

# Comparative Analysis of Generalized Sidelobe Cancellation and Multi-Channel Linear Prediction for Speech Dereverberation and Noise Reduction

Thomas Dietzen, Ann Spriet, *Senior Member*, Wouter Tirry, Simon Doclo, *Senior Member*, Marc Moonen, *Fellow*, and Toon van Waterschoot, *Member*

**Abstract**—For blind speech dereverberation, two frameworks are commonly used: on the one hand, the multi-channel linear prediction (MCLP) framework, and on the other hand, data-dependent beamforming, e.g., the generalized sidelobe canceler (GSC) framework. The MCLP framework is designed to perform deconvolution and hence has gained increased prominence in blind speech dereverberation. The GSC framework is commonly used for noise reduction, but may be applied for dereverberation as well. In previous work we have shown that for the noiseless case, MCLP and the GSC yield in theory mathematically equivalent results in terms of dereverberation. In this paper, we assume additional coherent- as well as incoherent-noise components and formally analyze and compare both frameworks in terms of dereverberation and noise reduction performance. Both the theoretical analysis and time domain simulation results demonstrate that unlike the GSC, MCLP expectably shows limited performance in terms of noise reduction, while both perform equally well in terms of dereverberation, provided that the GSC blocking matrix achieves complete blocking of the early reverberant-speech component and sufficiently many microphones are available. In case of incomplete blocking, however, the GSC performs inferior to MCLP in terms of dereverberation, as shown in short-time Fourier transform (STFT) domain simulations.

**Index Terms**—Multi-channel linear prediction, data-dependent beamforming, dereverberation, noise reduction.

## I. INTRODUCTION

IT is well known that reverberation, caused by reflections against room boundaries and objects, and background noise may have a deteriorating effect on the quality and intelligibility of a speech signal recorded by a microphone [1]. Speech dereverberation accompanied by noise reduction is therefore needed in many applications ranging from hands-free mobile telephony to distant automatic speech recognition.

Dereverberation approaches based on multiple microphones take advantage of spatial diversity and, according to the multiple input/output inverse theorem (MINT) [2], theoretically allow complete inversion of the (presumed time-invariant) room impulse responses (RIRs) between the speech source and the microphone array, provided that the corresponding transfer functions do not share common zeros. In practical applications

however, the RIRs are unknown – and since MINT is very sensitive to RIR estimation errors [3], which are unavoidable in practice, especially in noisy environments [4], [5], explicit inversion is not favorable. In recent years instead, assuming *no or limited prior knowledge* on the RIRs, multi-channel linear prediction (MCLP) [6]–[17], beamforming [18]–[25] and combinations thereof [26]–[28] have been most commonly and successfully used for (blind) speech dereverberation, while partly including noise reduction [17], [21], [22], [25]–[28]. In the following, we briefly review these approaches.

The MCLP framework is designed to perform deconvolution, and is hence suited for dereverberation, while noise reduction is *not* targeted. It operates blindly on the microphone signals, i.e. does not require any prior knowledge on the RIRs. A block diagram of MCLP is shown in Fig. 1. The framework relies on the premise that the reverberant component to be canceled can be modeled as a filtered version of the delayed microphone signals, i.e. as a linear prediction component. The prediction delay is a design parameter defining the number of early reflections to be maintained. The sole task in MCLP therefore consists in estimating the multi-channel prediction filter from the microphone signals. When the prediction filter is of sufficient order, MCLP is theoretically able to completely equalize the RIRs [7]. Nowadays, MCLP is commonly implemented in frequency sub-bands using the short-time Fourier transform (STFT) [8]–[17], [26]–[28]. Incorporating the power spectral density (PSD) of the speech-source signal in the cost function has been shown to be beneficial, as, e.g., in the weighted prediction error (WPE) method [10], [11], where the speech-source signal is modeled as time-varying Gaussian [8]–[11] or using sparse priors [13]. Adaptive approaches based on recursive least squares [12], [15] and the Kalman filter [14], [16], [17], [28] have been proposed. In [17], given *noisy* microphone signals, the reverberant-speech component and the prediction-filter coefficients are estimated in an alternating fashion. To reduce noise after dereverberation, it has been proposed to cascade MCLP with minimum-variance distortionless response (MVDR) beamforming [26], [27], which became a popular approach in the recent CHiME-5 challenge [29].

Beamforming is designed to perform spatial filtering, and is hence commonly used for noise reduction, but may *also* be applied for dereverberation [30]. One can distinguish between data-independent (e.g., superdirective) beamforming and data-dependent (e.g., MVDR) beamforming. Although

T. Dietzen is with KU Leuven, Dept. of Electrical Engineering (ESAT), STADIUS Center for Dynamical Systems, Signal Processing and Data Analytics, Leuven, Belgium and was previously with NXP Semiconductors Belgium NV, Leuven, Belgium. A. Spriet, and W. Tirry are with NXP Semiconductors Belgium NV. S. Doclo is with University of Oldenburg, Dept. of Medical Physics and Acoustics and the Cluster of Excellence Hearing4All, Oldenburg, Germany. M. Moonen and T. van Waterschoot are with KU Leuven, ESAT, STADIUS.

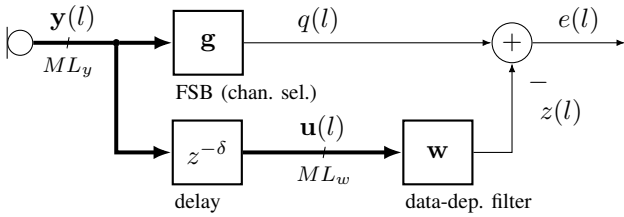


Fig. 1: The MCLP framework employing the prediction delay  $\delta$  in the data-dependent filter path.

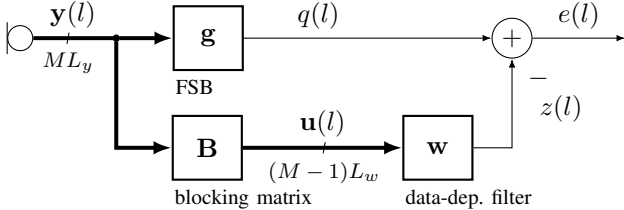


Fig. 2: The GSC framework employing the blocking matrix  $\mathbf{B}$  in the data-dependent filter path.

beamforming traditionally does *not* target channel inversion, it may be considered equivalent to MINT if the (presumably known) RIRs are incorporated in the filter design [18]. The so-called MINTFormer [20] provides a trade-off between the performance of MINT and the robustness of beamforming. In [25], a MINT-based multi-channel Wiener filter for joint dereverberation and noise reduction has been proposed. The analysis in [19] shows that for MVDR beamforming incorporating known RIRs, an inherent performance trade-off exists between dereverberation and noise reduction in case of incoherent as well as mixed coherent- and incoherent-noise fields. In this work, we are mainly concerned with the generalized sidelobe canceller (GSC) [31], [32], an implementation of the minimum-power distortionless response (MPDR) beamformer and widely employed in noise reduction. Fig. 2 depicts a block diagram of the GSC, which consists of three components: a filter-and-sum beamformer (FSB) steering a beam into the target direction, a blocking matrix blocking the speech component, and a data-dependent filter minimizing the output power and thereby suppressing residual noise components. The ability to block the speech component is essential to the GSC, as speech leakage through the blocking matrix may lead to partial speech cancellation. Employing the GSC for dereverberation, the blocking matrix should block the early (but not the late) reverberant-speech component. In [21] therefore a blocking matrix is used incorporating the relative early transfer functions (RETFs) of the speech source in order to jointly perform dereverberation and noise reduction. In the nested GSC [22], an inner GSC is employed for dereverberation and an outer GSC for noise reduction. In [28], we have proposed to integrate the GSC and MCLP in a parallel manner, and compared to the corresponding MCLP-GSC cascade, cf. also [26], [27].

A comparison of the block diagrams in Fig. 1 and 2 readily reveals the major difference between the two frameworks, which is due to their different objective. Where MCLP – designed for deconvolution – applies a simple delay to the microphone signals in the data-dependent filter path, the GSC

instead – designed for spatial filtering – applies a blocking matrix. On the one hand, regarding dereverberation, the need for a blocking matrix is certainly a drawback of the GSC as compared to MCLP, as its design requires prior knowledge. On the other hand, regarding noise reduction, the blocking matrix distinguishes the speech source from potential localized noise sources, which is not possible in MCLP. For the noiseless dereverberation task, we have shown in [24] that the MCLP and GSC framework theoretically lead to the mathematically equivalent results for stationary source signals. In practice, additional noise may always be present. In this paper therefore, using pre-whitened least squares (LS) filter estimates, we formally analyze and compare the behavior of both frameworks in case of noise, both in terms of dereverberation and noise reduction. The main intention is to provide a better understanding of the theoretical performance limitations of both frameworks depending on a number of boundary conditions, such as noise levels, filter length and number of microphones, which cannot be done by naive comparison. In our theoretical analysis, we assume *complete blocking* of the early reverberant-speech component in the GSC blocking matrix, which requires prior knowledge of the early part of the speech-source RIRs or the RETFs. We derive that if the number of microphones is sufficiently large, the GSC theoretically achieves complete coherent-noise cancellation if incoherent noise is absent, while MCLP cancels the late coherent-noise components only, as expected by design. Further, in case of *complete blocking*, the GSC performs equally well as MCLP in terms of dereverberation; theoretically achieving complete reverberation cancellation if incoherent noise is absent. These theoretical findings are confirmed by time domain simulations. In addition, in case of *incomplete blocking*, based on STFT domain simulations using estimated RETFs, we show that the GSC instead performs inferior to MCLP in terms of dereverberation.

In Sec. II, the signal model for both frameworks is presented. In Sec. III, the filter estimation is discussed. Sec. IV and V proceed with the performance analysis of the MCLP and the GSC framework, respectively. A comparative summary of the two frameworks is presented in Sec. VI, followed by simulation results in Sec. VII.

## II. SIGNAL MODEL

In this section, we define the signal model for both MCLP and the GSC. For simplicity, we employ the same notation for those signals and filters that correspond in both frameworks, cf. Fig. 1 and 2. As outlined before, the major difference between both consists in the use of a prediction delay  $\delta$  in MCLP (cf. Fig. 1) and a blocking matrix  $\mathbf{B}$  in the GSC (cf. Fig. 2) in the data-dependent filter path. In addition, the GSC speech reference is typically created by applying an FSB, whereas in MCLP a particular microphone signal is traditionally selected [6]–[16], [26], [27]. Both cases are covered generically by the filter  $\mathbf{g}$  (cf. Fig. 1 and 2). The signal model equivalently applies in the time domain and the STFT domain, where  $l$  respectively denotes the time or frame index. In case of the STFT domain, throughout the

paper, the frequency sub-band index is omitted as we treat all frequency sub-bands independently. Subsequently, vectors are denoted by lower-case boldface letters, matrices by upper-case boldface letters,  $\mathbf{I}^{L \times L}$  and  $\mathbf{0}^{L_1 \times L_2}$  denote identity and zero matrices with the (optional) superscript indicating their dimensions,  $\mathbf{A}^*$ ,  $\mathbf{A}^T$ ,  $\mathbf{A}^H$ ,  $\mathbf{A}^+$ , and  $\mathbb{E}[\mathbf{A}]$  denote the complex conjugate, the transpose, the complex conjugate transpose, the pseudoinverse and the expected value of a matrix  $\mathbf{A}$ ,  $\text{blkdiag}[\mathbf{A}_1, \dots, \mathbf{A}_N]$  constructs a block-diagonal matrix from its arguments, and  $\text{tplz}[\mathbf{a}, L]$  creates a Toeplitz matrix of  $L$  columns with the first column defined by the vector  $(\mathbf{a}^T \ \mathbf{0}^{1 \times (L-1)})^T$ .

