is done in accordance to the a priori assumptions set by $\tilde{\mathbf{h}}_a$ is why the notation is maintained as $\tilde{\mathbf{h}}$.

$$\hat{\tilde{\mathbf{w}}}_{\text{mvdr}} = \hat{\eta}_q^* \hat{\mathbf{q}}_p \quad (37)$$

which will be referred to as MVDR-DD. The corresponding speech estimate is then computed using (16):

$$\hat{z}_1 = \hat{\eta}_q \hat{\mathbf{q}}_p^H \mathbf{y} \quad (38)$$

where now all $(M_a + M_e)$ signals need to be filtered as opposed to only $(M_e + 1)$ signals in (33) when an a priori RTF vector is used. In general, the MVDR-DD would also be used for microphone configurations where there is no a priori knowledge available, such as those consisting of external microphones only.

4 Integrated MVDR beamformer

4.1 Quadratically constrained quadratic program

As opposed to relying on only an a priori RTF vector or a data-dependent RTF vector, the merging or integration of both RTF vectors into a single approach can be framed as a quadratically constrained quadratic program (QCQP), firstly with respect to the unprocessed microphone signals:

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimise}} \quad \mathbf{w}^H \hat{\mathbf{R}}_{nn} \mathbf{w} \\ & \text{s.t.} \quad \left| \mathbf{w}^H \tilde{\mathbf{h}} - 1 \right|^2 \leq \tilde{\epsilon}^2 \\ & \quad \quad \left| \mathbf{w}^H \hat{\mathbf{h}} - 1 \right|^2 \leq \hat{\epsilon}^2 \end{aligned} \quad (39)$$

where $\tilde{\epsilon}^2$ and $\hat{\epsilon}^2$ are maximum-tolerated squared deviations from a distortionless response due to $\tilde{\mathbf{h}}$ or $\hat{\mathbf{h}}$ respectively. The constraints of (39) can also be re-written in the standard form [27] as follows:

$$\mathbf{w}^H \tilde{\mathbf{h}} \tilde{\mathbf{h}}^H \mathbf{w} - 2\Re\{\tilde{\mathbf{h}}^H \mathbf{w}\} + 1 - \tilde{\epsilon}^2 \leq 0 \quad (40)$$

$$\mathbf{w}^H \hat{\mathbf{h}} \hat{\mathbf{h}}^H \mathbf{w} - 2\Re\{\hat{\mathbf{h}}^H \mathbf{w}\} + 1 - \hat{\epsilon}^2 \leq 0 \quad (41)$$

where $\Re\{\cdot\}$ denotes the real part of its argument. As the matrices $\hat{\mathbf{R}}_{nn}$, $\tilde{\mathbf{h}} \tilde{\mathbf{h}}^H$, and $\hat{\mathbf{h}} \hat{\mathbf{h}}^H$ are all positive semidefinite, it is then evident that the QCQP of (39) is convex [27]. In the PWT domain, (39) is equivalently:

$$\begin{aligned} & \underset{\mathbf{w}}{\text{minimise}} \quad \mathbf{w}^H \mathbf{w} \\ & \text{s.t.} \quad \left| \mathbf{w}^H \tilde{\mathbf{h}} - 1 \right|^2 \leq \tilde{\epsilon}^2 \\ & \quad \quad \left| \mathbf{w}^H \hat{\mathbf{h}} - 1 \right|^2 \leq \hat{\epsilon}^2 \end{aligned} \quad (42)$$

where $\tilde{\mathbf{h}}$ and $\hat{\mathbf{h}}$ are given in (31) and (36) respectively. Whereas in (20), the hard constraint of \mathbf{h} is replaced by either $\tilde{\mathbf{h}}$ or $\hat{\mathbf{h}}$, (42) can be interpreted as the relaxation of the hard constraints imposed by $\tilde{\mathbf{h}}$ or $\hat{\mathbf{h}}$ by the specified deviations $\tilde{\epsilon}^2$ and $\hat{\epsilon}^2$ respectively. In the following, the

quantities $\left| \underline{\mathbf{w}}^H \underline{\tilde{\mathbf{h}}} - 1 \right|^2$ and $\left| \underline{\mathbf{w}}^H \underline{\hat{\mathbf{h}}} - 1 \right|^2$ are referred to as speech distortions and $\tilde{\epsilon}^2$ and $\hat{\epsilon}^2$ are the respective maximum tolerable speech distortions. Furthermore, the first inequality constraint in (42) will be referred to as the a priori constraint (APC), and the second inequality constraint will be referred to as the data-dependent constraint (DDC).

The QCQP of (39) is in fact a subset of the more general QCQP considered in [28, 29] and as well as an extension to the parametrised multi-channel Wiener filter [30]. In [28, 29], the inequality constraints considered are a set of a priori measured RTF vectors, and in [30], only one inequality constraint is considered. The difference in (39) from both of these approaches is that two inequality constraints are considered, one that relies on a priori knowledge and the other which is fully estimated from the data.

The Lagrangian of (42) is given by:

$$\mathcal{L}(\underline{\mathbf{w}}, \alpha, \beta) = \underline{\mathbf{w}}^H \underline{\mathbf{w}} + \alpha \left(\left| \underline{\mathbf{w}}^H \underline{\tilde{\mathbf{h}}} - 1 \right|^2 - \tilde{\epsilon}^2 \right) + \beta \left(\left| \underline{\mathbf{w}}^H \underline{\hat{\mathbf{h}}} - 1 \right|^2 - \hat{\epsilon}^2 \right) \quad (43)$$

where α and β are Lagrangian multipliers. Taking the partial derivative of (43) with respect to $\underline{\mathbf{w}}$ and setting to zero results in what will be referred to as the integrated MVDR beamformer, MVDR-INT:

$$\underline{\mathbf{w}}_{\text{int}} = \left(\mathbf{I}_{(M_a+M_e)} + \alpha \underline{\tilde{\mathbf{h}}}\underline{\tilde{\mathbf{h}}}^H + \beta \underline{\hat{\mathbf{h}}}\underline{\hat{\mathbf{h}}}^H \right)^{-1} \left(\alpha \underline{\tilde{\mathbf{h}}}\underline{\tilde{\mathbf{h}}}^H + \beta \underline{\hat{\mathbf{h}}}\underline{\hat{\mathbf{h}}}^H \right) \mathbf{e}_1 \quad (44)$$

where the actual values of α and β depend on the prescribed maximum tolerable speech distortions $\tilde{\epsilon}^2$ and $\hat{\epsilon}^2$. It can also be observed that (44) is in fact identical (in the PWT domain) to the integrated MVDR beamformer considered in [22] and hence can be written as a linear combination of $\underline{\tilde{\mathbf{w}}}_{\text{mvdr}}$ and $\underline{\hat{\mathbf{w}}}_{\text{mvdr}}$ with complex weightings⁵ [22]:

$$\underline{\mathbf{w}}_{\text{int}} = g_{ap}(\alpha, \beta) \underline{\tilde{\mathbf{w}}}_{\text{mvdr}} + g_{dd}(\alpha, \beta) \underline{\hat{\mathbf{w}}}_{\text{mvdr}} \quad (45)$$

where $\underline{\tilde{\mathbf{w}}}_{\text{mvdr}}$ and $\underline{\hat{\mathbf{w}}}_{\text{mvdr}}$ are given in (32) and (37), respectively, and the complex weightings are given by:

$$g_{ap}(\alpha, \beta) = \left[\frac{\alpha k_{aa} [1 + \beta(k_{bb} - k_{ab})]}{D} \right] \quad (46)$$

$$g_{dd}(\alpha, \beta) = \left[\frac{\beta k_{bb} [1 + \alpha(k_{aa} - k_{ba})]}{D} \right] \quad (47)$$

where

$$D = \alpha k_{aa} + \beta k_{bb} + \alpha\beta(k_{aa}k_{bb} - k_{ab}k_{ba}) + 1 \quad (48)$$

and

$$k_{aa} = \underline{\tilde{\mathbf{h}}}\underline{\tilde{\mathbf{h}}}^H, \quad k_{bb} = \underline{\hat{\mathbf{h}}}\underline{\hat{\mathbf{h}}}^H, \quad (49)$$

$$k_{ab} = \underline{\tilde{\mathbf{h}}}\underline{\hat{\mathbf{h}}}^H, \quad k_{ba} = \underline{\hat{\mathbf{h}}}\underline{\tilde{\mathbf{h}}}^H. \quad (50)$$

Using the expressions for $\underline{\tilde{\mathbf{w}}}_{\text{mvdr}}$ and $\underline{\hat{\mathbf{w}}}_{\text{mvdr}}$ from (32) and (37) respectively, the resulting speech estimate from the MVDR-INT is then:

$$z_{\text{int}} = g_{ap}^*(\alpha, \beta) \tilde{z}_1 + g_{dd}^*(\alpha, \beta) \hat{z}_1 \quad (51)$$

where \tilde{z}_1 and \hat{z}_1 are defined in (33) and (38) respectively. Hence, the integrated beamformer output is simply a linear combination of the two speech estimates which relied on either a priori information or not.

Once appropriate values are chosen for $\tilde{\epsilon}^2$ and $\hat{\epsilon}^2$, then a package for specifying and solving convex programs such as CVX [31, 32] can be used for solving (42). Alternatively, more computationally efficient methods may be applied such as those proposed in [28, 29], one of which is highlighted in Algorithm 1. Here, a gradient ascent method [33] for solving (42) is described, which is based on solving the dual problem:

$$\begin{aligned} & \underset{(\alpha, \beta)}{\text{maximise}} \quad \mathcal{D}(\alpha, \beta) \\ & \text{s.t. } \alpha \geq 0; \beta \geq 0 \end{aligned} \quad (52)$$

where $\mathcal{D}(\alpha, \beta) = \inf_{\underline{\mathbf{w}}_{\text{int}}} \mathcal{L}(\underline{\mathbf{w}}_{\text{int}}, \alpha, \beta)$ is the infimum of $\mathcal{L}(\underline{\mathbf{w}}_{\text{int}}, \alpha, \beta)$ and referred to as the dual function. As the dual function is concave [27], a gradient ascent procedure can be used to update the values of α and β using the gradients, $\frac{\partial \mathcal{D}(\alpha, \beta)}{\partial \alpha} = \left| \underline{\mathbf{w}}_{\text{int}}^H \underline{\tilde{\mathbf{h}}} - 1 \right|^2 - \tilde{\epsilon}^2$ and $\frac{\partial \mathcal{D}(\alpha, \beta)}{\partial \beta} = \left| \underline{\mathbf{w}}_{\text{int}}^H \underline{\hat{\mathbf{h}}} - 1 \right|^2 - \hat{\epsilon}^2$, i.e. the gradients of the dual function with respect to the particular Lagrange multiplier are the respective constraints. This then gives rise to Algorithm 1 [29], which makes use of the simplified expressions for $\underline{\mathbf{w}}_{\text{int}}$ with the complex-valued weightings as opposed to computing (44) directly. The Lagrangian multipliers, α and β , are then updated via the gradient ascent procedure with the step size γ , whose value can be controlled using a backtracking method [34]. The algorithm continues until the respective gradients are within some specified tolerance, δ .

4.2 Effect of $\tilde{\epsilon}$ and $\hat{\epsilon}$

As the QCQP of (42) in principle is to be solved for every time frame and frequency bin, it can therefore lead to quite a versatile beamformer as the parameters,

⁵It can also be expressed as a convex combination of various beamformers as discussed in [22].