The acoustic scenario is presented in Sec. II-A, while in Sec. II-B the speech reference signal and its individual components are defined. In Sec. II-C and Sec. II-D, the data-dependent filter input signal is discussed for MCLP and the GSC, respectively. In Sec. II-E, the filter output and the enhanced signal are generically defined.

### A. Acoustic Scenario

We assume an acoustic scenario comprising one speech source emitting the signal  $s_1(l)$ , and  $N - 1$  localized noise sources emitting the signals  $s_n(l)$ ,  $n = 2 \dots N$ , in a reverberant environment with  $M$  microphones. The  $m^{\text{th}}$  microphone signal  $y_m(l)$ ,  $m = 1 \dots M$ , consists of the reverberant-speech component, reverberant-noise components, referred to as coherent-noise components hereafter, as well as an incoherent-noise component (originating from spatially uncorrelated noise, e.g., sensor noise), i.e.

$$y_m(l) = \sum_{n=1}^N \underbrace{\sum_{k=0}^{L_h-1} h_{n,m}^*(k) s_n(l-k)}_{x_{n,m}(l)} + v_m(l), \quad (1)$$

with  $h_{n,m}(k)$  denoting the time-invariant (sub-band) RIR between the  $n^{\text{th}}$  source and the  $m^{\text{th}}$  microphone of length  $L_h$  (neglecting the dead time common to all RIRs),  $k$  the tap index,  $x_{n,m}(l)$  the reverberant components (reverberant-speech and coherent-noise components), and  $v_m(l)$  the incoherent-noise component. Note that in the STFT case, the sub-band convolution model in (1) poses an approximation of the time-domain convolution [33], where the sub-band RIR length  $L_h$  is roughly  $R_{\text{STFT}}$  times smaller than the corresponding time domain RIR length, with  $R_{\text{STFT}}$  denoting the hop size in the STFT analysis [33]. We define the stacked multi-microphone vector  $\mathbf{y}(l) \in \mathbb{C}^{ML_y}$ ,

$$\mathbf{y}(l) = (\mathbf{y}_1^T(l) \ \dots \ \mathbf{y}_M^T(l))^T, \quad (2)$$

$$\mathbf{y}_m(l) = (y_m(l) \ \dots \ y_m(l - L_y + 1))^T, \quad (3)$$

with  $L_y$  the number of samples/frames per microphone. With  $\mathbf{x}_n(l)$  and  $\mathbf{v}(l)$  defined in a similar manner as in (2), we obtain

$$\mathbf{y}(l) = \sum_{n=1}^N \mathbf{x}_n(l) + \mathbf{v}(l) = \mathbf{x}(l) + \mathbf{v}(l). \quad (4)$$

With the blockwise Toeplitz matrix  $\mathbf{H} \in \mathbb{C}^{NL_s \times ML_y}$  and the stacked source-signal vector  $\mathbf{s}(l) \in \mathbb{C}^{NL_s}$  defined by

$$\mathbf{H} = \begin{pmatrix} \mathbf{H}_{1,1} & \dots & \mathbf{H}_{1,M} \\ \vdots & & \vdots \\ \mathbf{H}_{N,1} & \dots & \mathbf{H}_{N,M} \end{pmatrix} = \begin{pmatrix} \mathbf{H}_1 \\ \vdots \\ \mathbf{H}_N \end{pmatrix}, \quad (5)$$

$$\mathbf{H}_{n,m} = \text{tplz} \left[ (h_{n,m}(0) \ \dots \ h_{n,m}(L_h - 1))^T, L_y \right], \quad (6)$$

$$\mathbf{s}(l) = (\mathbf{s}_1^T(l) \ \dots \ \mathbf{s}_N^T(l))^T, \quad (7)$$

$$\mathbf{s}_n(l) = (s_n(l) \ \dots \ s_n(l - L_s + 1))^T, \quad (8)$$

$$L_s = L_h + L_y - 1, \quad (9)$$

the vector  $\mathbf{x}(l)$  can then be written as

$$\mathbf{x}(l) = \sum_{n=1}^N \mathbf{H}_n^H \mathbf{s}_n(l) = \mathbf{H}^H \mathbf{s}(l). \quad (10)$$

We assume  $\mathbf{s}_n(l)$  and  $\mathbf{v}(l)$  to be mutually independent, i.e. with the correlation matrices  $\Psi_{s_n}(l) = \mathbb{E}[\mathbf{s}_n(l) \mathbf{s}_n^H(l)]$  and  $\Psi_v(l)$  equivalently, using (4), (7), (10), we find

$$\Psi_s(l) = \text{blkdiag}[\Psi_{s_1}(l), \dots, \Psi_{s_N}(l)], \quad (11)$$

$$\Psi_x(l) = \mathbf{H}^H \Psi_s(l) \mathbf{H}, \quad (12)$$

$$\Psi_y(l) = \Psi_x(l) + \Psi_v(l). \quad (13)$$

The matrices  $\Psi_{s_n}(l)$  are assumed to be invertible, such that  $\Psi_s^{-1}(l) = \text{blkdiag}[\Psi_{s_1}^{-1}(l), \dots, \Psi_{s_N}^{-1}(l)]$ . Note that in the STFT domain, it is commonly assumed that  $\mathbb{E}[s_1(l) s_1^*(l - k)] = 0$  for  $k \neq 0$  if the STFT hop size is sufficiently large, i.e.  $\Psi_{s_1}(l)$  becomes a diagonal matrix. Similar assumptions could be made for other source signals, but are not required in our analysis.

### B. Speech Reference Signal

With the filter  $\mathbf{g} \in \mathbb{C}^{ML_y}$ , we define the speech reference signal  $q(l)$  for both frameworks as

$$\begin{aligned} q(l) &= \mathbf{g}^H \mathbf{y}(l) \\ &= \sum_{n=1}^N \underbrace{(\mathbf{H}_n \mathbf{g})^H \mathbf{s}_n(l)}_{q_{s_n}(l)} + \underbrace{\mathbf{g}^H \mathbf{v}(l)}_{q_v(l)}, \end{aligned} \quad (14)$$

where  $q_{s_n}(l)$  and  $q_v(l)$  denote the individual source components of  $q(l)$ . Defining the parameter  $d \ll L_h$  as the boundary between early and late reverberation, the reverberant-speech and coherent-noise components  $q_{s_n}(l)$  may further be decomposed into early and late components  $\dot{q}_{s_n}(l)$  and  $\ddot{q}_{s_n}(l)$ , i.e.

$$q_{s_n}(l) = \underbrace{(\dot{\mathbf{C}} \mathbf{H}_n \mathbf{g})^H \mathbf{s}_n(l)}_{\dot{q}_{s_n}(l)} + \underbrace{(\ddot{\mathbf{C}} \mathbf{H}_n \mathbf{g})^H \mathbf{s}_n(l)}_{\ddot{q}_{s_n}(l)}, \quad (15)$$

with  $\dot{\mathbf{C}} \in \mathbb{C}^{L_s \times L_s}$  and its complement  $\ddot{\mathbf{C}}$  defined as

$$\dot{\mathbf{C}} = \begin{pmatrix} \mathbf{I}^{d \times d} & \mathbf{0}^{d \times (L_s - d)} \\ \mathbf{0}^{(L_s - d) \times d} & \mathbf{0}^{(L_s - d) \times (L_s - d)} \end{pmatrix}, \quad (16)$$

$$\ddot{\mathbf{C}} = \mathbf{I}^{L_s \times L_s} - \dot{\mathbf{C}}. \quad (17)$$

For later derivations throughout Sec. IV and Sec. V, we note that  $\ddot{q}_{s_n}(l)$  in (15) may alternatively be expressed as

$$\ddot{q}_{s_n}(l) = (\ddot{\mathbf{C}}_d \mathbf{H}_n \mathbf{g})^H \mathbf{s}_n(l-d), \quad (18)$$

where  $\ddot{\mathbf{C}}_d$  is derived from  $\ddot{\mathbf{C}}$  by shifting  $d$  rows upwards, i.e.

$$\ddot{\mathbf{C}}_d = \begin{pmatrix} \mathbf{0}^{(L_s-d) \times d} & \mathbf{I}^{(L_s-d) \times (L_s-d)} \\ \mathbf{0}^{d \times d} & \mathbf{0}^{d \times (L_s-d)} \end{pmatrix}. \quad (19)$$

Based on the above definitions, the (sub-band) impulse response (IR) relating  $s_n(l)$  and  $q_{s_n}(l)$  is graphically represented in Fig. 3. The parameter  $d$  can be controlled by design choices in the MCLP and the GSC framework, as shown in Sec. II-C and Sec. II-D, respectively.

We now define the early reverberant-speech component  $\dot{q}_{s_1}(l)$  as the target component to be maintained, and the remaining late reverberant-speech component plus all noise components  $\dot{q}_{s_1}(l) + \sum_{n=2}^N q_{s_n}(l) + q_v(l)$  as the component to be canceled. Note that in a *different* acoustic scenario, e.g., with  $N$  speech sources instead of one speech source plus  $N-1$  noise sources, the target component could be defined differently, e.g., by  $\sum_{n=1}^N \dot{q}_{s_n}(l)$ .

### C. MCLP Filter Input

In the MCLP framework, the filter input signal  $\mathbf{u}(l) \in \mathbb{C}^{ML_w}$  is a delayed version of the microphone signals  $\mathbf{y}(l)$ . The prediction delay  $\delta$  is chosen as  $\delta = d$ , i.e.

$$\begin{aligned} \mathbf{u}(l) &= \mathbf{y}(l-d) \\ &= \mathbf{H}^H \mathbf{s}(l-d) + \mathbf{v}(l-d). \end{aligned} \quad (20)$$

Hence, the length  $L_y$  in (9) equals the length  $L_w$  of a single filter channel of the data-dependent filter  $\mathbf{w}$ , i.e.

$$L_y = L_w. \quad (21)$$

With (9), (21), we determine that  $\mathbf{H} \in \mathbb{C}^{NL_s \times ML_y}$  is a fat matrix if the MCLP filter length  $L_w$  satisfies the condition

$$L_w \geq \frac{N(L_h - 1)}{M - N}, \quad (22)$$

which obviously requires  $M > N$  microphones. If  $L_w$  is chosen according to (22) and the (sub-band) RIRs meet the MINT requirements (i.e. no common zeros), which is commonly assumed [2], [7], then the system is invertible and  $\mathbf{H}$  has full row rank [2]. As it is crucial for our derivations in Sec. IV-B, full row rank of  $\mathbf{H}$  is assumed in the remainder. Since our simulation results in Sec. VII support our theoretical conclusions in Sec. IV, we consider this assumption to be reasonable.

### D. GSC Filter Input

In the GSC framework, the filter input signal  $\mathbf{u}(l) \in \mathbb{C}^{(M-1)L_w}$  is constructed by applying a blocking matrix  $\mathbf{B} \in \mathbb{C}^{ML_y \times (M-1)L_w}$  to the microphone signal, i.e.

$$\begin{aligned} \mathbf{u}(l) &= \mathbf{B}^H \mathbf{y}(l) \\ &= (\mathbf{H}\mathbf{B})^H \mathbf{s}(l) + \mathbf{B}^H \mathbf{v}(l), \end{aligned} \quad (23)$$

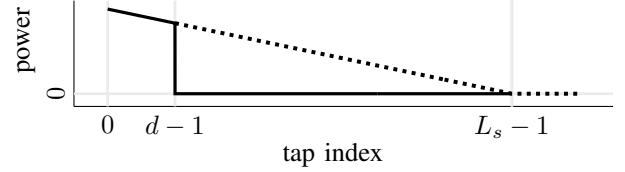


Fig. 3: Schematic of the (sub-band) IR relating  $s_n(l)$  and  $q_{s_n}(l)$ , separated in early part  $\dot{\mathbf{C}}\mathbf{H}_n\mathbf{g}$  [—] applied to  $s_n(l)$  and late part  $\ddot{\mathbf{C}}_d\mathbf{H}_n\mathbf{g}$  [····] applied to  $s_n(l-d)$ .

where  $L_w$  again describes the length of a single filter channel of the filter  $\mathbf{w}$ . Eq. (23) is the GSC counterpart to (20) for MCLP. We intend to completely block all components in  $\mathbf{x}_1(l) = \mathbf{H}_1^H \mathbf{s}_1(l)$  correlated to the target component  $\dot{q}_{s_1}(l)$  as defined in (15). A matrix  $\mathbf{B}$  satisfying<sup>1</sup> this condition may be defined in the following manner,

$$\mathbf{B} = \begin{pmatrix} -\dot{\mathbf{H}}_{1,2} & \cdots & -\dot{\mathbf{H}}_{1,M} \\ \text{blkdiag} [\dot{\mathbf{H}}_{1,1}, \dots, \dot{\mathbf{H}}_{1,1}] \end{pmatrix}, \quad (24)$$

$$\dot{\mathbf{H}}_{1,m} = \text{tplz} \left[ (h_{1,m}(0) \cdots h_{1,m}(L_b - 1))^T, L_w \right], \quad (25)$$

where  $L_b$  denotes the length of the blocking filters, such that in the GSC, we find for  $L_y$  in (9),

$$L_y = L_b + L_w - 1. \quad (26)$$

The definition in (24)–(25) ensures that all components corresponding to the first  $d \geq L_b$  taps of the speech-source (sub-band) RIRs  $h_{1,m}(k)$  are nullified, where the case  $d > L_b$  occurs if the first  $L_b$  taps of  $h_{1,m}(k)$  are succeeded by one or more zeros. The product  $\mathbf{H}\mathbf{B}$  takes the form

$$\mathbf{H}\mathbf{B} = \begin{pmatrix} \mathbf{0}^{d \times (M-1)L_w} \\ \mathbf{H}_B \end{pmatrix}, \quad (27)$$

where, using (9), (26),  $\mathbf{H}_B \in \mathbb{C}^{NL_s - d \times (M-1)L_w}$  is a fat matrix if the GSC filter length  $L_w$  satisfies the condition

$$L_w \geq \frac{N(L_h - 2) + (N - 1)d}{M - N - 1}, \quad (28)$$

which obviously requires  $M > N + 1$  microphones. If  $L_w$  is chosen according to (28) and the  $M - 1$  (sub-band) impulse responses in  $\mathbf{H}\mathbf{B}$  meet the MINT requirements, i.e. the nullity of  $(\mathbf{H}\mathbf{B})^H$  does not exceed  $d$ , then  $\mathbf{H}_B$  has full row rank according to the rank-nullity theorem [34]. As it is crucial for our derivations in Sec. V-B1, full row rank of  $\mathbf{H}_B$  is assumed in the remainder. Since our simulation results in Sec. VII support our theoretical conclusions in Sec. V, we consider this assumption to be reasonable. Comparing  $L_w$  for the GSC and MCLP in (28) and (22), respectively, we find that the GSC requires longer filters. Note however that the GSC employs one filter channel less.

<sup>1</sup>Many definitions of  $\mathbf{B}$  achieving complete blocking exist. In the STFT domain, for  $L_b = 1$ , the blocking matrix may also be defined using RETFs [21]. If the target component is instead defined as  $\sum_{n=1}^N \dot{q}_{s_n}(l)$  as in the *different* acoustic scenario mentioned in Sec. II-B, then also the definition of  $\mathbf{B}$  needs to change accordingly.

### E. Enhanced signal

For both frameworks, the filter output signal  $z(l)$  and the enhanced signal  $e(l)$  are given by

$$z(l) = \mathbf{w}^H \mathbf{u}(l), \quad (29)$$

$$e(l) = q(l) - z(l), \quad (30)$$

with  $\mathbf{u}(l)$  given by (20) in MCLP or (23) in the GSC. For MCLP,  $z(l)$  is the linear prediction of  $q(l)$ , and  $e(l)$  accordingly the linear prediction residual. The estimation of the filter  $\mathbf{w}$  is discussed in Sec. III.

## III. FILTER ESTIMATION

We now present the pre-whitened LS estimate of  $\mathbf{w}$  in Sec. III-A and discuss the choice of the pre-whitening matrix in Sec. III-B. In Sec. III-C, we present the corresponding Wiener solution, which is then used in the theoretical analysis in the subsequent Sec. IV and Sec. V.

### A. Pre-whitened LS

With  $l = 0 \dots L_{obs} - 1$  and  $L_{obs}$  denoting the number of observations used in the filter estimation, let  $\mathbf{q}_{(\cdot)} \in \mathbb{C}^{1 \times L_{obs}}$  and  $\mathbf{U}_{(\cdot)} \in \mathbb{C}^{ML_w \times L_{obs}}$  denote correspondingly stacked versions of  $q(l)$  and  $\mathbf{u}(l)$ , i.e.

$$\mathbf{q}_{(\cdot)} = (q(0) \ \cdots \ q(L_{obs} - 1)), \quad (31)$$

$$\mathbf{U}_{(\cdot)} = (\mathbf{u}(0) \ \cdots \ \mathbf{u}(L_{obs} - 1)), \quad (32)$$

and let  $\mathbf{Y}_{(\cdot)}$ ,  $\mathbf{S}_{(\cdot)}$ ,  $\mathbf{S}_{n|(\cdot)}$ , and  $\mathbf{V}_{(\cdot)}$  be defined equivalently to  $\mathbf{U}_{(\cdot)}$  in (32). Further, let  $\Omega_{(\cdot)}^{-1/2} \in \mathbb{C}^{L_{obs} \times L_{obs}}$  denote some pre-whitening matrix with  $\Omega_{(\cdot)} = \Omega_{(\cdot)}^{H/2} \Omega_{(\cdot)}^{1/2}$  to be defined explicitly in Sec. III-B. Let  $\tilde{\mathbf{U}}_{(\cdot)}$  and  $\tilde{\mathbf{q}}_{(\cdot)}$  denote correspondingly pre-whitened versions of  $\mathbf{U}_{(\cdot)}$  and  $\mathbf{q}_{(\cdot)}$ , i.e.

$$\tilde{\mathbf{q}}_{(\cdot)} = \mathbf{q}_{(\cdot)} \Omega_{(\cdot)}^{-1/2} = (\tilde{q}(0) \ \cdots \ \tilde{q}(L_{obs} - 1))^T, \quad (33)$$

$$\tilde{\mathbf{U}}_{(\cdot)} = \mathbf{U}_{(\cdot)} \Omega_{(\cdot)}^{-1/2} = (\tilde{\mathbf{u}}(0) \ \cdots \ \tilde{\mathbf{u}}(L_{obs} - 1)), \quad (34)$$

and let  $\tilde{\mathbf{Y}}_{(\cdot)}$ ,  $\tilde{\mathbf{S}}_{(\cdot)}$ ,  $\tilde{\mathbf{S}}_{n|(\cdot)}$ , and  $\tilde{\mathbf{V}}_{(\cdot)}$  as well as their respective column vectors  $\tilde{\mathbf{y}}(l)$ ,  $\tilde{\mathbf{s}}(l)$ ,  $\tilde{\mathbf{s}}_n(l)$ , and  $\tilde{\mathbf{v}}(l)$  be defined equivalently to (34). Based on these definitions, the pre-whitened data  $\tilde{q}(l)$  and  $\tilde{\mathbf{u}}(l)$  may be expressed equivalently to  $q(l)$  in (14) and  $\mathbf{u}(l)$  in (20), (23), where  $\tilde{\mathbf{y}}(l)$ ,  $\tilde{\mathbf{s}}(l)$ , and  $\tilde{\mathbf{v}}(l)$  replace  $\mathbf{y}(l)$ ,  $\mathbf{s}(l)$ , and  $\mathbf{v}(l)$ , respectively. Based on (29)–(30), (33)–(34), we generically define the LS cost function,

$$J_{LS}(\mathbf{w}) = \left\| \underbrace{\tilde{\mathbf{q}}_{(\cdot)} - \mathbf{w}^H \tilde{\mathbf{U}}_{(\cdot)}}_{\tilde{\mathbf{e}}_{(\cdot)}} \right\|_2^2, \quad (35)$$

$$\tilde{\mathbf{e}}_{(\cdot)} = (\tilde{e}(0) \ \cdots \ \tilde{e}(L_{obs} - 1))$$

leading to the the LS filter estimate  $\mathbf{w}_{LS}$ ,

$$\begin{aligned} \mathbf{w}_{LS} &= \arg \min_{\mathbf{w}} J_{LS}(\mathbf{w}) \\ &= (\tilde{\mathbf{U}}_{(\cdot)} \tilde{\mathbf{U}}_{(\cdot)}^H)^{-1} \tilde{\mathbf{U}}_{(\cdot)} \tilde{\mathbf{q}}_{(\cdot)}^H \\ &= (\mathbf{U}_{(\cdot)} \Omega_{(\cdot)}^{-1} \mathbf{U}_{(\cdot)}^H)^{-1} \mathbf{U}_{(\cdot)} \Omega_{(\cdot)}^{-1} \mathbf{q}_{(\cdot)}^H. \end{aligned} \quad (36)$$