Algorithm 1 Gradient ascent method for solving the QCQP of (42)

- 1: Initialise $\alpha, \beta, \mathbf{w}_{\text{int}}$. Set tolerance, $\delta, n = 0$
- 2: **while** $\left(|\mathbf{w}_{\text{int}}^H(n) \tilde{\mathbf{h}} - 1|^2 - \tilde{\epsilon}^2 \right) > \delta$ OR $\left(|\mathbf{w}_{\text{int}}^H(n) \hat{\mathbf{h}} - 1|^2 - \hat{\epsilon}^2 \right) > \delta$ **do**
- 3: $g_{ap}(n) = g_{ap}(\alpha(n-1), \beta(n-1))$ from (46)
- 4: $g_{dd}(n) = g_{dd}(\alpha(n-1), \beta(n-1))$ from (47)
- 5: $\mathbf{w}_{\text{int}}(n) = g_{ap}(n) \tilde{\mathbf{w}}_{\text{mvdr}} + g_{dd}(n) \hat{\mathbf{w}}_{\text{mvdr}}$ from (45)
- 6: Set γ according to a backtracking method.
- 7: $\alpha(n) = \max \left\{ \alpha(n-1) + \gamma \left(|\mathbf{w}_{\text{int}}^H(n) \tilde{\mathbf{h}} - 1|^2 - \tilde{\epsilon}^2 \right), 0 \right\}$
- 8: $\beta(n) = \max \left\{ \beta(n-1) + \gamma \left(|\mathbf{w}_{\text{int}}^H(n) \hat{\mathbf{h}} - 1|^2 - \hat{\epsilon}^2 \right), 0 \right\}$
- 9: $n = n + 1$
- 10: **end while**

$\tilde{\epsilon}$ and $\hat{\epsilon}$ can be set independently for each frequency in every time frame in order to define the inequality constraints. So although (42) is a well-known QCQP for which there are several methods available to find the solution, it still remains unclear as to what would be a reasonable strategy for setting or tuning $\tilde{\epsilon}$ and $\hat{\epsilon}$ in practice. As opposed to [22], where tuning rules were developed for the Lagrangian multipliers, here a strategy is outlined for tuning $\tilde{\epsilon}$ and $\hat{\epsilon}$, which will in turn compute the appropriate Lagrangian multipliers (for instance as outlined in Algorithm 1), as this is believed to be a more insightful procedure.

In order to develop a strategy for tuning $\tilde{\epsilon}$ and $\hat{\epsilon}$, it will be useful to observe the constraints of (42) in more detail. The derivations that follow will reveal that the space spanned by $\tilde{\epsilon}$ and $\hat{\epsilon}$ can be divided into four distinct regions as illustrated in Fig. 1, where each of these regions corresponds to a particular set of constraints being active.

Firstly, substitution of $\mathbf{w}_{\text{int}} = \mathbf{0}$ into the APC and DDC from (42) shows that when $\tilde{\epsilon} > 1$ and $\hat{\epsilon} > 1$, both the APC and the DDC are inactive. This condition therefore defines the upper-right region (region I) of Fig. 1 and indeed corresponds to a complete attenuation of the microphone signals, i.e. a zero output signal.

For the case when $\hat{\epsilon} \rightarrow \infty$, i.e. when the DDC is inactive, then $\beta \rightarrow 0$. If the APC is still active however, it becomes⁶:

$$\left| \mathbf{w}_{\text{int}}^H \tilde{\mathbf{h}} - 1 \right| \leq \tilde{\epsilon} \quad (53)$$

Furthermore, if $0 \leq \tilde{\epsilon} \leq 1$, then it can be deduced that:

$$\lim_{\substack{\hat{\epsilon} \rightarrow \infty \\ 0 \leq \tilde{\epsilon} \leq 1}} \mathbf{w}_{\text{int}} = (1 - \tilde{\epsilon}) \tilde{\mathbf{w}}_{\text{mvdr}} \quad (54)$$

Substitution of (54) into (53) readily makes this evident, recalling that $\tilde{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} = 1$. It is worthwhile to also note that by using (46), the relationship between α and $\tilde{\epsilon}$ for $0 \leq \tilde{\epsilon} \leq 1$ is then given as:

$$\lim_{\substack{\hat{\epsilon} \rightarrow \infty \\ 0 \leq \tilde{\epsilon} \leq 1}} \alpha = \frac{1}{k_{aa}} \frac{(1 - \tilde{\epsilon})}{\tilde{\epsilon}} \quad (55)$$

In regard to the DDC, as $\hat{\epsilon}$ is decreased (from $\hat{\epsilon} \rightarrow \infty$), it remains inactive until $\left| \mathbf{w}_{\text{int}}^H \hat{\mathbf{h}} - 1 \right| = \hat{\epsilon}$. By substitution of (54) into the DDC of (42), the value of $\hat{\epsilon}$ at which the DDC becomes active, $\hat{\epsilon}_o$, is given by:

$$\hat{\epsilon}_o = \left| \tilde{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} (1 - \tilde{\epsilon}) - 1 \right| \quad (56)$$

In the limits of $\tilde{\epsilon}$, when $\tilde{\epsilon} = 1$, $\hat{\epsilon}_o = 1$, and when $\tilde{\epsilon} = 0$, $\hat{\epsilon}_o = \left| \tilde{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right|$, where depending on $\tilde{\mathbf{h}}$, $\left| \tilde{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right| < 1$ or $\left| \tilde{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right| \geq 1$. The range of values obtained for $\hat{\epsilon}_o$ from (56) within the domain $0 \leq \tilde{\epsilon} \leq 1$ define what will be referred to as the DDC bounding curve as depicted in Fig. 1. Hence, region II in Fig. 1 is enclosed by the DDC bounding curve, $\tilde{\epsilon} = 0$ and $\tilde{\epsilon} = 1$, representing the space where the APC is active and the DDC is inactive.

A similar analysis can be followed starting from the case when $\tilde{\epsilon} \rightarrow \infty$, i.e. when the APC is inactive and hence $\alpha \rightarrow 0$. If the DDC is still active however, it becomes:

$$\left| \mathbf{w}_{\text{int}}^H \hat{\mathbf{h}} - 1 \right| \leq \hat{\epsilon} \quad (57)$$

When $0 \leq \hat{\epsilon} \leq 1$, then the following relationships can be deduced:

$$\lim_{\substack{\tilde{\epsilon} \rightarrow \infty \\ 0 \leq \hat{\epsilon} \leq 1}} \mathbf{w}_{\text{int}} = (1 - \hat{\epsilon}) \hat{\mathbf{w}}_{\text{mvdr}} \quad (58)$$

$$\lim_{\substack{\tilde{\epsilon} \rightarrow \infty \\ 0 \leq \hat{\epsilon} \leq 1}} \beta = \frac{1}{k_{bb}} \frac{(1 - \hat{\epsilon})}{\hat{\epsilon}} \quad (59)$$

Finally, for the APC, as $\tilde{\epsilon}$ is decreased (from initially $\tilde{\epsilon} \rightarrow \infty$), the value, $\tilde{\epsilon}_o$, at which this constraint becomes active is given by:

$$\tilde{\epsilon}_o = \left| \hat{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} (1 - \hat{\epsilon}) - 1 \right| \quad (60)$$

⁶The square root has been taken on both sides of the inequality from (42) in order to simplify the derivations that follow.

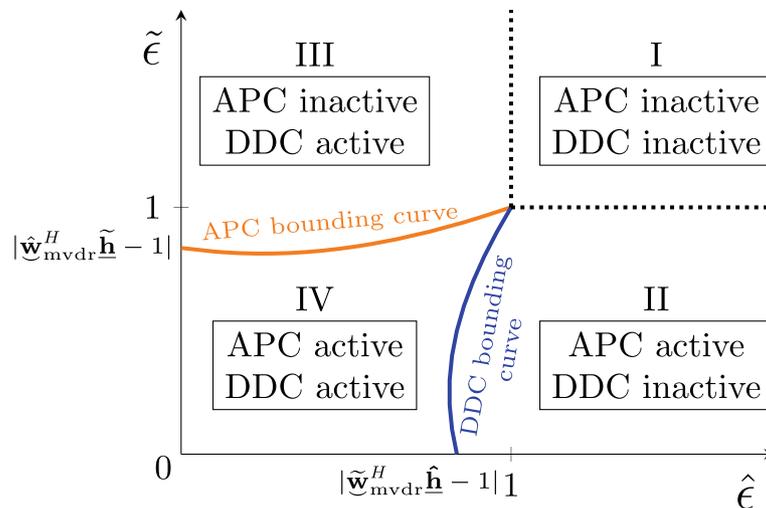


Fig. 1 Depiction of the four regions for which the a priori constraint (APC) and the data-dependent constraint (DDC) may be active or inactive within the space spanned by the maximum tolerable speech distortion parameters, $\tilde{\epsilon}$ and $\hat{\epsilon}$. The curve dividing the regions II and IV is the DDC bounding curve defined when the equality is satisfied from (56). The curve dividing the regions III and IV is the APC bounding curve defined when the equality is satisfied from (60)

In the limits of $\hat{\epsilon}$, when $\hat{\epsilon} = 1$, $\tilde{\epsilon}_o = 1$, and when $\hat{\epsilon} = 0$, $\tilde{\epsilon}_o = \left| \hat{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right|$, where depending on $\tilde{\mathbf{h}}$, $\left| \hat{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right| < 1$ or $\left| \hat{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right| \geq 1$. The range of values obtained for $\tilde{\epsilon}_o$ from (60) within the domain $0 \leq \hat{\epsilon} \leq 1$ define what will be referred to as the APC bounding curve as depicted in Fig. 1. Hence, region III in Fig. 1 is enclosed by the APC bounding curve, $\hat{\epsilon} = 0$ and $\hat{\epsilon} = 1$, representing the space where the APC is inactive and the DDC is active.

Finally, in the lower-left region, region IV, both the APC and the DDC become active within the area enclosed by the APC and DDC bounding curve. It should be kept in mind that Fig. 1 is only an illustration and that the shape of the area for which the APC and DDC are both active can change depending on the RTF vectors, $\tilde{\mathbf{h}}$ and $\hat{\mathbf{h}}$. For instance, Fig. 1 shows $\left| \hat{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}} - 1 \right| < 1$ and $\left| \hat{\mathbf{w}}_{\text{mvdr}}^H \hat{\mathbf{h}} - 1 \right| < 1$ (points on the axes), whereas it is possible that either of these points may be greater than or equal to one.

5 Confidence metric and tuning

5.1 Confidence metric

One of the ingredients towards developing a tuning strategy for setting appropriate values for $\tilde{\epsilon}$ and $\hat{\epsilon}$ is that of a confidence metric, which is indicative of the confidence in the accuracy of the data-dependent RTF vector. In [22], it was proposed that a principal generalised eigenvalue

resulting from the data-dependent estimation procedure be used as such a confidence metric. In the following, it is proposed again to use such a metric; however, due to the formulation in the PWT domain, the principal eigenvalue, $\hat{\lambda}_1$ from the EVD in (34) will be used. It can be shown that $\hat{\lambda}_1$ is equivalent to the resulting posterior SNR when the MVDR-DD is applied and therefore serves as a reasonable metric for making a decision with respect to the accuracy of the data-dependent RTF. For the MVDR-DD in (37), the resulting posterior SNR is given by:

$$\widehat{\text{SNR}}_{\text{DD}} = \frac{\hat{\mathbf{w}}_{\text{mvdr}}^H \hat{\mathbf{R}}_{\text{yy}} \hat{\mathbf{w}}_{\text{mvdr}}}{\hat{\mathbf{w}}_{\text{mvdr}}^H \hat{\mathbf{R}}_{\text{nn}} \hat{\mathbf{w}}_{\text{mvdr}}} \quad (61)$$

where it is recalled that $\hat{\mathbf{R}}_{\text{nn}} = \mathbf{I}_{(M_s+M_e)}$. Substitution of (34) and (37)⁷ into (61) results in $\widehat{\text{SNR}}_{\text{DD}} = \hat{\lambda}_1$.