Note that with (33)–(34),  $\mathbf{w}_{LS}$  in (36) may alternatively be written as

$$\mathbf{w}_{LS} = \left( \sum_{l=0}^{L_{obs}-1} \tilde{\mathbf{u}}(l) \tilde{\mathbf{u}}^H(l) \right)^{-1} \sum_{l=0}^{L_{obs}-1} \tilde{\mathbf{u}}(l) \tilde{q}^*(l). \quad (37)$$

### B. Choice of pre-whitening matrix

The pre-whitening matrix  $\Omega_{(\cdot)}^{-1/2}$  may be used to mitigate an estimation bias due to the speech source-signal statistics. E.g., for the two cases  $\varepsilon_{(\cdot)} \in \{\hat{\mathbf{q}}_{s_1|(\cdot)}, \mathbf{s}_{1|(\cdot)}\}$  with  $\hat{\mathbf{q}}_{s_1|(\cdot)}$  and  $\mathbf{s}_{1|(\cdot)}$  defined equivalently to (31), an unbiased estimate is achieved if

$$\Omega_{(\cdot)} = \Psi_{\varepsilon|(\cdot)}, \quad (38)$$

where  $\Psi_{\varepsilon|(\cdot)} = \mathbb{E}[\varepsilon_{(\cdot)}^H \varepsilon_{(\cdot)}]$ . Eq. (36) then corresponds to the (unbiased) generalized LS estimator of  $\mathbf{w}$  for the data model  $\mathbf{q}_{(\cdot)} = \mathbf{w}^H \mathbf{U}_{(\cdot)} + \varepsilon_{(\cdot)}$ , where  $\varepsilon_{(\cdot)}$  resembles the observation noise. In the time domain,  $\Psi_{\hat{q}_{s_1|(\cdot)}}$  and  $\Psi_{s_1|(\cdot)}$  are generally non-diagonal. The choice  $\Omega_{(\cdot)} = \Psi_{s_1|(\cdot)}$  here corresponds to the pre-whitening paradigms proposed in [6] for MCLP and [23] for the GSC. The choice  $\Omega_{(\cdot)} = \mathbf{I}$  generally yields a biased estimate, as demonstrated in the simulations in Sec. VII-B1.

In the STFT domain,  $\Psi_{\hat{q}_{s_1|(\cdot)}}$  may be modeled as a matrix with  $2d-1$  non-zero diagonals and  $\Psi_{s_1|(\cdot)}$  may be modeled as a fully diagonal matrix, cf. Sec. II-A, where the  $l^{\text{th}}$  diagonal element of  $\Psi_{s_1|(\cdot)}$  corresponds to the PSD  $\psi_{s_1}(l) = \mathbb{E}[|s_1(l)|^2]$ . In this case, with (31)–(32),  $\mathbf{w}_{LS}$  in (36) may therefore be written as

$$\mathbf{w}_{LS} = \left( \sum_{l=0}^{L_{obs}-1} \frac{\mathbf{u}(l) \mathbf{u}^H(l)}{\psi_{s_1}(l)} \right)^{-1} \sum_{l=0}^{L_{obs}-1} \frac{\mathbf{u}(l) q^*(l)}{\psi_{s_1}(l)}, \quad (39)$$

i.e. each frame  $q(l)$  and  $\mathbf{u}(l)$  is weighted by the inverse of  $\psi_{s_1}(l)$ , which, in case of MCLP, corresponds to the WPE criterion [10], [11]. Note that  $\psi_{s_1}(l)$  varies over time for non-stationary source signals.

In Sec. VII-B2, we present STFT-domain simulations for  $d = 1$  and  $\Omega_{(\cdot)} = \Psi_{\hat{q}_{s_1|(\cdot)}} \propto \Psi_{s_1|(\cdot)}$ . Herein, prior to estimating  $\mathbf{w}$  according to (36), the PSDs  $\psi_{\hat{q}_{s_1}}(l)$  on the diagonal of  $\Psi_{\hat{q}_{s_1|(\cdot)}}$  are estimated as proposed in [28], [35], i.e. by applying the generalized eigenvalue decomposition (GEVD) to the spatial correlation matrix of the microphone signals in each frame  $l$ , where it is assumed that the spatial coherence matrix of the late reverberant component may be modeled as diffuse, cf. Sec. VII-A3.

Note that in the *different* acoustic scenario mentioned in Sec. II-B with the target component defined as  $\sum_{n=1}^N \hat{q}_{s_n}(l)$  instead of  $\hat{q}_{s_1}(l)$ , in order to achieve an unbiased filter estimate, one has to change  $\Omega_{(\cdot)}$  accordingly, e.g., using  $\varepsilon_{(\cdot)} = \sum_{n=1}^N \hat{\mathbf{q}}_{s_n|(\cdot)}$  in (38).

### C. Convergence to Wiener filter solution

For the purpose of the analysis in Sec. IV and Sec. V, we assume wide-sense stationarity for the pre-whitened signals  $\tilde{\mathbf{u}}(l)$  and  $\tilde{q}(l)$ , i.e. their statistics are independent of  $l$ . Then, for  $L_{obs} \rightarrow \infty$ , the estimate  $\mathbf{w}_{LS}$  in (37) converges to the Wiener filter solution  $\mathbf{w}_{WF}$ ,

$$\mathbf{w}_{WF} = \Psi_{\tilde{\mathbf{u}}}^+ \psi_{\tilde{\mathbf{u}}\tilde{q}}, \quad (40)$$

with  $\Psi_{\tilde{\mathbf{u}}} = \mathbb{E}[\tilde{\mathbf{u}}(l) \tilde{\mathbf{u}}^H(l)]$  and  $\psi_{\tilde{\mathbf{u}}\tilde{q}} = \mathbb{E}[\tilde{\mathbf{u}}(l) \tilde{q}^*(l)]$ . Here, the inverse in (37) is replaced by the pseudoinverse, as in the GSC,  $\Psi_{\tilde{\mathbf{u}}}$  becomes rank-deficient in absence of incoherent noise and in case of complete blocking, i.e. if (27) holds, cf. Sec. V-B1 and Appendix A-1.



(a) correlation matrix  $\Psi_{\tilde{s}_n}$  (b) correlation matrix  $\Psi_{\tilde{s}_n|d}$

Fig. 4: Schematic of the correlation matrices  $\Psi_{\tilde{s}_n}$  and  $\Psi_{\tilde{s}_n|d}$  as different submatrices of a larger correlation matrix.

#### IV. MCLP ANALYSIS

For  $\mathbf{w} = \mathbf{w}_{WF}$ , we now derive the MCLP filter output signal  $z(l)$  in Sec. IV-A and then derive and discuss the enhanced signal  $e(l)$  under different noise conditions in Sec. IV-B.

##### A. MCLP Filter Output

Using (14), (20), and noting that  $\mathbb{E}[\tilde{\mathbf{y}}(l-d)\tilde{\mathbf{y}}^H(l-d)] = \mathbb{E}[\tilde{\mathbf{y}}(l)\tilde{\mathbf{y}}^H(l)]$ , the terms  $\Psi_{\tilde{u}}$  and  $\psi_{\tilde{u}\tilde{q}}$  in (40) become

$$\Psi_{\tilde{u}} = \Psi_{\tilde{y}}, \quad (41)$$

$$\psi_{\tilde{u}\tilde{q}} = \Psi_{\tilde{y}|d}\mathbf{g}, \quad (42)$$

$$\Psi_{\tilde{y}|d} = \mathbb{E}[\tilde{\mathbf{y}}(l-d)\tilde{\mathbf{y}}^H(l)]. \quad (43)$$

Inserting (20) in (29) and substituting  $\mathbf{w}$  by  $\mathbf{w}_{WF}$  in (40), we obtain for the filter output signal,

$$z(l) = (\Psi_{\tilde{y}}^+ \Psi_{\tilde{y}|d}\mathbf{g})^H \mathbf{y}(l-d). \quad (44)$$

Let the shifted correlation matrices  $\Psi_{\tilde{s}_n|d}$ ,  $\Psi_{\tilde{s}_1|d}$ , and  $\Psi_{\tilde{v}|d}$  be defined equivalently to  $\Psi_{\tilde{y}|d}$  in (43), with relations equivalent to (11)–(13). We now introduce the following relation between  $\Psi_{\tilde{s}_n|d}$  and  $\Psi_{\tilde{s}_n}$ , which is used in the subsequent derivations in Sec. IV-B. For this, note that we can interpret  $\Psi_{\tilde{s}_n|d}$  and  $\Psi_{\tilde{s}_n}$  as different submatrices of a larger correlation matrix, as shown in Fig. 4. The submatrix defining  $\Psi_{\tilde{s}_n|d}$  is shifted left by  $d$  columns as compared to the submatrix defining  $\Psi_{\tilde{s}_n}$ . Noting that the autocorrelation width of  $\tilde{s}_n(l)$  is typically much smaller than  $L_h$  in both the time and STFT domain, we assume that the autocorrelation of  $\tilde{s}_n(l)$  is zero for lags greater than  $L_s - d$ , where  $L_s - d \geq \frac{M}{M-N}(L_h - 1) - d$  and  $d \ll L_h$ , cf. (9), (21)–(22). Using (16), (19), we then express  $\Psi_{\tilde{s}_n|d}$  in terms of  $\Psi_{\tilde{s}_n}$  by

$$\Psi_{\tilde{s}_n|d} = \Psi_{\tilde{s}_n} \check{\mathbf{C}}_d + \check{\mathbf{C}}_d \Psi_{\tilde{s}_n} \dot{\mathbf{C}}. \quad (45)$$

The product  $\Psi_{\tilde{s}_n} \check{\mathbf{C}}_d$  shifts the elements in  $\Psi_{\tilde{s}_n}$  right by  $d$  columns. The product  $\check{\mathbf{C}}_d \Psi_{\tilde{s}_n} \dot{\mathbf{C}}$  replaces the resulting zero columns by the first  $d$  columns of  $\Psi_{\tilde{s}_n}$  shifted up by  $d$  rows.

##### B. MCLP Enhancement

We now analyze the behavior of MCLP considering two scenarios: absence and presence of incoherent noise.

1) *Absence of Incoherent Noise*: The absence of incoherent noise corresponds to  $\mathbf{v}(l) = \mathbf{0}$ , i.e.  $\mathbf{y}(l) = \mathbf{x}(l)$ . In this case, using (10) and relations equivalent to (11)–(13), the individual terms in (44) are equal to  $\mathbf{y}(l-d) = \mathbf{H}^H \mathbf{s}(l-d)$ ,  $\Psi_{\tilde{y}} = \mathbf{H}^H \Psi_{\tilde{s}} \mathbf{H}$ , and  $\Psi_{\tilde{y}|d} = \mathbf{H}^H \Psi_{\tilde{s}|d} \mathbf{H}$ . Inserting these in (44) and noting that  $\mathbf{H}^+ = \mathbf{H}^H (\mathbf{H} \mathbf{H}^H)^{-1}$  and hence  $\mathbf{H} \mathbf{H}^+ = \mathbf{I}$  since  $\mathbf{H}$  is assumed to have full row rank yields

$$z(l) = (\Psi_{\tilde{s}}^+ \Psi_{\tilde{s}|d} \mathbf{H} \mathbf{g})^H \mathbf{s}(l-d), \quad (46)$$

which, using (11), (45), (14)–(15), may be written as

$$z(l) = \sum_{n=1}^N \left( \dot{q}_{s_n}(l) + \vartheta_{n|d}^H \mathbf{s}_n(l-d) \right), \quad (47)$$

$$\text{with } \vartheta_{n|d} = \Psi_{\tilde{s}_n}^{-1} \check{\mathbf{C}}_d \Psi_{\tilde{s}_n} \dot{\mathbf{C}} \mathbf{H}_n \mathbf{g}. \quad (48)$$

As apparent from (47)–(48), all reverberant source components are treated *mutually independently* and *equally*. This holds as long as (22) is satisfied and  $\mathbf{H}$  has full row rank. Inserting (47) into (30) yields the MCLP output signal,

$$e(l) = \sum_{n=1}^N \underbrace{\left( \dot{q}_{s_n}(l) - \vartheta_{n|d}^H \mathbf{s}_n(l-d) \right)}_{e_{s_n}(l)}. \quad (49)$$

From (49), we observe that  $e(l)$  equals the sum of the early components  $\dot{q}_{s_n}(l)$  and a (potential) bias term  $-\vartheta_{n|d}^H \mathbf{s}_n(l-d)$  per source, with  $\vartheta_{n|d} \in \mathbb{C}^{L_s}$  and  $L_s = L_h + L_w - 1$  according to (9), (26). Therefore, as only the late components  $\dot{q}_{s_n}(l)$  are canceled, the MCLP framework suits best in the *different* acoustic scenario mentioned in Sec. II-B with the target component defined as  $\sum_{n=1}^N \dot{q}_{s_n}(l)$  instead of  $\dot{q}_{s_1}(l)$ . Combining (15) and (18), we can compare the individual components  $e_{s_n}(l)$  in (49) to

$$q_{s_n}(l) = \dot{q}_{s_n}(l) + (\check{\mathbf{C}}_d \mathbf{H}_n \mathbf{g})^H \mathbf{s}_n(l-d), \quad (50)$$

i.e. the bias term replaces the late component  $\dot{q}_{s_n}(l) = (\check{\mathbf{C}}_d \mathbf{H}_n \mathbf{g})^H \mathbf{s}_n(l-d)$ . Similarly as for the (sub-band) IR  $\mathbf{H}_n \mathbf{g}$  relating  $s_n(l)$  and  $q_{s_n}(l)$  in Fig. 3, we visualize the (sub-band) IR relating  $s_n(l)$  and  $e_{s_n}(l)$ , composed of the early part  $\dot{\mathbf{C}} \mathbf{H}_n \mathbf{g}$  and the bias part  $-\vartheta_{n|d}$ , in Fig. 5. In the following we interpret the bias term in more detail, which has also partly been done in our previous work [24]. Firstly, from  $\dot{\mathbf{C}} \mathbf{H}_n \mathbf{g}$  in (48), we observe that the bias term  $-\vartheta_{n|d}^H \mathbf{s}_n(l-d)$  depends on the first  $d$  taps of  $\mathbf{H}_n \mathbf{g}$  only, i.e. on its early part, but *not* its late part. Secondly, we note that  $\vartheta_{n|d}$  depends on the correlation matrix  $\Psi_{\tilde{s}_n}$  of the pre-whitened version  $\tilde{s}_n(l)$  of  $s_n(l)$ , cf. Sec. III-A. We can hence argue that for  $\Omega_{(\cdot)} = \Psi_{\dot{q}_{s_1}|\cdot}$  as defined in (38), with  $\dot{q}_{s_1}(l) = (\dot{\mathbf{C}} \mathbf{H}_1 \mathbf{g})^H \mathbf{s}_1(l)$ , cf. (15), the coloration of the pre-whitened speech-source signal  $\tilde{s}_1(l)$  is inverse to the filter  $\dot{\mathbf{C}} \mathbf{H}_1 \mathbf{g}$ , such that only the first element of the vector  $\Psi_{\tilde{s}_1} \dot{\mathbf{C}} \mathbf{H}_1 \mathbf{g}$  is non-zero. Similarly, we can argue that for  $\Omega_{(\cdot)} = \Psi_{s_1|\cdot}$ , the matrix  $\Psi_{\tilde{s}_1}$  becomes diagonal, such that only the first  $d$  elements of  $\Psi_{\tilde{s}_1} \dot{\mathbf{C}} \mathbf{H}_1 \mathbf{g}$  are non-zero. In both cases, with  $\check{\mathbf{C}}_d$  as in (19), we find  $\check{\mathbf{C}}_d \Psi_{\tilde{s}_1} \dot{\mathbf{C}} \mathbf{H}_1 \mathbf{g} = \mathbf{0}$ , and therefore  $\vartheta_{1|d} = \mathbf{0}$  in (48) and finally  $e_{s_1}(l) = \dot{q}_{s_1}(l)$  in (49). Hence, the estimator is indeed unbiased for  $\Omega_{(\cdot)} \in \{\Psi_{\dot{q}_{s_1}|\cdot}, \Psi_{s_1|\cdot}\}$ , as anticipated in Sec. III.

Note that the remaining early components may still be biased, i.e.  $\vartheta_{n|d} \neq \mathbf{0}$  for  $n \neq 1$ . In general, for  $\Omega_{(\cdot)} = \mathbf{I}$ , the term  $\vartheta_{n|d}^H \mathbf{s}_n(l-d)$  in (49) represents a (delayed) linear prediction component of  $\dot{q}_{s_n}(l)$ , i.e. the output signal component  $e_{s_n}(l)$  may be understood as a (partially) whitened version of  $\dot{q}_{s_n}(l)$ . This effect is also known as excessive whitening [7].

2) *Presence of Incoherent Noise*: If additional incoherent noise  $\mathbf{v}(l) \neq \mathbf{0}$  is present, the pseudoinverse of the sum  $\Psi_{\tilde{y}} = \Psi_{\tilde{x}} + \Psi_{\tilde{v}}$  in the filter  $\Psi_{\tilde{y}}^+ \Psi_{\tilde{y}|d} \mathbf{g}$  in (44) cannot be

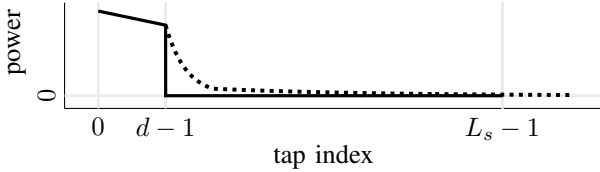


Fig. 5: Schematic of the (sub-band) IR relating  $s_n(l)$  and  $e_{s_n}(l)$ , separated in early part  $\dot{\mathbf{C}}\mathbf{H}_n\mathbf{g}$  [—] applied to  $s_n(l)$  and bias part  $-\vartheta_{n|d}$  [·····] applied to  $s_n(l-d)$ .

decomposed into its individual components, such that further simplification of (44) is not possible. In this more general case, MCLP cancels the linear prediction of the sum of  $\mathbf{g}^H\mathbf{x}(l)$  and  $\mathbf{g}^H\mathbf{v}(l)$ . Noting that the incoherent noise acts as  $M$  additional independent sources, we find that the condition (22) for complete linear prediction in the MCLP framework, where  $M$  is required to exceed the number of independent sources  $N$ , cannot be fulfilled, resulting in decreased performance.

## V. GSC ANALYSIS

Similarly to Sec. IV, for  $\mathbf{w} = \mathbf{w}_{\text{WF}}$ , we now derive the GSC filter output signal  $z(l)$  in Sec. V-A and then derive and discuss the enhanced signal  $e(l)$  under different noise conditions in Sec. V-B.

### A. GSC Filter Output

Following a derivation similar to Sec. IV-A, using (14) and (23),  $\Psi_{\tilde{u}}$  and  $\psi_{\tilde{u}\tilde{q}}$  in (40) can be written as

$$\Psi_{\tilde{u}} = \mathbf{B}^H \Psi_{\tilde{y}} \mathbf{B}, \quad (51)$$

$$\psi_{\tilde{u}\tilde{q}} = \mathbf{B}^H \Psi_{\tilde{y}} \mathbf{g}. \quad (52)$$

Inserting (23) in (29) and substituting  $\mathbf{w}$  by  $\mathbf{w}_{\text{WF}}$  in (40), we obtain for the filter output signal,

$$z(l) = (\mathbf{B}(\mathbf{B}^H \Psi_{\tilde{y}} \mathbf{B}) + \mathbf{B}^H \Psi_{\tilde{y}} \mathbf{g})^H \mathbf{y}(l). \quad (53)$$

### B. GSC Enhancement

Similarly to Sec. IV-B, we now analyze the behavior of the GSC, again considering two scenarios: absence and presence of incoherent noise.

1) *Absence of Incoherent Noise:* For the case  $\mathbf{v}(l) = \mathbf{0}$ , i.e.  $\mathbf{y}(l) = \mathbf{x}(l)$ , using (10) and relations equivalent to (11)–(13), the individual terms in (53) are equal to  $\mathbf{y}(l) = \mathbf{H}^H \mathbf{s}(l)$  and  $\Psi_{\tilde{y}} = \mathbf{H}^H \Psi_{\tilde{s}} \mathbf{H}$ . Inserting in (53) yields

$$z(l) = \left( \mathbf{H}\mathbf{B}((\mathbf{H}\mathbf{B})^H \Psi_{\tilde{s}} \mathbf{H}\mathbf{B}) + (\mathbf{H}\mathbf{B})^H \Psi_{\tilde{s}} \mathbf{H}\mathbf{g} \right)^H \mathbf{s}(l). \quad (54)$$

In Appendix A-1, assuming complete blocking such that (27) holds, it is shown that (54) can be reformulated as

$$z(l) = \ddot{q}_{s_1}(l) + \vartheta_{1|d}^H \mathbf{s}_1(l-d) + \sum_{n=2}^N q_{s_n}(l), \quad (55)$$

$$\text{with } \vartheta_{1|d} = (\ddot{\mathbf{C}}_d \Psi_{\tilde{s}_1} \ddot{\mathbf{C}}_d^H)^+ \ddot{\mathbf{C}}_d \Psi_{\tilde{s}_1} \dot{\mathbf{C}} \mathbf{H}_0 \mathbf{g}. \quad (56)$$