As in [22], $\hat{\lambda}_1$ can then be used in a logistic function to define the confidence metric, $F(l)$ ⁸:

$$F(l) = \frac{1}{1 + e^{-\rho(10 \log_{10}(\hat{\lambda}_1(l)) - \lambda_t)}} \quad (62)$$

where $F(l) \in [0, 1]$, ρ controls the gradient of the transition from 0 to 1, and λ_t is a threshold (in dB), beyond which $F(l) \rightarrow 1$. Hence, as $10 \log_{10}(\hat{\lambda}_1(l))$ increases beyond λ_t , then $F(l) \rightarrow 1$, indicating high confidence in the accuracy of the data-dependent RTF vector. On the other hand, as

⁷Recall that $\hat{\mathbf{w}}_{\text{mvdr}}$ can be equivalently expressed as $\hat{\mathbf{w}}_{\text{mvdr}} = \hat{\eta}_0^* \hat{\mathbf{Q}} \mathbf{e}_1$.

⁸The time index is reintroduced here to reinforce that these quantities are to be computed in each time frame. All frequencies are still treated equivalently.

$10 \log_{10}(\hat{\lambda}_1(l))$ decreases below λ_t , then $F(l) \rightarrow 0$, indicating low confidence in the accuracy of the data-dependent RTF vector.

5.2 Tuning strategy

With the depiction of the space spanned by $\tilde{\epsilon}$ and $\hat{\epsilon}$ from Fig. 1 in mind, a general two-step procedure can be followed to establish a particular tuning strategy:

1. Choose two points on the $\{\hat{\epsilon}, \tilde{\epsilon}\}$ plane: ϵ_{AP} and ϵ_{DD} . The coordinates of ϵ_{AP} , $\{\hat{\epsilon}_{AP}, \tilde{\epsilon}_{AP}\}$, will specify the maximum tolerable speech distortions for the case when there is no confidence in the accuracy of the data-dependent RTF vector. The coordinates of ϵ_{DD} , $\{\hat{\epsilon}_{DD}, \tilde{\epsilon}_{DD}\}$, on the other hand, will specify the maximum tolerable speech distortions for the case when there is complete confidence in the accuracy of the data-dependent RTF vector.
2. Define an appropriate path in order to connect ϵ_{AP} and ϵ_{DD} , where the variation along this path would be a function of the confidence metric, $F(l)$. As $F(l)$ changes in each time-frequency segment, different values of $\hat{\epsilon}$ and $\tilde{\epsilon}$ will be chosen along this path and subsequently used in the QCQP from (42).

Figure 2 depicts three examples of how such a general tuning strategy can be interpreted in the $\{\hat{\epsilon}, \tilde{\epsilon}\}$ plane, where a linear path has been used to connect the points, ϵ_{AP} and ϵ_{DD} . Before further elaborating on Fig. 2, however, one possible tuning strategy will be briefly outlined. In this strategy, ϵ_{AP} and ϵ_{DD} are chosen by making use of the relationship between the integrated MVDR and the so-called speech distortion weighted multi-channel Wiener filter (SDW-MWF) [35, 36]. Although ϵ_{AP} and ϵ_{DD} can in general be chosen without making use of this relation, it is done to highlight how the speech distortion parameter, μ , from the SDW-MWF is related to the maximum tolerable speech distortion parameters of the integrated MVDR, especially as this μ is a well-established trade-off parameter. For the path connecting ϵ_{AP} and

ϵ_{DD} , a linear path will be defined using the confidence metric, $F(l)$.

In the PWT domain, the cost function for the SDW-MWF is given by:

$$\underset{\underline{\mathbf{w}}}{\text{minimise}} \quad \mu \underline{\mathbf{w}}^H \underline{\mathbf{w}} + \sigma_{s_{a,1}}^2 \left| \underline{\mathbf{w}}^H \underline{\mathbf{h}} - 1 \right|^2 \quad (63)$$

which consists of two terms, the first corresponding to the noise power spectral density after filtering and the second corresponding to the speech distortion. The speech distortion parameter $\mu \in (0, \infty)$ is used to trade-off between the amount of noise reduction and speech distortion, where larger values of μ put more emphasis on reducing the noise and smaller values put more emphasis on reducing the speech distortion. Two separate SDW-MWF formulations can then be considered for $\tilde{\underline{\mathbf{h}}}$ and $\hat{\underline{\mathbf{h}}}$ respectively:

$$\underset{\underline{\mathbf{w}}}{\text{minimise}} \quad \tilde{\mu} \underline{\mathbf{w}}^H \underline{\mathbf{w}} + \sigma_{s_{a,1}}^2 \left| \underline{\mathbf{w}}^H \tilde{\underline{\mathbf{h}}} - 1 \right|^2 \quad (64)$$

$$\underset{\underline{\mathbf{w}}}{\text{minimise}} \quad \hat{\mu} \underline{\mathbf{w}}^H \underline{\mathbf{w}} + \sigma_{s_{a,1}}^2 \left| \underline{\mathbf{w}}^H \hat{\underline{\mathbf{h}}} - 1 \right|^2 \quad (65)$$

where $\tilde{\mu} \in (0, \infty)$ and $\hat{\mu} \in (0, \infty)$ are the separate speech distortion parameters for each cost function. The solutions to (64) and (65) are then respectively given by:

$$\tilde{\underline{\mathbf{w}}}_{\text{sdw}} = \left(\tilde{\mu} \mathbf{I}_{(M_a+M_e)} + \hat{\sigma}_{s_{a,1}}^2 \tilde{\underline{\mathbf{h}}} \tilde{\underline{\mathbf{h}}}^H \right)^{-1} \hat{\sigma}_{s_{a,1}}^2 \tilde{\underline{\mathbf{h}}} \tilde{\underline{\mathbf{h}}}^H \mathbf{e}_1 \quad (66)$$

$$\hat{\underline{\mathbf{w}}}_{\text{sdw}} = \left(\hat{\mu} \mathbf{I}_{(M_a+M_e)} + \hat{\sigma}_{s_{a,1}}^2 \hat{\underline{\mathbf{h}}} \hat{\underline{\mathbf{h}}}^H \right)^{-1} \hat{\sigma}_{s_{a,1}}^2 \hat{\underline{\mathbf{h}}} \hat{\underline{\mathbf{h}}}^H \mathbf{e}_1 \quad (67)$$

where $\hat{\sigma}_{s_{a,1}}^2$ is an estimate of $\sigma_{s_{a,1}}^2$. On comparing the $\underline{\mathbf{w}}_{\text{int}}$ in (44) to (66) and (67), it can be observed that there is a relationship between the integrated MVDR beamformer and the SDW-MWF. By considering the expressions written as an MVDR beamformer followed by a single-channel post filter [36], it can be deduced that [22]:

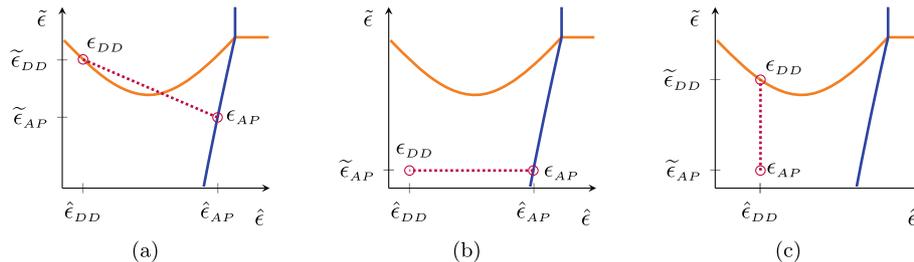


Fig. 2 Depiction of three different tuning strategies **(a)** trading off the maximum tolerable speech distortions between the APC and DDC, **(b)** fixed maximum tolerable speech distortion for the APC but variable maximum tolerable speech distortion for the DDC, and **(c)** fixed maximum tolerable speech distortion for the DDC but variable maximum tolerable speech distortion for the APC

$$\alpha = \frac{\hat{\sigma}_{s_{a,1}}^2}{\tilde{\mu}} \text{ when } \beta = 0 \quad (68)$$

$$\beta = \frac{\hat{\sigma}_{s_{a,1}}^2}{\tilde{\mu}} \text{ when } \alpha = 0 \quad (69)$$

Proceeding to define the coordinates of ϵ_{AP} , (68) is substituted into (55) to obtain a value for $\tilde{\epsilon}$ as:

$$\tilde{\epsilon}_{AP} = \frac{\tilde{\mu}}{\tilde{\mu} + \hat{\sigma}_{s_{a,1}}^2 k_{aa}} \quad (70)$$

Hence, the range of values for $\tilde{\mu}$ are essentially compressed into a range of values for $\tilde{\epsilon}_{AP}$ such that $0 \leq \tilde{\epsilon}_{AP} \leq 1$. This means that $\tilde{\epsilon}_{AP}$ can be chosen to be within this range without having to specify $\tilde{\mu}$. However, (70) serves to clarify how the choice of $\tilde{\epsilon}_{AP}$ is related to the cost function of (64).

Using the value of $\tilde{\epsilon}_{AP}$ in (56) then yields a range of choices for $\hat{\epsilon}_{AP}$ such that $\hat{\epsilon}_{AP} \leq \hat{\epsilon}_o$:

$$\hat{\epsilon}_{AP} \leq \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \hat{\mathbf{h}}(1 - \tilde{\epsilon}_{AP}) - 1 \right| \quad (71)$$

If $\hat{\epsilon}_{AP} = \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \hat{\mathbf{h}}(1 - \tilde{\epsilon}_{AP}) - 1 \right|$, then ϵ_{AP} lies on the DDC bounding curve of Fig. 1. For all values of $\hat{\epsilon}$ such that $\hat{\epsilon} > \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \hat{\mathbf{h}}(1 - \tilde{\epsilon}_{AP}) - 1 \right|$, the DDC remains inactive and hence setting a value of $\hat{\epsilon}$ within this region will always result in the same achievable⁹ speech distortion defined by $\left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \hat{\mathbf{h}}(1 - \tilde{\epsilon}_{AP}) - 1 \right|$. Furthermore, when the DDC is inactive, then (68) holds, so that values of $\tilde{\epsilon}$ and $\hat{\epsilon}$ in region II from Fig. 1 would result in the SDW-MWF from (66).

Similarly, by firstly substituting (69) in (59) and making use of (60), the coordinates $\{\hat{\epsilon}_{DD}, \tilde{\epsilon}_{DD}\}$ of ϵ_{DD} can be defined as:

$$\hat{\epsilon}_{DD} = \frac{\hat{\mu}}{\hat{\mu} + \hat{\sigma}_{s_{a,1}}^2 k_{bb}} \quad (72)$$

$$\tilde{\epsilon}_{DD} \leq \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}_{DD}) - 1 \right| \quad (73)$$

Now if $\tilde{\epsilon}_{DD} = \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}_{DD}) - 1 \right|$, then ϵ_{DD} lies on the APC bounding curve of Fig. 1. Additionally, for all values of $\tilde{\epsilon}$ such that $\tilde{\epsilon} > \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}_{DD}) - 1 \right|$, the APC remains inactive and hence setting a value of $\tilde{\epsilon}$ within this region will always result in the same achievable speech distortion defined by $\left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}_{DD}) - 1 \right|$. Furthermore, when the APC is inactive, then (69) holds, so that values of $\tilde{\epsilon}$ and $\hat{\epsilon}$ in region III from Fig. 1 would result in the SDW-MWF from (67).