Inserting (55) into (30) yields the GSC output signal,

$$e(l) = \dot{q}_{s_1}(l) - \vartheta_{1|d}^H \mathbf{s}_1(l-d). \quad (57)$$

Eqs. (55)–(57) form the GSC counterpart to (47)–(49) for MCLP. From (57), we observe that  $e(l)$  consists of two terms: the target component  $\dot{q}_{s_1}(l)$  and a bias term  $-\vartheta_{1|d}^H \mathbf{s}_1(l-d)$ . This implies that not only the late reverberant-speech component, but also the coherent-noise components are completely canceled in the GSC, which is in contrast to MCLP, where only the late, but not the early coherent-noise components could be canceled. In Appendix A-2, it is shown that  $\vartheta_{1|d}$  in (56) for the GSC is indeed equal to  $\vartheta_{1|d}$  in (48) for MCLP. Hence, the discussion on the bias term in MCLP in Sec. IV-B1 similarly applies to the GSC, implying that for  $\Omega_{(\cdot,\cdot)} \in \{\Psi_{\dot{q}_{s_1}|\cdot}, \Psi_{\mathbf{s}_1|\cdot}\}$  in (36), we find  $\vartheta_{1|d} = \mathbf{0}$  and  $e(l) = \dot{q}_{s_1}(l)$ , i.e. we achieve complete and unbiased cancellation. Note that if the target component was defined as  $\sum_{n=1}^N \dot{q}_{s_n}(l)$  instead of  $\dot{q}_{s_1}(l)$  as in the *different* acoustic scenario mentioned in Sec. II-B and B was changed accordingly, for the same  $\Omega_{(\cdot,\cdot)}$ , the GSC would yield the same result as MCLP. Also note that the above conclusions hold for complete blocking, i.e. if (27) is satisfied. For *incomplete blocking*, partial speech cancellation may appear. In Sec. VII, we simulate both cases.

2) *Presence of Incoherent Noise:* Similarly to the MCLP framework, if additional incoherent noise  $\mathbf{v}(l) \neq \mathbf{0}$  is present, a simplification of (53) is not possible. We may therefore apply the same reasoning as in Sec. IV-B2. Noting that the incoherent noise acts as  $M$  additional independent noise sources, we find that the condition for complete cancellation in the GSC framework (28), where  $M-1$  is required to exceed the number of independent sources  $N$ , cannot be fulfilled, resulting in decreased performance. These conclusions are compliant with the analysis in [19], which demonstrates that in MVDR beamforming, there is an inherent trade-off between dereverberation and noise reduction for incoherent and mixed-coherent-plus-incoherent noise fields.

## VI. COMPARATIVE SUMMARY

Table I summarizes the theoretical findings from Sec. IV and Sec. V. For a given pre-whitening matrix  $\Omega_{(\cdot,\cdot)}^{-1/2}$ , MCLP does not further distinguish between the speech source and the localized noise sources, while the GSC does so by means of the spatial pre-processing in the blocking matrix. MCLP hence treats all source components  $q_{s_n}(l)$  the same, suppressing the late components  $\ddot{q}_{s_n}(l)$ , but none of the early components  $\dot{q}_{s_n}(l)$ . By contrast, the GSC suppresses all coherent-noise components  $q_{s_n}(l)$  on the one hand, and the late reverberant-speech component  $\ddot{q}_{s_1}(l)$  on the other hand, provided that the blocking matrix  $\mathbf{B}$  achieves *complete blocking* of the early reverberant-speech component. In both frameworks, an unbiased speech component with  $\vartheta_{1|d} = \mathbf{0}$  may be obtained by pre-whitening  $q(l)$  and  $\mathbf{u}(l)$  accordingly. The presence of incoherent noise decreases the performance.

The spatial pre-processing of the GSC naturally comes at a cost. The blocking matrix requires spatial information, which needs to be acquired in practice. Further, as to be demonstrated in simulations, cf. Sec. VII-B, in case of *incomplete blocking*, the GSC performs inferior in terms of dereverberation as compared to MCLP. Compared to MCLP, the minimum number

Property	MCLP framework	GSC framework
spatial knowledge	not required	required in $\mathbf{B}$ , cf. (24)–(25)
filter input signal	$\mathbf{u}(l) = \mathbf{y}(l-d)$	$\mathbf{u}(l) = \mathbf{B}^H \mathbf{y}(l)$
filter length required for complete* cancellation	$L_w \geq \frac{N(L_h - 1)}{M - N}$ , requires $M > N$	$L_w \geq \frac{N(L_h - 2) + (N - 1)d}{M - N - 1}$ , requires $M > N + 1$
output signal*	$e(l) = \sum_{n=1}^N \hat{q}_{s_n}(l) - \boldsymbol{\vartheta}_{n d}^H \mathbf{s}_n(l-d)$	$e(l) = \hat{q}_{s_1}(l) - \boldsymbol{\vartheta}_{1 d}^H \mathbf{s}_0(l-d)$

\* if incoherent noise absent (reduced performance otherwise)

TABLE I: Comparative summary of the MCLP framework versus the GSC framework.

of required microphones is increased by one, as the blocking matrix creates  $M - 1$  independent output signals only from  $M$  input signals. In the GSC, the number of filter channels is accordingly decreased by one, where a higher filter length  $L_w$  is required per channel.

## VII. SIMULATIONS

In this section, we present simulation results comparing MCLP and the GSC in terms of dereverberation and noise reduction performance. The simulation setup is described in VII-A, and the results are discussed in VII-B.

### A. Simulation Setup

In order to confirm the theory in Sec. IV–V and to further assess the practical relevance we respectively perform simulations for time domain and STFT domain implementations. The reasons for this are as follows: in the time domain, using oracle knowledge on the early RIRs, complete blocking can be simulated for the GSC, cf. Sec. VII-A3a. Further, unweighted global power-ratio measures can be well defined and evaluated, cf. Sec. VII-A4a. Therefore, in order to confirm the theory, we perform simulations on the time domain implementation, employing an ideal blocking matrix in the GSC yielding *complete blocking*, and evaluate the performance using unweighted global power-ratio measures. In the STFT domain, complete blocking cannot be simulated, since the sub-band convolution model in (1) poses an approximation of the time-domain convolution only [33]. Instead of using oracle knowledge in the blocking matrix, we estimate the RETFs from the microphone signals, such that the GSC performance also depends on the estimation quality of the RETFs, cf. Sec. VII-A3b. As, due to incomplete blocking, power-ratio measures equivalent to those used in the time domain cannot be well defined, and as unweighted global power-ratio measures are further known to relate comparably poorly to the perceived speech quality, we instead use perceptually motivated frequency-weighted segmental power-ratio measures [36]. Therefore, in order to address the practical relevance, we perform simulations on the STFT domain implementation, employing an estimated blocking matrix in the GSC yielding *incomplete blocking*, and evaluate the performance using weighted segmental power-ratio measures.

1) *Acoustic Scenario*: In order to generate multi-channel RIRs, the randomized image method [37] is used at a sampling frequency of 16 kHz, whereby the image sources are randomly displaced within a sphere of 8 cm. Multi-channel RIRs are generated using the randomized image method [37] at a sampling frequency of 16 kHz, with the image sources randomly displaced within a sphere of 8 cm. A fractional delay low-pass filter with a relative cut-off frequency of 0.9 and a length of 11 taps is applied, such that the energy of each acoustic wave, i.e. of the direct component and each reflection, is spread over 11 samples. The room dimensions are  $5 \times 4 \times 3$  m, the reverberation time is 0.5 s. The room impulse responses are truncated after 8000 taps. A linear array of 8 microphones with inter-microphone distances of (4, 4, 4, 8, 4, 4, 4) cm is used. The simulations comprise one speech source and one localized noise source. In total, 128 scenarios are generated. In each scenario, the position and orientation of the microphone array is randomized. The speech source is located at a random position in broadside direction at 2 m distance to the center of the microphone array (i.e. on a circle around its axis). The position of the localized noise source is randomized, with the constraint that the distance to the center of the microphone array is at least 1 m and the angle between the localized noise source and the speech source, seen from the center of the microphone array, is at least  $15^\circ$ .

2) *Source Signals*: We define two different source signal settings. In the first one, both the speech-source signal and the localized noise source signal are chosen to be temporally correlated, i.e. *colored* signals. The (non-stationary) speech-source signal  $s_1(l)$  is composed of male and female speech of in total 51 s duration [38], while for the localized noise source signal  $s_2(l)$ , stationary pink noise is used. This setting is evaluated for both the time domain and the STFT domain implementations. In the time domain, for the chosen setup, cf. Sec. VII-A3, the coloration of the source signals causes a biased filter estimate. In the second setting, both  $s_1(l)$  and  $s_2(l)$  are chosen to be temporally uncorrelated, i.e. *white*, stationary signals, and have been generated independently from the source signals in the first setting. This setting is evaluated in the time domain implementations only, leading to an unbiased filter estimate and serving as a reference in order to illustrate the effect of the bias in the first setting. Since sensor noise is always present in practice, spatially and temporally uncorrelated noise  $\mathbf{v}(l)$  is added in all simulations. Note that due to the incoherent noise, the time domain simu-



lation results may at most approximately reach the theoretical limits discussed in Sec. IV-B1 and Sec. V-B1.

The power of the noise components is defined via the signal-to-coherent-noise ratio  $SNR_y^{coh}$  and the signal-to-incoherent-noise ratio  $SNR_y^{inc}$  in the first microphone, i.e.

$$SNR_y^{coh} = 10 \log_{10} \frac{\sum_l |x_{1,1}(l)|^2}{\sum_l |x_{2,1}(l)|^2} \text{ dB}, \quad (58)$$

$$SNR_y^{inc} = 10 \log_{10} \frac{\sum_l |x_{1,1}(l)|^2}{\sum_l |v_1(l)|^2} \text{ dB}, \quad (59)$$

where the reverberant-speech component  $x_{1,1}(l)$  in the first microphone is considered to be the useful signal.

### 3) MCLP and GSC implementation:

a) *Time Domain:* In the time domain, we define the direct speech component to be the target component, i.e. we choose  $\delta = L_b = 11$  samples, corresponding to the energy spread of a single acoustic wave, cf. Sec. VII-A1, yielding  $d = 11$  for MCLP and  $d \geq 11$  for the GSC, cf. Sec. II-C and Sec. II-D. An ideal GSC blocking matrix  $\mathbf{B}$  was designed, cf. (24)–(25). The filter  $\mathbf{g}$  is chosen to be a matched filter (MF) such that  $\mathbf{B}^H \mathbf{g} = \mathbf{0}$ , both for the GSC and MCLP. We choose  $\Omega_{(\cdot)} = \mathbf{I}$  in (36), leading to a biased filter estimate for colored source signals. The effect of the bias is shown by comparing the performance for both colored and white source signals.

b) *STFT Domain:* In the STFT domain, using square-root-Hann windows of 512 samples with 50% overlap, we choose  $\delta = L_b = 1$  frame. The GSC blocking matrix  $\mathbf{B}$  uses an estimate of the RETFs, which we obtain as presented in [28], [35], [39]: we estimate the *average* spatial correlation matrix of the microphone signals using the whole batch, and the spatial correlation matrix of the stationary noise components using 5 s noise-only frames, such that the spatial speech-component correlation matrix can be estimated by subtraction. Then, from the spatial speech-component correlation matrix estimate, the RETF relative to the first microphone is estimated using the GEVD, assuming that the spatial coherence matrix of the late reverberant-speech component in frame  $l$  may be modeled as diffuse [28], [35]. Again, the filter  $\mathbf{g}$  is chosen to be an MF with  $\mathbf{B}^H \mathbf{g} = \mathbf{0}$ , i.e.  $\mathbf{g}$  is a normalized version of the RETF estimate. For  $\Omega_{(\cdot)}$  in (36), we use an estimate of  $\Psi_{\hat{q}_{s_1} | (\cdot)}$ , which in the STFT domain is diagonal for  $d = 1$ , cf. Sec. III-B. Since  $\mathbf{g}$  is a normalized version of the RETF estimate, estimating the PSDs  $\psi_{\hat{q}_{s_1}}(l)$  in  $\Psi_{\hat{q}_{s_1} | (\cdot)}$  corresponds to estimating the early-reverberant speech component in the first microphone. Again, this can be done using the GEVD [28], [35], now applied to a *recursive* estimate of the spatial speech-component correlation matrix in frame  $l$ . See, e.g., [28], [35] for a detailed and more formal discussion on GEVD-based RETF and PSD estimation.

### 4) Performance Measures:

a) *Time Domain:* In the time domain, equivalently to  $q_{s_1}(l)$ ,  $q_{s_2}(l)$ ,  $q_v(l)$  and  $\hat{q}_{s_1}(l)$ , we define the individual source components of  $e(n)$  as  $e_{s_1}(l)$ ,  $e_{s_2}(l)$ ,  $e_v(l)$  and  $\hat{e}_{s_1}(l)$ , where  $\hat{e}_{s_1}(l) = \hat{q}_{s_1}(l)$ , cf. (49), (57). With  $\sigma \in \{q, e\}$ , the signal-to-coherent-noise ratio  $SNR_\sigma^{coh}$ , the signal-to-incoherent-noise ratio  $SNR_\sigma^{inc}$ , the signal-to-total-noise ratio  $SNR_\sigma^{tot}$ , and the

signal-to-reverberation ratio  $SRR_\sigma$  at the MF output and the MCLP and GSC output are defined as

$$SNR_\sigma^{tot} = 10 \log_{10} \frac{\sum_l |\hat{\sigma}_{s_1}(l)|^2}{\sum_l |\sigma_{s_2}(l) + \sigma_v(l)|^2} \text{ dB}, \quad (60)$$

$$SRR_\sigma = 10 \log_{10} \frac{\sum_l |\hat{\sigma}_{s_1}(l)|^2}{\sum_l |\sigma_{s_1}(l) - \hat{\sigma}_{s_1}(l)|^2} \text{ dB}, \quad (61)$$

where the component  $\hat{\sigma}_{s_1}(l)$  is considered to be the useful signal. Please note that  $q(l)$ , and hence for  $\sigma = q$  also the measures in (58)–(61), are independent of the particular framework. Further, note that in the denominator of (61), for  $\sigma = q$ , the difference  $q_{s_1}(l) - \hat{q}_{s_1}(l)$  equals the late reverberant-speech component  $\hat{q}_{s_1}(l)$ , while for  $\sigma = e$ , the difference  $e_{s_1}(l) - \hat{e}_{s_1}(l)$  comprises not only residual reverberation, but also a bias term in the general case. For evaluation, we use the improvement in  $SNR^{tot}$  and  $SRR$ , i.e.

$$\Delta SNR^{tot} = SNR_e^{tot} - SNR_q^{tot}, \quad (62)$$

$$\Delta SRR = SRR_e - SRR_q. \quad (63)$$

b) *STFT Domain:* In the STFT domain, the target component  $\hat{q}_{s_1}(l)$  cannot be observed separately, since the sub-band convolution model in (1) poses an approximation of the time-domain convolution only [33]. Further, due to overlapping frames in the STFT processing and incomplete blocking in the GSC, the target component  $\hat{q}_{s_1}(l)$  may not be completely maintained in  $e(l)$ , such that the measures in (60)–(63) are not suitable in the STFT domain. Instead, we define the direct-component in  $q(l)$  as a reference signal, which *cannot* be assumed to be equivalent to  $\hat{q}_{s_1}(l)$ . Then, with  $\sigma \in \{q, e\}$ , for  $\sigma(l)$  and  $\sigma_{s_1}(l)$ , respectively, we compute the frequency-weighted segmental signal-to-noise-plus-reverberation ratio and the frequency-weighted segmental signal-to-reverberation ratio [36], denoted as  $SNRR^{fwseg}$  and  $SRR^{fwseg}$  and indicating the dereverberation-plus-noise-reduction performance and the dereverberation-only performance.

5) *Varied Parameters:* In the time domain, simulations are carried out for different values of the following parameters:  $SNR_y^{coh}$ ,  $SNR_y^{inc}$ ,  $L_w$ , and  $M$ . The filter length  $L_w$  is presented relatively to the theoretical minimum given in (22), (28), denoted by  $L_w^{rel}$ . While one parameter is varied, the others are fixed at  $SNR_q^{coh} = 0$  dB,  $SNR_q^{inc} = 90$  dB,  $L_w^{rel} = 1$ , and  $M = 8$ , i.e. all simulations intersect at this point. For  $N = 2$ , the minimum number of microphones required by MCLP and the GSC is given by  $M = 3$  and  $M = 4$ , respectively, cf. (22), (28). If the number of microphones  $M$  falls below this required minimum, the filter length is computed setting the denominators in (22), (28) to one. Simulations posing nearly ideal conditions, i.e. sufficiently high  $SNR_q^{inc}$ ,  $L_w^{rel} \geq 1$  and sufficiently many microphones  $M$ , validate the theoretical results in Sec. IV and Sec. V, with minor deviations occurring due to the LS approximation in (36) of the Wiener solution in (40) and remaining low-level incoherent noise.

In the STFT domain, simulations are carried out for different values of  $SNR_y^{coh}$  only, with  $SNR_q^{inc} = 90$  dB,  $L_w^{rel} = 1$ , and  $M = 8$ . Since complete blocking is not achieved in the STFT domain, decreased performance is expected for the GSC.

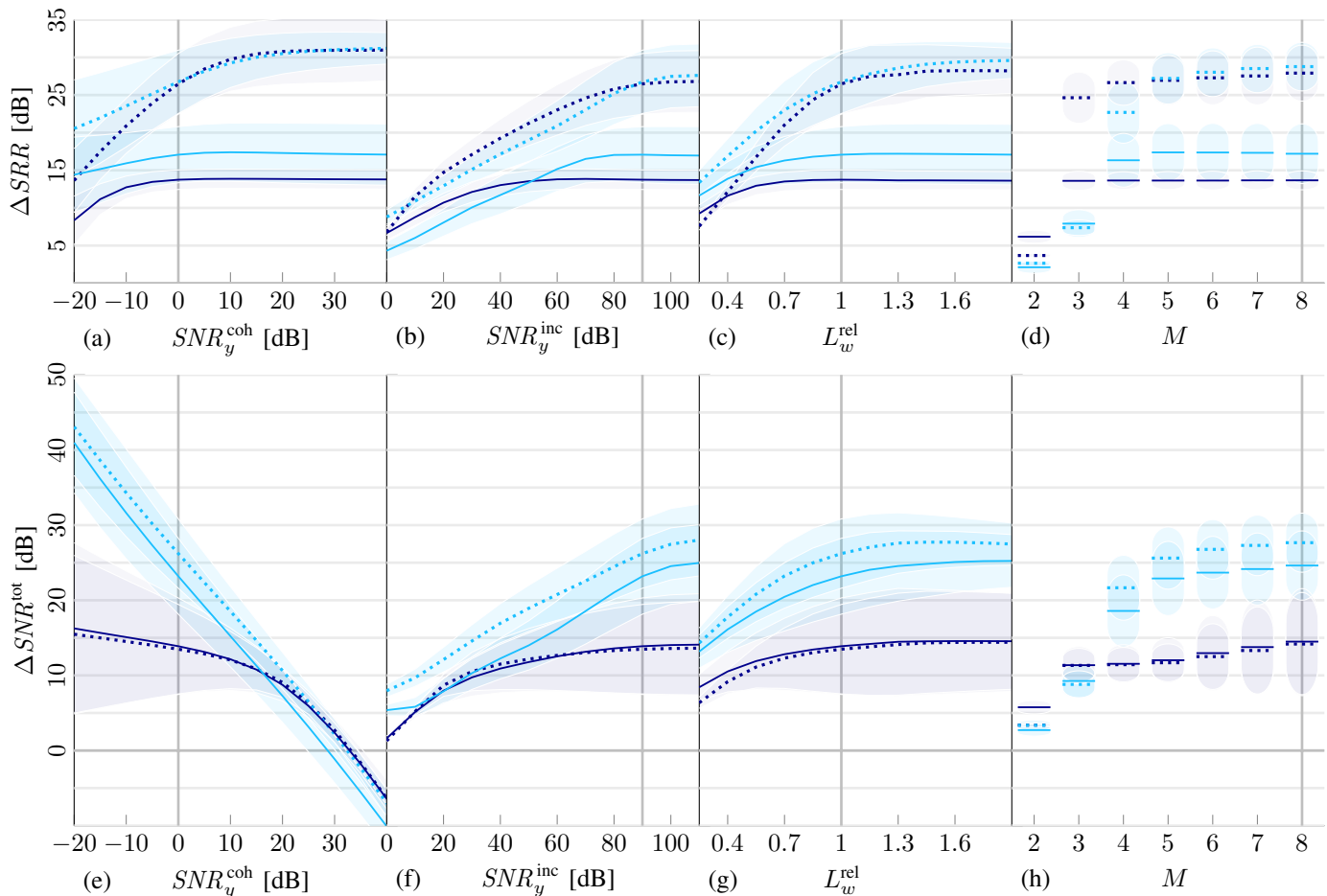


Fig. 6: Dereverberation/noise reduction performance  $\Delta SRR/\Delta SNR^{tot}$  versus (a)/(e) the signal-to-coherent-noise ratio  $SNR_q^{coh}$ , (b)/(f) the signal-to-incoherent-noise ratio  $SNR_q^{inc}$ , (c)/(g) the relative filter length  $L_w^{rel}$  and (d)/(h) the number of microphones  $M$  for colored and white source signals of the MCLP framework, respectively denoted by [—] and [⋯⋯], and the GSC framework, respectively denoted by [—] and [⋯⋯]. The vertical grid lines indicate the intersection point of the individual subplots. The shaded areas represent the standard deviation.