⁹Achievable here is meant to differentiate between the actual speech distortion that is obtained and the maximum tolerable value that was specified.

The insight of Fig. 1 and additional value of the MVDR-INT as compared to the SDW-MWF is now apparent. Given the two SDW-MWF solutions from (66) and (67), it is not immediately clear how to optimally interpolate between them by using a linear combination of the filters themselves. In Fig. 1, however, it can be seen that an optimal interpolation between (66) and (67), i.e. between regions II and III, can be achieved through the specification of the maximum tolerable speech distortion parameters, $\tilde{\epsilon}$ and $\hat{\epsilon}$ along some path from region II to region III. In essence, the MVDR-INT has introduced region IV, which serves as a bridge for connecting regions II and III, thereby facilitating the use of both the priori and data-dependent RTF vectors. This then corresponds to the second step of the general procedure for tuning, where ϵ_{AP} and ϵ_{DD} are to be connected. Here, it is proposed to use the confidence metric, $F(l)$ to perform a linear interpolation between ϵ_{AP} and ϵ_{DD} to yield the values for $\tilde{\epsilon}$ and $\hat{\epsilon}$ respectively as:

$$\hat{\epsilon} = (1 - F(l)) \hat{\epsilon}_{AP} + F(l) \hat{\epsilon}_{DD} \quad (74)$$

$$\tilde{\epsilon} = (1 - F(l)) \tilde{\epsilon}_{AP} + F(l) \tilde{\epsilon}_{DD} \quad (75)$$

which are subsequently squared to be used in the QCQP from (42). Consequently, as the confidence in the accuracy of the data-dependent RTF vector increases, the maximum tolerable speech distortions will be specified by values tending towards $\{\hat{\epsilon}_{DD}, \tilde{\epsilon}_{DD}\}$. On the contrary, as this confidence decreases, maximum tolerable speech distortions will be specified by values tending towards $\{\hat{\epsilon}_{AP}, \tilde{\epsilon}_{AP}\}$.

Returning focus to Fig. 2, the three examples of a tuning strategy can now be understood. A particular realisation of the APC and the DDC bounding curves has been plotted and the intersecting point of both curves corresponds to the $\{1, 1\}$ coordinate (recall Fig. 1). In the tuning of Fig. 2a, as $F(l)$ increases, the path along the dotted line is taken from ϵ_{AP} to arrive at ϵ_{DD} which gradually sets a larger value of $\tilde{\epsilon}$ for the APC and a smaller value of $\hat{\epsilon}$ for the DDC. Depending on the particular realisation of the APC and DDC bounding curves, it may be that such a path can entirely lie within the area enclosed by these curves or part of it may lie outside as shown in Fig. 2a. The latter is in fact a fortunate circumstance because the achieved speech distortion corresponding to the inactive constraint will actually be lower than what was prescribed by the tuning. In the case of Fig. 2a for instance, when the linear path is above the APC bounding curve, it means that $\tilde{\epsilon} > \left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}) - 1 \right|$ (recall (60)). Since beyond $\left| \underset{\sim}{\hat{\mathbf{w}}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}) - 1 \right|$ the APC continues to be inactive, the actual speech distortion that would

be achieved in relation to this constraint would correspond to $\left| \tilde{\mathbf{w}}_{\text{mvdr}}^H \tilde{\mathbf{h}}(1 - \hat{\epsilon}) - 1 \right|$, which is by definition less than $\tilde{\epsilon}$. Hence, although there is a linear path from ϵ_{AP} to ϵ_{DD} , at the point where this linear path intersects with the APC bounding curve, the actual speech distortions that would be achieved are those that continue along the APC bounding curve in order to arrive at ϵ_{DD} .

The tunings depicted in Fig. 2b and c are representative of strategies where the maximum tolerable speech distortion is fixed for one of the constraints, and only the maximum tolerable speech distortion for the other constraint is tuned. In Fig. 2b, ϵ_{DD} is defined by setting $\tilde{\epsilon}_{DD} = \tilde{\epsilon}_{AP}$, so that the maximum tolerable speech distortion for the APC is fixed. $\hat{\epsilon}$ is then tuned according to (74). This is representative of a case where the APC is always active and the DDC is only included if there is confidence in the accuracy of the data-dependent RTF vector. Figure 2c depicts an opposite strategy, where now ϵ_{AP} is set by setting $\hat{\epsilon}_{AP} = \hat{\epsilon}_{DD}$, so that the maximum tolerable speech distortion for the DDC is fixed.

6 Evaluation and discussion

In order to gain further insight into the behaviour of the integrated MVDR beamformer using the QCQP formulation, a simulation was firstly considered involving only an LMA without XMs. As will be demonstrated, observing such a scenario facilitates the visualisation of the theoretical beam patterns that would be generated under different tuning strategies. Following this simulation, recorded data from an acoustic scenario involving behind-the-ear dummy¹⁰ hearing aid microphones along with XMs in a cocktail party scenario was then analysed and evaluated.

6.1 Beam patterns for a linear microphone array

As the notion of a traditional beam pattern is not immediately extended to the case of an LMA with XMs¹¹, the following beam patterns are generated using an LMA only.

For visualising the beam patterns, a linear LMA consisting of 4 microphones and 5-cm spacing was considered. Two anechoic RTF vectors, simulating an a priori RTF vector, $\hat{\mathbf{h}}_{\mathbf{a}}$, and a data-dependent RTF vector, $\tilde{\mathbf{h}}_{\mathbf{a}}$, were computed according to a far-field approximation, i.e. $\left[1 e^{-j2\pi f \tau_2(\theta)} e^{-j2\pi f \tau_3(\theta)} e^{-j2\pi f \tau_4(\theta)} \right]^T$, where f is the frequency (Hz) which was set to 3 kHz, $\tau_m(\theta) = \frac{(m-1)0.05 \cos(\theta)}{c}$ is the relative time delay between the m^{th} microphone and the reference microphone (the microphone closest to the desired speech source) of the LMA, θ is the angle of the desired speech source, and $c = 345 \text{ m s}^{-1}$

is the speed of sound. For $\tilde{\mathbf{h}}_{\mathbf{a}}, \theta = 0^\circ$ and for $\hat{\mathbf{h}}_{\mathbf{a}}, \theta = 60^\circ$. Using this definition of $\tilde{\mathbf{h}}_{\mathbf{a}}, \hat{\mathbf{C}}_{\mathbf{a}}$, and $\tilde{\mathbf{f}}_{\mathbf{a}}$ were defined accordingly from (7) and $\tilde{\mathbf{Y}}_{\mathbf{a}}$ from (8). With $\mathbf{R}_{\mathbf{n}_a \mathbf{n}_a} = \mathbf{I}_{M_a}$, the pre-whitening operation from (10) was then computed but with $\tilde{\mathbf{Y}}_{\mathbf{a}}$ instead of $\tilde{\mathbf{Y}}$, and hence denoted as $\tilde{\mathbf{L}}_{\mathbf{a}}$. In the PWT domain, the respective RTF vectors are given by $\tilde{\mathbf{h}}_{\mathbf{a}} = \tilde{\mathbf{L}}_{\mathbf{a}}^{-1} \tilde{\mathbf{Y}}_{\mathbf{a}}^H \tilde{\mathbf{h}}_{\mathbf{a}}$ and $\hat{\mathbf{h}}_{\mathbf{a}} = \tilde{\mathbf{L}}_{\mathbf{a}}^{-1} \tilde{\mathbf{Y}}_{\mathbf{a}}^H \hat{\mathbf{h}}_{\mathbf{a}}$. The optimal PWT domain filters, $\tilde{\mathbf{w}}_{\text{mvdr}}$, and $\hat{\mathbf{w}}_{\text{mvdr}}$ were then computed as in (21), but using either $\tilde{\mathbf{h}}_{\mathbf{a}}$ or $\hat{\mathbf{h}}_{\mathbf{a}}$. Finally, (74) and (75) were used to $\tilde{\epsilon}$ and $\hat{\epsilon}$, after which (42) was then solved using CVX [31, 32] to yield the integrated MVDR beamformer for the LMA only, denoted as $\tilde{\mathbf{w}}_{\text{int}}$. The beam patterns were computed as $|\tilde{\mathbf{w}}_{\text{int}}^H \mathbf{h}(\theta)|$, where $\mathbf{h}(\theta)$ is the PWT domain RTF vector corresponding to an angle, θ .

Figure 3 illustrates the resulting beam patterns for two tuning strategies for different values of $F(l)$ (in this case $l = 1$ and hence the dependence on l is omitted). The left-hand plot of Fig. 3 corresponds to a tuning strategy similar to that depicted in Fig. 2a, where there is a trade-off between the two constraints. For this strategy, $\tilde{\mu} = \hat{\mu} = 0.2$ and $\hat{\sigma}_{\text{sa},1}^2 = 1$, which means that ϵ_{AP} and ϵ_{DD} were fairly close to the x -axis and y -axis respectively. As F increases, the beam pattern is clearly seen to evolve from focusing on the a priori direction of 0° to eventually that of the data-dependent direction of 60° . As a linear path is followed, at the midpoint, both $\tilde{\epsilon}$ and $\hat{\epsilon}$ are of a similarly larger values, which explains the nature of the lower magnitude in the beam pattern during the transition.

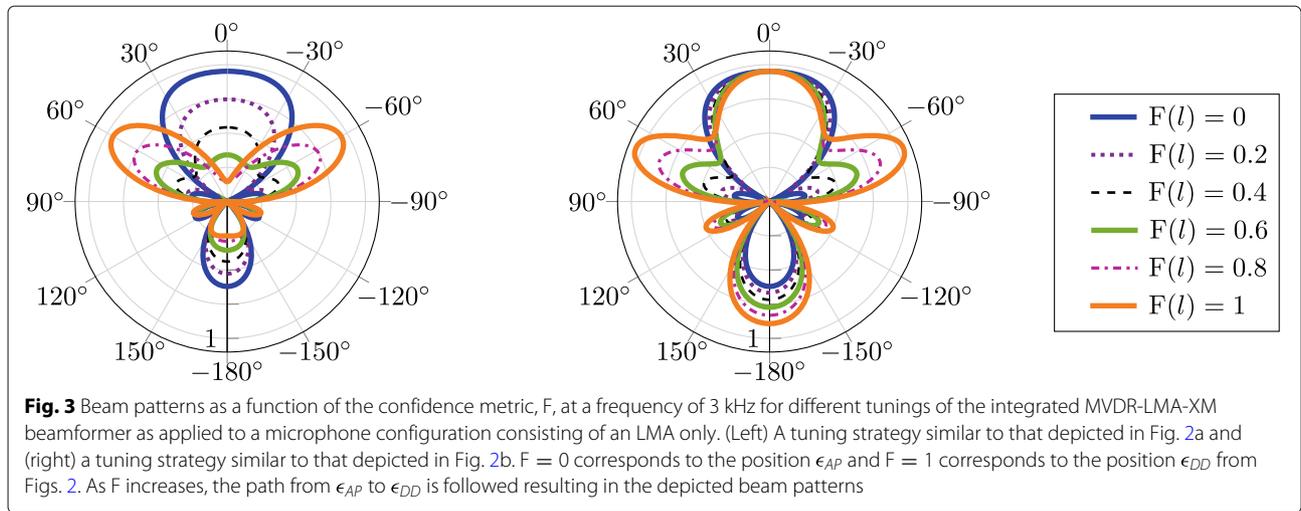
The right-hand plot of Fig. 3 corresponds to a tuning strategy as depicted in Fig. 2b, i.e. when the APC is always active. As F increases, it can be observed that the beam in the a priori direction of 0° is maintained, while more gain is attributed to the data-dependent direction of 60° . In this particular case, however, it is noted that although the response at 60° is in accordance with the maximum tolerable speech distortion prescribed, there is a slight tilt of the beam towards 68° as compared to if only the DDC was active. Nevertheless, this can still be a useful tuning strategy for cases when a high confidence is placed on the a priori RTF vector.

6.2 Effect of $\tilde{\epsilon}$ and $\hat{\epsilon}$

In this section, the effect of $\tilde{\epsilon}$ and $\hat{\epsilon}$ on the behaviour of the integrated MVDR beamformer for the case of an LMA and XMs is further investigated using recorded audio data. A batch processing framework will be applied so as to observe an average performance at a single frequency. In the following section, the processing will be done using a Weighted Overlap and Add (WOLA) framework [37] and a broadband performance will be assessed.

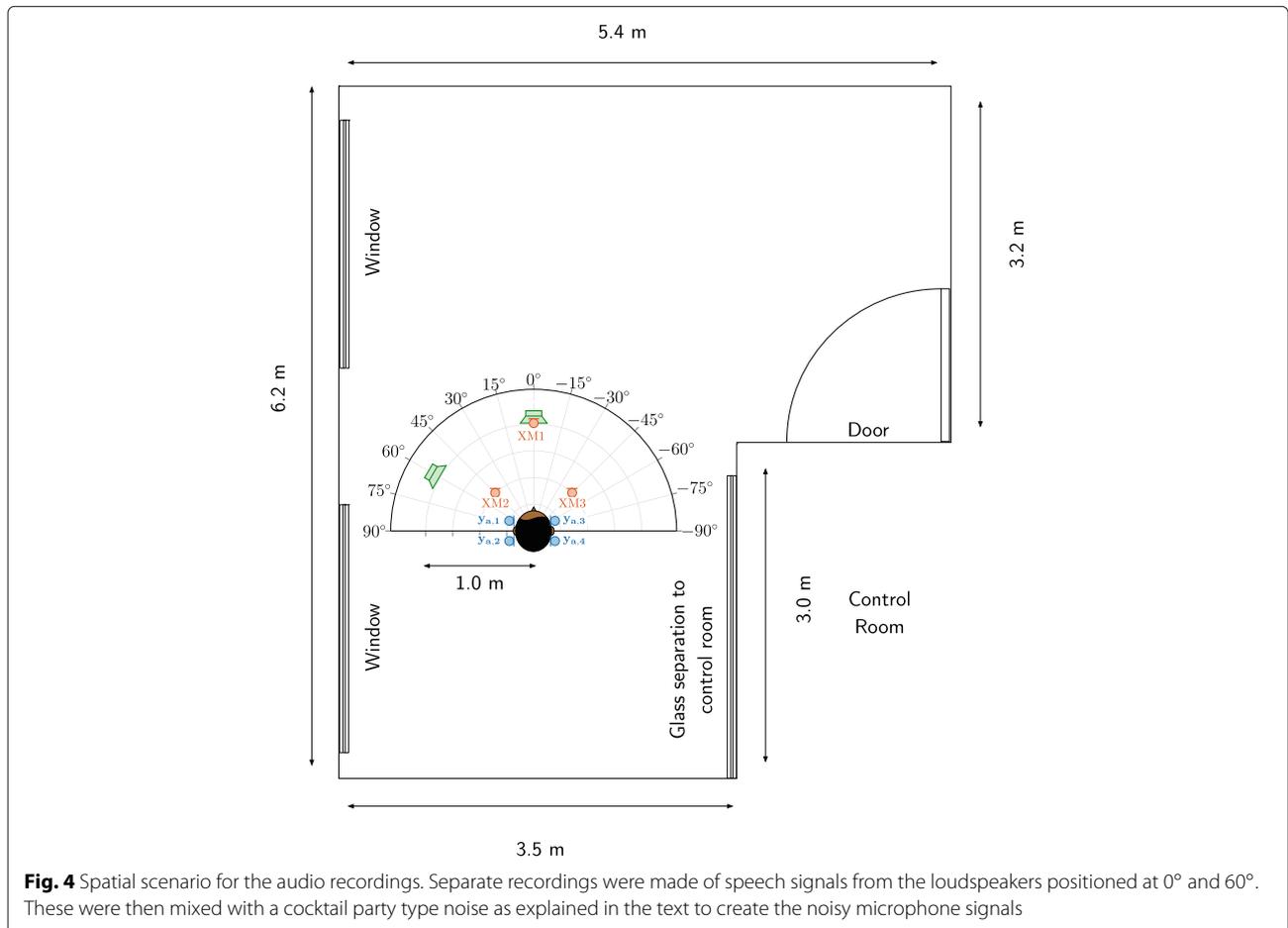
¹⁰This means that only the microphone signals alone, without any processing, are captured.

¹¹The complication arises in that some of the XMs can be in the nearfield with respect to the desired source. A visualisation can nevertheless be created, but will have to be considered within a plane or volume with Cartesian coordinates.



Audio recordings of speech and noise were made in the laboratory room as depicted in Fig. 4, which has a reverberation time of approximately 1.5 s. A Neumann KU-100 dummy head was placed in a central location of the room and equipped with two (i.e. left and right) behind-the-ear

hearing aids, each consisting of two microphones spaced approximately 1.3 cm apart. Hence, in the following, the LMA is considered as having a total of four microphones, i.e. the stacked left ear and right ear microphones. The first microphone of the left ear hearing aid was used as the



reference microphone. Three omnidirectional XMs (two AKG CK32 microphones and one AKG CK970 microphone) were placed at heights of 1 m from the floor and at varying distances from the dummy head as shown in Fig. 4. A Genelec 8030C loudspeaker was placed at 1 m and different azimuth angles from the dummy head to generate a speech signal from a male speaker [38]. The loudspeaker and the dummy head were placed at a height of approximately 1.3 m from the floor (only angles 0° and 60° were used as shown in Fig. 4). For the noise, a cocktail party scenario was re-created. With the same configuration of the dummy head and external microphones from Fig. 4, participants stood outside of a 1-m circumference from the dummy head in a random manner (i.e. all participants were not confined to a particular corner in the room). Beverages in glasses as well as snacks were served while the participants engaged in conversation. At any given time, there were nine male participants and six female participants present in the room. A recording of such a scenario was made for approximately 1 h, but a random sample was used in the following analysis.

As opposed to a free-field a priori RTF vector, a more suitable a priori RTF vector for the behind-the-ear hearing aid microphones was obtained from pre-measured impulse responses in the scenario as depicted in Fig. 4. The impulse responses were computed from an exponential sine-sweep measurement with the loudspeaker position at 0° (the azimuth direction directly in front of the dummy head) and 1 m so that the a priori RTF vector would be defined in accordance with a source located at 0° and 1 m from the dummy head. The initial section of these impulse responses corresponding to the direct component was extracted, with a length according to the size of the Discrete Fourier Transform (DFT) window to be used in the STFT domain processing. This direct component was then smoothed with a Tukey window and converted to the frequency domain. In each frequency bin, these smoothed frequency domain impulse responses were then scaled with respect to the smoothed frequency domain impulse response of the reference microphone. This was then used as $\mathbf{h}_a(k)$ and was kept the same for each time frame.

A scenario was firstly considered for the desired speech source located at 0° in Fig. 4, i.e. the location where the a priori RTF vector was defined. A 4s sample of the desired speech signal was mixed with a random sample of the cocktail party noise at a broadband input SNR of 0 dB. For the batch processing framework with a DFT size of 256 samples, \mathbf{R}_{yy} and \mathbf{R}_{nn} were estimated by time averaging across the entire length of the signal in the respective speech-plus-noise or noise-only frames. Using the SPP [25] from the first microphone of the left ear hearing aid, frames for which the speech was active were chosen if the resulting SPP > 0.85 . The RTF vectors, $\tilde{\mathbf{h}}$ and $\hat{\mathbf{h}}$, were computed according to the procedures described in

Sections 3.1 and 3.2. Using CVX [31, 32], the MVDR-INT from (42) was then evaluated for a range of $0 < \tilde{\epsilon} < 1.5$ and $0 < \hat{\epsilon} < 1.5$ at a frequency of 2 kHz.

Figure 5a and b display the resulting (base-10 log) values of the Lagrangian multipliers α and β respectively as a function of $\tilde{\epsilon}$ and $\hat{\epsilon}$, along with the APC and DDC bounding curves. These plots support the theoretical analysis of the space spanned by $\tilde{\epsilon}$ and $\hat{\epsilon}$ from Fig. 1. In Fig. 5a, it is clearly observed that as the value of $\tilde{\epsilon}$ exceeds the APC bounding curve, then $\alpha \rightarrow 0$ so that the APC is inactive while the DDC remains active. Similarly, in Fig. 5b, as the value of $\hat{\epsilon}$ exceeds the DDC bounding curve, then $\beta \rightarrow 0$ so that the APC remains active and the DDC is inactive. The regions where both constraints are active, and when neither are active can also be observed.

Figure 5c and d are plots of the corresponding change in SNR (Δ SNR) from the reference microphone as well as the speech distortion which was computed as follows:

$$\Delta\text{SNR} = 10 \log_{10} \left(\frac{|\underbrace{\mathbf{w}^H \mathbf{h}}_{\text{int}}|^2}{\underbrace{\mathbf{w}^H \mathbf{w}}_{\text{int}}} \right) - 10 \log_{10} \left(\frac{1}{\mathbf{e}_1^T \hat{\mathbf{R}}_{nn} \mathbf{e}_1} \right) \quad (76)$$

$$\text{SD} = \left| \underbrace{\mathbf{w}^H \mathbf{h}}_{\text{int}} - 1 \right|^2 \quad (77)$$

where the first term of the Δ SNR is the output SNR and the second term is the input SNR at the unprocessed reference microphone¹² and in this scenario $\mathbf{h} = \tilde{\mathbf{h}}$. The true value of \mathbf{h} is unknown; hence, the results of Fig. 5c and d are suggestive for the case when the true RTF vector corresponds to that of the a priori assumed RTF vector. In Fig. 5c, since $\tilde{\mathbf{w}} \rightarrow \mathbf{0}$ in the region where $\hat{\epsilon} \geq 1$ and $\tilde{\epsilon} \geq 1$, it is purposefully hatched so as to indicate that in this region an output SNR is undefined.

As expected, it can be observed that the best Δ SNR is achieved for the region where the DDC is inactive and the APC is active, with a compromise within the region where the two constraints are active. An interesting observation here is the poor Δ SNR in the region where $\tilde{\epsilon} \rightarrow 0$ and $\hat{\epsilon} \rightarrow 0$. Even though the maximum tolerable speech distortions have been specified to be quite small, in this case $\tilde{\mathbf{h}}$ and $\hat{\mathbf{h}}$ can be parallel, which can lead to redundant constraints and an ill-conditioning problem as discussed in [22]. In terms of the SD, fairly low distortions are achieved when either of the constraints are active or when both are active. As both $\tilde{\epsilon} \rightarrow 1$ and $\hat{\epsilon} \rightarrow 1$, the speech distortion increases, which is expected from (70) and (72), i.e. the SDW-MWF parameters, $\tilde{\mu}$ and $\hat{\mu}$. As $\tilde{\mu} \rightarrow \infty$, $\tilde{\epsilon} \rightarrow 1$, and as $\hat{\mu} \rightarrow \infty$, $\hat{\epsilon} \rightarrow 1$, which accounts for the increasing

¹²The numerator of this term is 1 since the first component of the RTF vector for the unprocessed microphone signals is 1.

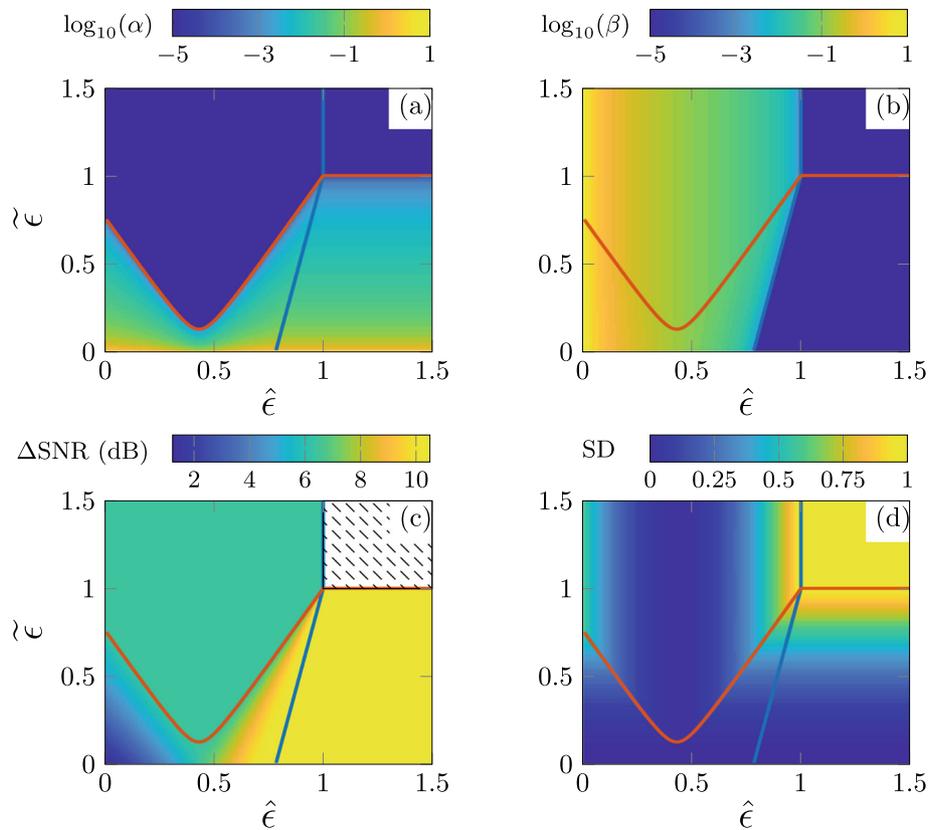


Fig. 5 Behaviour of the integrated MVDR-LMA-XM beamformer at a frequency of 2 kHz as a function of $\tilde{\epsilon}$ and $\hat{\epsilon}$ for the case when the desired speech source is at 0° , i.e. in the direction of the a priori constraint. **(a)** Lagrangian multiplier, $\log_{10}(\alpha)$. **(b)** Lagrangian multiplier, $\log_{10}(\beta)$. **(c)** Δ SNR. **(d)** speech distortion (SD). The APC and DDC bounding curves analogous to those from Fig. 1 are also shown

speech distortion from Fig. 5d. Another point to highlight in Fig. 5d is that a low speech distortion is also achieved in the region where the APC bounding curve is a minimum, regardless of the value of $\tilde{\epsilon}$. As discussed in Section 5.2, for a value of $\tilde{\epsilon} > \tilde{\epsilon}_o$ (where $\tilde{\epsilon}_o$ is the value of $\tilde{\epsilon}$ on the APC bounding curve from (60)), the achievable distortion would in fact correspond to $\tilde{\epsilon}_o$ on the APC bounding curve, which is quite low in this minimum region.

Figure 6 now displays a similar set of results, however for the case when the desired speech source was located at 60° as depicted in Fig. 4. As the a priori RTF vector was based on a speaker located at 0° , this scenario represented a mismatch between the a priori RTF vector and the true RTF vector. The same procedure as previously described was also followed to obtain the MVDR-INT filters.

Figure 6a and b display the resulting values of the (base-10 log) Lagrangian multipliers α and β respectively as a function of $\tilde{\epsilon}$ and $\hat{\epsilon}$, along with the APC and DDC bounding curves. The nature of these plots is quite similar to that of Fig. 5a and b in terms of how α and β vary with respect to the bounding curves. In comparison to Fig. 5a and b, Fig. 6a and b also highlight the fact that these bounding curves can have quite different appearances.

Figure 6c and d display the corresponding Δ SNR and SD respectively, however with $\mathbf{h} = \hat{\mathbf{h}}$ in (76), and hence, the results are suggestive for the case when the true RTF vector corresponds to that of the data-dependent RTF vector. Now it can be observed that the best Δ SNR is achieved for the region where the APC is inactive and the DDC is active, with a compromise within the region where the two constraints are active. For the SD, fairly low speech distortions are achieved for small values of $\hat{\epsilon}$ as expected. For small values of $\tilde{\epsilon}$ and large values of $\hat{\epsilon}$, i.e. toward the region where only the APC is active, it can be observed that the speech distortion increases, which is a direct result of the speech source not being in the a priori defined direction of 0° . Once again, it can also be seen that the speech distortion generally increases as both $\tilde{\epsilon} \rightarrow 1$ and $\hat{\epsilon} \rightarrow 1$.

The results of Figs. 5 and 6 provide some more insight into the behaviour of the MVDR-INT and demonstrate that in some scenarios a better performance can be achieved when either only the APC or only the DDC is active. Furthermore, it was observed that there were transition regions where a compromise could be achieved between these limits of performance when either only the

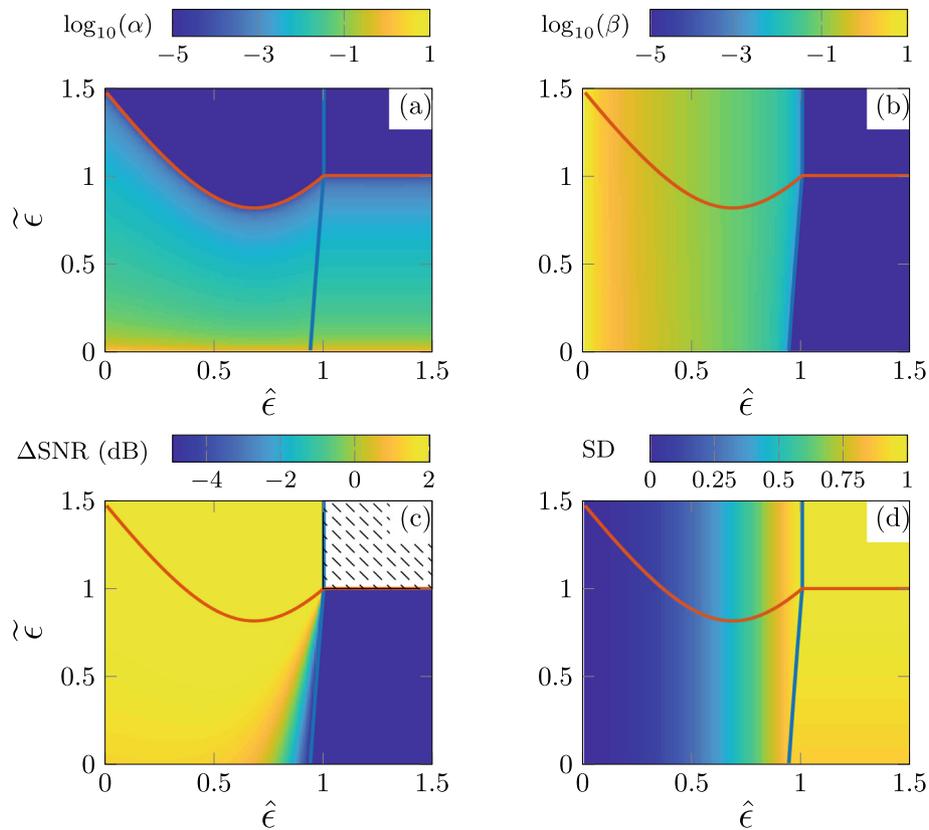


Fig. 6 Behaviour of the integrated MVDR-LMA-XM beamformer at a frequency of 2 kHz as a function of $\hat{\zeta}$ and $\hat{\epsilon}$ for the case when the source is at 60° , i.e. not in the direction of the a priori constraint. **(a)** Lagrangian multiplier, $\log_{10}(\alpha)$. **(b)** Lagrangian multiplier, $\log_{10}(\beta)$. **(c)** Δ SNR. **(d)** Speech distortion (SD). The APC and DDC bounding curves analogous to those from Fig. 1 are also shown

APC or only the DDC is active. Therefore, it suggests that tuning strategies such as those depicted in Fig. 2 would indeed be appropriate means of obtaining an optimal filter as opposed to relying on only an APC or DDC.

6.3 Performance of tuning strategies

The audio recordings as previously described for the scenario depicted in Fig. 4 are also used to observe the performance of the tuning strategies. A desired speech signal was created where the desired speech source was initially located at 0° for a duration of 5 s and then instantaneously moved to 60° for another 6 s. This was then mixed with a random sample of the cocktail party noise at a broadband input SNR of 2 dB. The same a priori RTF vector pertaining to the hearing aid microphones, $\tilde{\mathbf{h}}_{\mathbf{a}}(k)$, as previously described was used, i.e. $\tilde{\mathbf{h}}_{\mathbf{a}}(k)$, was computed for a source located at 0° and 1 m from the dummy head.

For the STFT processing, the WOLA method, with a DFT size of 256 samples, 50% overlap, a square-root hanning window, and a sampling frequency of 16 kHz were used. By using the SPP [25] computed on XM2, frames

were classified as containing speech if the $SPP > 0.8$; otherwise, the frames were classified as noise only. All RTF vector estimates were performed in frames which were classified as containing speech. All the relevant correlation matrices were also estimated using a forgetting factor corresponding to an averaging time of 300 ms. $\mathbf{R}_{\mathbf{nn}}$ was only estimated when the $SPP < 0.8$.

For the MVDR-INT, two tuning strategies were considered—(i) the trade-off between the maximum tolerable speech distortions for the APC and DDC, corresponding to Fig. 2a, which will be referred to as MVDR-INT-3a and (ii) where the maximum tolerable speech distortion for the APC is constant, but the maximum tolerable speech distortion for the DDC varies, corresponding to Fig. 2b, and which will be referred to as MVDR-INT-3b. For both tunings, $\tilde{\mu} = \hat{\mu} = 0.001$, and $\hat{\sigma}_{s_{a,1}}^2$ was computed using the method from [39] as implemented in [40] but with the noise estimation update computed as in [25]. A different setting was used for the confidence metric, $F(l)$ in (62) for each of the tunings such that for the MVDR-INT-3a, $\rho = 1$ and $\lambda_t = 5$ dB, and

for MVDR-INT-3b, $\rho = 1$ and $\lambda_t = 10$ dB, i.e. a higher thresholding was used for the MVDR-AP tuning. With all parameters assigned, the QCQP problem from (42) was solved using the gradient ascent procedure as described in Algorithm 1.

The metrics used to evaluate the following experiments were the speech intelligibility-weighted SNR [41] (SI-SNR), the short-time objective intelligibility (Δ STOI) [42], and the normalised speech-to-reverberation modulation energy ratio for cochlear implants (SRMR-CI) [43]. The SI-SNR improvement in relation to the reference microphone was calculated as:

$$\Delta \text{SI-SNR} = \sum_i I_i (\text{SNR}_{i,\text{out}} - \text{SNR}_{i,\text{in}}) \quad (78)$$

where the band importance function I_i expresses the importance of the i th one-third octave band with centre frequency, f_i^c for intelligibility, $\text{SNR}_{i,\text{in}}$ is the input SNR (dB), and $\text{SNR}_{i,\text{out}}$ is the output SNR (dB) in the i th one-third octave band. The centre frequencies, f_i^c , and the values for I_i are defined in [44]. The input SNR was computed accordingly using the unprocessed speech-only and unprocessed noise-only components (in the discrete time domain) at the reference microphone, and the output SNR from the individually processed speech-only and processed noise-only components (in the discrete time domain) resulting from the particular algorithm. For the STOI metric, the reference signal used was the unprocessed desired speech source convolved with 256 samples (i.e. same length as the DFT size) of the (pre-measured) impulse response from the desired speech signal location to the reference microphone. As the room was quite reverberant, however, a true reference signal is somewhat ambiguous to define, and hence, the non-intrusive metric, SRMR-CI, suitable for hearing instruments, in particular cochlear implants, was also used.

Figure 7 displays the performance of the various algorithms, where all the metrics have been computed in 2-s time frames with a 25% overlap. The relative improvements of the SI-SNR and the STOI metrics in relation to the reference microphone have been plotted. The metrics for XM1 and XM2 from Fig. 4 are also plotted. In order to contextualise the values of the SRMR-CI metric, an additional plot of the performance for the reference signal (that which was used for the STOI metric) is displayed. From all the metrics, as expected, the MVDR-AP performs better than the MVDR-DD in the first 5 s as the speech source was at 0° , i.e. the a priori direction. However, in the latter 6 s, when the speech source was at 60° , the MVDR-DD achieves a better performance.

With respect to the XMs, it can also be seen that the performance of XM1 decreases after 5 s as the source moves to the location of 60° , while XM2 has more of a consistent performance across the different speech locations. In

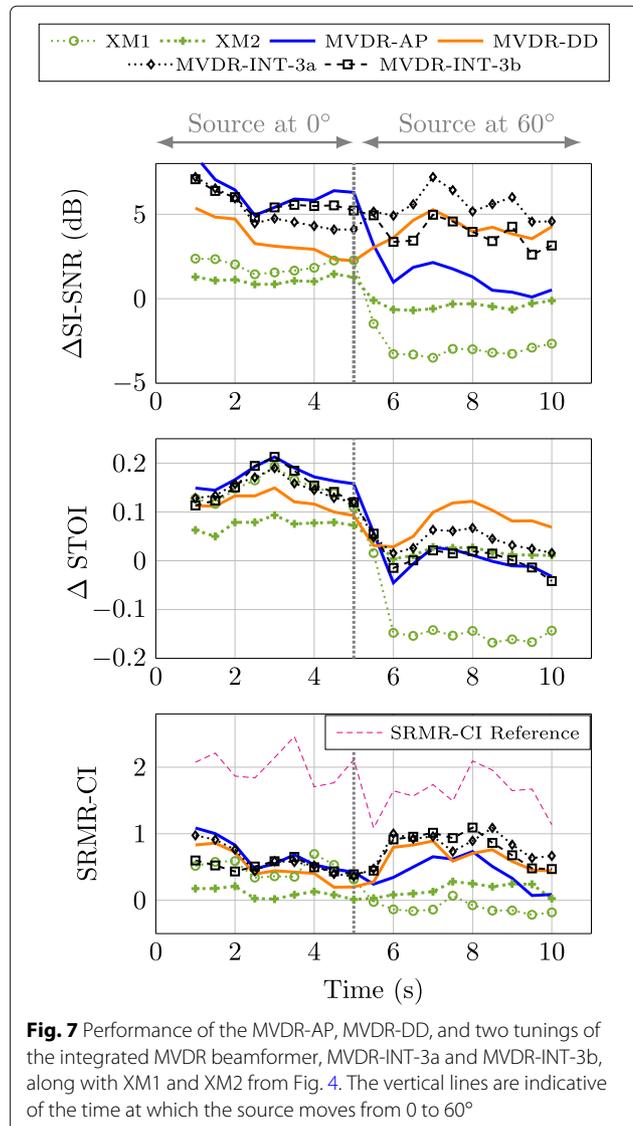


Fig. 7 Performance of the MVDR-AP, MVDR-DD, and two tunings of the integrated MVDR beamformer, MVDR-INT-3a and MVDR-INT-3b, along with XM1 and XM2 from Fig. 4. The vertical lines are indicative of the time at which the source moves from 0 to 60°

terms of the Δ SI-SNR, the performance of all of the other algorithms is better than either of the XMs, which demonstrates that simply listening to the XM only would not always immediately yield satisfactory performance.

Within the first 5 s, the MVDR-INT-3a is able to find a compromise between the MVDR-AP and MVDR-DD in terms of all metrics. In the final 6 s, although the Δ STOI is once again in between the MVDR-AP and MVDR-DD, the performance in terms of Δ SI-SNR and SRMR-CI is in fact better than either of the MVDR-AP or the MVDR-DD. This is a direct consequence of the nature of the integrated MVDR-LMA-XM beamformer as different linear combinations of the MVDR-AP and the MVDR-DD are effectively applied to different time-frequency segments, yielding a broadband SI-SNR that could be better than either the MVDR-AP or MVDR-DD.

For the MVDR-INT-3b, within the first 5 s, the performance in terms of all the metrics is closer to that of the MVDR-AP which is expected as the APC is kept active at all times. In the following 6 s, the STOI metric indicates that the speech intelligibility has not changed from that of the MVDR-AP. However, an improvement can be observed in both Δ SI-SNR and SRMR-CI metrics as some frequency bins would have also had the DDC active.

The corresponding confidence metrics across all time frames and frequencies for the MVDR-INT-3a and the MVDR-INT-3b are displayed in Fig. 8. The upper plot corresponds to the confidence metric of MVDR-INT-3a and reveals that much of the confidence has been placed on the higher frequencies, presumably because there was less noise in this region. Therefore, a smaller value of $\hat{\epsilon}$ and a larger value of $\tilde{\epsilon}$ would have been assigned to the DDC and APC respectively, i.e. the MVDR-INT-3a in this region would have tended toward the MVDR-DD. Several regions of uncertainty are also observed where the MVDR-INT-3a would then find a compromise between the MVDR-AP and the MVDR-DD. In the lower plot of Fig. 8, the confidence metric for the MVDR-INT-3b shows a much more conservative behaviour due to the larger threshold of λ_t . It is observed that there are now many regions where there is little confidence, and hence a larger value of $\hat{\epsilon}$ and a smaller value of $\tilde{\epsilon}$

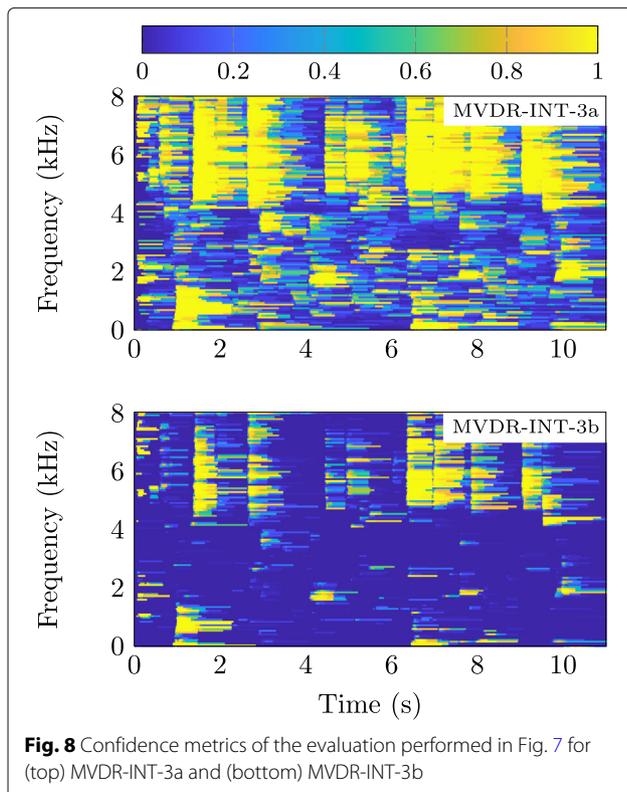
would have been assigned to the DDC and APC, respectively, i.e. the MVDR-INT-3b in these regions would have tended toward the MVDR-AP. More confidence is now only placed in the higher frequency region and there are still some regions of uncertainty so that a compromise can be achieved. The resulting audio signals from this section¹³ may also be listened to for a subjective evaluation at [45].

7 Conclusion

An integrated MVDR beamformer that merges the benefits from using an available a priori relative transfer function (RTF) vector and a data-dependent RTF vector was developed for a microphone configuration consisting of a local microphone array (LMA) and multiple external microphones (XMs). The framework has been presented in a pre-whitened-transformed (PWT) domain, which consists of an initial transformation of the microphone signals through a blocking matrix and a fixed beamformer, followed by a pre-whitening operation, facilitating convenient processing operations. In the PWT domain, procedures for obtaining an a priori RTF vector and data-dependent RTF vector have also been derived, where the a priori RTF vector is based on an a priori RTF vector pertaining to the LMA only.

With the two RTF vectors, an integrated MVDR beamformer was proposed by formulating a quadratically constrained quadratic program (QCQP), with two constraints, one of which is related to the maximum tolerable speech distortion for the imposition of the a priori RTF vector and the other related to the maximum tolerable speech distortion for the imposition of the data-dependent RTF vector. It was shown how the space spanned by each of these maximum tolerable speech distortions could be divided into four separate regions, each of which corresponded to a particular set of constraints being active or inactive. This insight then facilitated the development of a general tuning framework where the maximum tolerable speech distortions are chosen in accordance with the confidence in the accuracy of the data-dependent RTF vector. A particular set of tuning rules was also proposed, which made use of a relationship to the speech distortion weighted multi-channel Wiener filter.

The potential of the integrated MVDR beamformer was demonstrated by using audio data from an LMA of behind-the-ear hearing aid microphones and three XMs for a single desired speech source within a re-created cocktail party scenario. A narrowband evaluation confirmed the theoretical behaviour of the integrated MVDR



¹³Audio samples are also uploaded for the case when the SPP was computed on the reference microphone.

as a function of the maximum tolerable speech distortion parameters. A broadband evaluation has shown that the integrated MVDR beamformer can be tuned to yield different enhanced speech signals, which may be suitable for improving speech intelligibility despite changes in the desired speech source position and imperfectly estimated spatial correlation matrices.

Abbreviations

APC: A priori constraint; DDC: Data-dependent constraint; EVD: Eigenvalue decomposition; GEVD: Generalised eigenvalue decomposition; LMA: Local microphone array; MVDR: Minimum variance distortionless response; MVDR-DD: Fully data-dependent MVDR beamformer; MVDR-AP: MVDR beamformer based on a priori knowledge; MVDR-INT: Integrated MVDR beamformer; MWF: Multi-channel Wiener filter; PWT: Pre-whitened-transformed; QCQP: Quadratically constrained quadratic program; RTF: Relative transfer function; SNR: Signal-to-noise ratio; SI-SNR: Speech intelligibility-weighted SNR; SPP: Speech presence probability; SRMR-CI: Speech-to-reverberation modulation energy ratio for cochlear implants; STOI: Short-time objective intelligibility; XM: External microphone; WOLA: Weighted Overlap and Add

Acknowledgements

This research work was carried out at the ESAT Laboratory of KU Leuven, in the frame of IWT O&O Project nr. 150432 'Advances in Auditory Implants: Signal Processing and Clinical Aspects', KU Leuven Impulsfonds IMP/14/037, KU Leuven C2-16-00449 'Distributed Digital Signal Processing for Ad-hoc Wireless Local Area Audio Networking', and KU Leuven Internal Funds VES/16/032. The research leading to these results has also 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.

Authors' contributions

RA, TvW, and MM conceptualised and analysed the QCQP framework and tuning strategy. RA drafted the manuscript, implemented the algorithms in software, and conducted the experiments. All authors have interpreted the results and reviewed the final manuscript. The authors read and approved the final manuscript.

Availability of data and materials

The microphone data analysed in the current study as well as audio samples of the processed signals are available at [45]. Further materials are also available from the corresponding author upon request.

Competing interests

The authors declare that they have no competing interests.

Received: 9 June 2020 Accepted: 16 December 2020

Published online: 10 February 2021

References

1. M. Brandstein, D. B. Ward, *Microphone Arrays: Signal Processing, Techniques and Applications*. (Springer, New York, 2001)
2. S. Gannot, E. Vincent, S. Markovich-Golan, A. Ozerov, A consolidated perspective on multi-microphone speech enhancement and source separation. *IEEE/ACM Trans. Audio Speech Lang. Process.* **25**(4), 692–730 (2017)
3. E. Vincent, T. Virtanen, S. Gannot, *Audio Source Separation and Speech Enhancement*. (Wiley, Chichester, West Sussex, 2018)
4. J. Szurley, A. Bertrand, B. van Dijk, M. Moonen, Binaural noise cue preservation in a binaural noise reduction system with a remote microphone signal. *IEEE/ACM Trans. Audio Speech Lang. Process.* **24**(5), 952–966 (2016)
5. N. Gößling, S. Doclo, in *2018 16th Int. Workshop on Acoustic Signal Enhancement (IWAENC)*. Relative transfer function estimation exploiting spatially separated microphones in a diffuse noise field, (Tokyo, 2018), pp. 146–150
6. N. Gößling, S. Doclo, in *Speech Communication; 13th ITG Symposium*. RTF-based binaural MVDR beamformer exploiting an external microphone in a diffuse noise field, (Oldenburg, Germany, 2018), pp. 1–5
7. N. Cvijanovic, O. Sadiq, S. Srinivasan, Speech enhancement using a remote wireless microphone. *IEEE Trans. Consum. Electron.* **59**(1), 167–174 (2013)
8. D. Yee, H. Kamkar-Parsi, R. Martin, H. Puder, A noise reduction post-filter for binaurally-linked single-microphone hearing aids utilizing a nearby external microphone. *IEEE/ACM Trans. Audio Speech Lang. Process.* **26**(1), 5–18 (2017)
9. R. Ali, G. Bernardi, T. van Waterschoot, M. Moonen, Methods of extending a generalized sidelobe canceller with external microphones. *IEEE/ACM Trans. Audio Speech Lang. Process.* **27**(9), 1349–1364 (2019)
10. A. Bertrand, M. Moonen, Robust distributed noise reduction in hearing aids with external acoustic sensor nodes. *EURASIP J. Adv. Signal Process.* **2009**, 1–14 (2009)
11. A. Hassani, Distributed signal processing algorithms for multi-task wireless acoustic sensor networks. PhD thesis, KU Leuven (2017)
12. S. Markovich-Golan, S. Gannot, I. Cohen, Distributed multiple constraints generalized sidelobe canceler for fully connected wireless acoustic sensor networks. *IEEE Trans. Audio Speech Lang. Process.* **21**(2), 343–356 (2013)
13. J. Capon, High-resolution frequency-wavenumber spectrum analysis. *Proc. IEEE.* **57**(8), 1408–1418 (1969)
14. H. L. Van Trees, *Optimum Array Processing*. (Wiley, Hoboken, 2001)
15. S. Gannot, D. Burshtein, E. Weinstein, Signal enhancement using beamforming and nonstationarity with applications to speech. *IEEE Trans. Signal Process.* **49**(8), 1614–1626 (2001)
16. J. E. Greenberg, P. M. Zurek, Evaluation of an adaptive beamforming method for hearing aids. *J. Acoust. Soc. Amer.* **91**(3), 1662–1676 (1992)
17. J. M. Kates, M. R. Weiss, A comparison of hearing-aid array-processing techniques. *J. Acoust. Soc. Amer.* **99**(5), 3138–3148 (1996)
18. M. Kompis, N. Dillier, Performance of an adaptive beamforming noise reduction scheme for hearing aid applications. I. Prediction of the signal-to-noise-ratio improvement. *J. Acoust. Soc. Amer.* **109**(3), 1123–1133 (2001)
19. A. Spriet, L. Van Deun, K. Eftaxiadis, J. Laneau, M. Moonen, B. van Dijk, A. van Wieringen, J. Wouters, Speech understanding in background noise with the two-microphone adaptive beamformer BEAM in the Nucleus Freedom cochlear implant system. *Ear Hear.* **28**(1), 62–72 (2007)
20. I. Cohen, Relative transfer function identification using speech signals. *IEEE Trans. Speech Audio Process.* **12**(5), 451–459 (2004)
21. S. Markovich-Golan, S. Gannot, in *2015 IEEE Int. Conf. Acoust., Speech and Signal Process. (ICASSP)*. Performance analysis of the covariance subtraction method for relative transfer function estimation and comparison to the covariance whitening method, (Brisbane, 2015), pp. 544–548
22. R. Ali, T. Van Waterschoot, M. Moonen, Integration of a priori and estimated constraints into an MVDR beamformer for speech enhancement. *IEEE/ACM Trans. Audio Speech Lang. Process.* **27**(12), 2288–2300 (2019)
23. L. Griffiths, C. Jim, An alternative approach to linearly constrained adaptive beamforming. *IEEE Trans. Antennas Propag.* **30**(1), 27–34 (1982)
24. S. Van Gerven, F. Xie, in *Proc. EUROSPEECH, vol. 3, Ródos*. A comparative study of speech detection methods, (Greece, 1997), pp. 1095–1098
25. T. Gerkmann, R. C. Hendriks, in *Proc. 2011 IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA '11)*. Noise power estimation based on the probability of speech presence, (2011), pp. 145–148
26. I. Markovsky, *Low Rank Approximation: Algorithms, Implementation, Applications*. (Springer, Heidelberg, 2012)
27. S. Boyd, L. Vandenberghe, *Convex Optimization*. (Cambridge University Press, New York, 2004)
28. W. C. Liao, Z. Q. Luo, I. Merks, T. Zhang, in *Proc. 2015 IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA '15)*. An effective low complexity binaural beamforming algorithm for hearing aids (IEEE, New Paltz, 2015), pp. 1–5
29. W. C. Liao, M. Hong, I. Merks, T. Zhang, Z. Q. Luo, in *2015 IEEE Int. Conf. Acoust., Speech and Signal Process. (ICASSP)*. Incorporating spatial information in binaural beamforming for noise suppression in hearing aids, (Brisbane, QLD, 2015), pp. 5733–5737
30. M. Souden, J. Benesty, S. Affes, On optimal frequency-domain multichannel linear filtering for noise reduction. *IEEE Trans. Audio Speech Lang. Process.* **18**(2), 260–276 (2010)

31. M. Grant, S. Boyd, CVX: Matlab Software for Disciplined Convex Programming, version 2.1 (2014). <http://cvxr.com/cvx>. Accessed May 2020
32. M. Grant, S. Boyd, ed. by V. Blondel, S. Boyd, and H. Kimura. Recent Advances in Learning and Control. Lecture Notes in Control and Information Sciences (Springer, Springer-Verlag London, 2008), pp. 95–110
33. S. Boyd, N. Parikh, E. Chu, B. Peleato, J. Eckstein, Distributed optimization and statistical learning via the alternating direction method of multipliers. *Found. Trends Mach. Learn.* **3**(1), 1–122 (2011)
34. J. Nocedal, S. J. Wright, *Numerical Optimization, 2nd edn.* (Springer, New York, 2006)
35. A. Spriet, M. Moonen, J. Wouters, Spatially pre-processed speech distortion weighted multi-channel Wiener filtering for noise reduction. *Signal Process.* **84**(12), 2367–2387 (2004)
36. S. Doclo, S. Gannot, M. Moonen, A. Spriet, *Handbook on Array Processing and Sensor Networks.* (Wiley, Hoboken, 2010), pp. 269–302. Chap. 10: acoustic beamforming for hearing aid applications
37. R. Crochiere, A weighted overlap-add method of short-time Fourier analysis/synthesis. *IEEE Trans. Acoust. Speech Signal Process.* **28**(1), 99–102 (1980)
38. C. Veaux, J. Yamagishi, K. MacDonald, CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit (2016). <http://homepages.inf.ed.ac.uk/jyamagis/page3/page58/page58.html>. Accessed Dec 2019
39. Y. Ephraim, D. Malah, Speech enhancement using a minimum mean-square error log-spectral amplitude estimator. *IEEE Trans. Acoust. Speech Signal Process.* **33**(2), 443–445 (1985)
40. M. Brookes, et al., *Voicebox: speech processing toolbox for Matlab.* (Imperial College, London, 1997). <http://www.ee.ic.ac.uk/hp/staff/dmb/voicebox/voicebox>
41. J. E. Greenberg, P. M. Peterson, P. M. Zurek, Intelligibility-weighted measures of speech-to-interference ratio and speech system performance. *J. Acoust. Soc. Amer.* **94**(5), 3009–3010 (1993)
42. C. H. Taal, R. C. Hendriks, R. Heusdens, J. Jensen, An algorithm for intelligibility prediction of time – frequency weighted noisy speech. *IEEE Trans. Audio Speech Lang. Process.* **19**(7), 2125–2136 (2011)
43. J. F. Santos, T. H. Falk, Updating the SRMR-CI metric for improved intelligibility prediction for cochlear implant users. *IEEE/ACM Trans. Audio Speech Lang. Process.* **22**(12), 2197–2206 (2014)
44. American National Standards Institute, *American National Standard Methods for calculation of the speech intelligibility index.* (Acoustical Society of America, 1997). <https://webstore.ansi.org/standards/asa/ansiasas31997r2017>. Accessed 6 June 1997
45. R. Ali (2020). ftp://ftp.esat.kuleuven.be/pub/SISTA/rali/Reports/public_data_mvdrprint. Accessed May 2020

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Submit your manuscript to a SpringerOpen[®] journal and benefit from:

- Convenient online submission
- Rigorous peer review
- Open access: articles freely available online
- High visibility within the field
- Retaining the copyright to your article

Submit your next manuscript at ► [springeropen.com](https://www.springeropen.com)