## B. Simulation Results

We now discuss the time and STFT domain simulation results in Sec. VII-B1 and Sec. VII-B2, respectively.

1) *Time Domain*: The performance of both frameworks in terms of  $\Delta SRR$  and  $\Delta SNR^{tot}$  are shown in Fig. 6. We first discuss the dereverberation performance, followed by the noise reduction performance.

a) *Dereverberation*: From Fig. 6 (a)–(d) we observe that under favorable conditions with predominantly late-reverberant-speech interference, i.e. for high  $SNR_y^{coh}$ , high  $SNR_y^{inc}$ ,  $L_w^{rel} \geq 1$  and sufficiently high  $M$ , the  $SRR$  improvement of both MCLP and GSC converge to the same value of around 31 dB for white source signals, respectively denoted by [⋯⋯] and [⋯⋯]. This upper limit is determined by the LS approximation (36) of the Wiener Solution of the (40). In all conditions, for colored source signals, the target component  $\hat{q}_{s_1}(l)$  is partially whitened due to the biased filter estimate, leading to a performance drop for both MCLP [—] and the GSC [—], cf. Sec. IV-B1 and Sec. V-B1 and Table I. The GSC reaches up to 18 dB  $\Delta SRR$ , outperforming MCLP by

3 dB. This is due to the potential delay between the direct component and the first reflection, increasing  $d$  for the GSC and thereby decreasing the bias, cf. Sec. VII-A3. The higher standard deviation of  $\Delta SRR$  for the GSC is a result of the variation of this delay over different source and microphone array positions.

As can be seen in Fig. 6 (a), for white source signals, MCLP shows a rather high sensitivity towards coherent noise for  $SNR_y^{coh} < 10$  dB, whereas the GSC is somewhat less sensitive. This can be explained by the limited number of observations  $L_{obs}$  used in the LS approximation (36) of the Wiener solution (40), causing LS to focus on the suppression of components with higher power, i.e. here on (late) coherent-noise suppression. For colored source signals, the effect is less pronounced in both frameworks.

While it can be observed from Fig. 6 (b) that both MCLP and GSC are highly sensitive to incoherent noise, the results also indicate an up to 2.5 dB lower performance of the GSC as compared to MCLP for  $SNR_y^{inc} < 90$  dB and  $< 55$  dB for white and colored source signals, respectively. The reason for this may lie in the GSC blocking matrix, which by construction

causes a cross-correlation of the incoherent-noise components in the data-dependent filter input  $\mathbf{u}(l)$ , as opposed to the mere delay in MCLP. Hence, for the GSC, not only the autocorrelation submatrices of  $\Psi_{\tilde{\mathbf{u}}}$  are affected by incoherent noise, but also the cross-correlation submatrices.

As shown in Fig. 6 (c),  $\Delta SRR$  saturates for both MCLP and GSC above  $L_w^{\text{rel}} = 1$ , both for white and colored source signals, as expected from theory. The GSC performs slightly better than MCLP if  $L_w^{\text{rel}} < 1$ . We may however state that for both frameworks, undermodeling is not extremely critical, as even at  $L_w^{\text{rel}} = 0.7$ , we obtain  $\Delta SRR$  values above 22 dB for white source signals, while the performance is hardly affected for colored source signals.

From Fig. 6 (d), we note that  $\Delta SRR$  drops sharply for both MCLP and GSC if the number of microphones is smaller than required, i.e.  $M < 3$  and  $M < 4$ , respectively. This holds for both white and colored source signals. MCLP reaches saturation at  $M = 3$ , while the GSC saturates at  $M = 5$  only instead of  $M = 4$ . This may be caused by remaining low-level incoherent noise and possibly nearly common zeros in the transfer functions corresponding to  $\mathbf{H}_B$  in (27) for  $M = 4$ .

*b) Noise Reduction:* From Fig. 6 (e)–(h) we observe that under favorable conditions with predominantly coherent-noise interference, i.e. for low  $SNR_y^{\text{coh}}$ , high  $SNR_y^{\text{inc}}$ ,  $L_w^{\text{rel}} \geq 1$  and sufficiently high  $M$ , the GSC [·····, —] shows increasing improvement in terms of  $\Delta SNR^{\text{tot}}$  for decreasing values of  $SNR_y^{\text{coh}}$ , while for MCLP [·····, —],  $\Delta SNR^{\text{tot}}$  is limited to at most 15 dB. the GSC [·····, —] clearly outperforms MCLP [·····, —] in terms of  $\Delta SNR^{\text{tot}}$ . This is due to the GSC suppressing the entire coherent-noise component  $q_{s_2}(l)$ , while MCLP suppresses the late coherent-noise component  $\tilde{q}_{s_2}(l)$  only, cf. (47)–(49), (55)–(57) and Table I. Note that MCLP exhibits a stronger standard deviation in  $\Delta SNR^{\text{tot}}$  than the GSC. This is caused by the varying power of the early coherent-noise component  $\tilde{q}_{s_2}(l)$ , as the power of the individual direct components at the output of the MF may be distributed over a range potentially larger than  $d$ , depending on the angle between the speech source and the coherent-noise source. In all conditions, the GSC performs somewhat worse for colored signals than white signals, while no significant difference is found for MCLP.

Fig. 6 (e) indicates that the GSC exceeds MCLP for  $SNR_y^{\text{coh}} < 20$  dB, while both frameworks perform similarly for high  $SNR_y^{\text{coh}}$  values. For the GSC,  $\Delta SNR^{\text{tot}}$  decreases at a rate of slightly less than 10 dB  $\Delta SNR^{\text{tot}}$  per 10 dB  $SNR_y^{\text{coh}}$ , such that the noise power at the output is almost constant throughout the simulated range. This implies that for  $SNR_y^{\text{coh}} \geq 35$  dB, the total noise power is in fact even boosted as compared to the output of the MF, both for MCLP and the GSC. Again, this effect can be explained by the limited number of observations  $L_{\text{obs}}$  in the LS estimate (36), here causing LS to focus on reverberant-speech suppression.

As can be seen from Fig. 6 (f), for both white and colored source signals, the GSC exceeds MCLP for higher  $SNR_y^{\text{inc}}$  values, while the difference reduces for lower values.

As shown in Fig. 6 (g), for both white and colored source signals, both MCLP and the GSC again saturate for  $L_w^{\text{rel}} \geq 1$ . Again, undermodeling does not appear to be extremely critical.

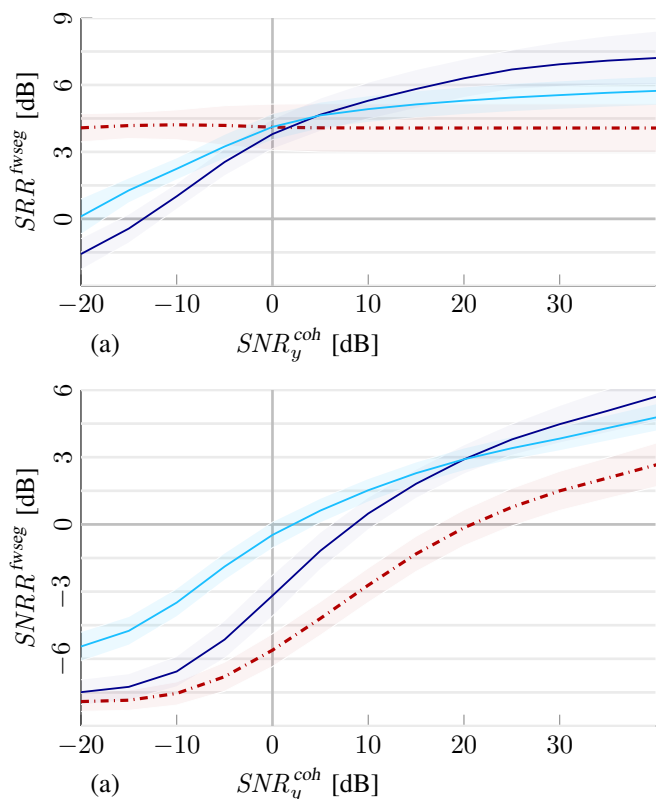


Fig. 7: (a) dereverberation-only/(b) dereverberation-plus-noise-reduction performance in terms of  $SRR^{\text{fwseg}}/SNR^{\text{fwseg}}$  versus the signal-to-coherent-noise ratio  $SNR_y^{\text{coh}}$  of the MF, the MCLP and the GSC framework, respectively denoted by [·····], [—], and [—]. The shaded areas represent the standard deviation.

From Fig. 6 (h), for both white and colored source signals, we once more find a sharp performance drop for both MCLP and the GSC if  $M < 3$  and  $M < 4$ , respectively. Again, saturation is reached at  $M = 3$  and  $M = 5$ , respectively.

*2) STFT Domain:* The performance of both frameworks in terms of  $SRR^{\text{fwseg}}$  and  $SNR^{\text{fwseg}}$  are shown in Fig. 7, where the performance of the MF serves as a reference.

From Fig. 7 (a), we note that the dereverberation-only performance of the MF [·····] in terms of  $SRR^{\text{fwseg}}$  remains almost constant around 4.1 dB. At high  $SNR_y^{\text{coh}}$  values with predominantly late-reverberant-speech interference, MCLP [—] and the GSC [—] outperform the MF by up to 3.1 dB and 4.1 dB, respectively. Note that in theory, for *complete blocking*, i.e. if (27) is satisfied, the GSC is expected to perform as effectively as MCLP in terms of dereverberation, cf. (47)–(49), (55)–(57), Table I, and the time domain simulation results in Sec. VII-B1. However, since the sub-band convolution model in (1) poses an approximation of the time-domain convolution only, and since the RETFs used in the GSC blocking matrix are subject to estimation errors, cf. Sec. VII-A3b, complete blocking is not achieved in the STFT domain. Due to *incomplete blocking*, the GSC hence suffers from some amount of early reverberant-speech cancellation and in addition from incomplete prediction of the late reverberant-speech component  $\tilde{q}_{s_1}(l)$ , leading to reduced

dereverberation performance in comparison to MCLP. At low  $SNR_y^{coh}$  values with predominantly coherent-noise interference, we find that both MCLP and the GSC perform worse than the MF, indicating speech distortion. Again, this effect can be explained by the limited number of observations  $L_{obs}$  used in the LS estimate (36), causing LS to focus on (late) coherent-noise suppression, cf. Sec. VII-B1. Since the GSC is able to suppress the early coherent-noise component  $\hat{q}_{s_2}(l)$  also, the GSC performance is less affected.

From Fig. 7 (b), we note that the dereverberation-plus-noise-reduction performance of the MF in terms of  $SNRR^{fwseg}$  ranges between  $-8$  dB for  $SNR_y^{coh} = -20$  dB and  $2.6$  dB for  $SNR_y^{coh} = 40$  dB, where the upper limit is still affected by the noise component, as can be seen by comparison to the dereverberation-only performance in Fig. 7 (a). At high  $SNR_y^{coh}$  values MCLP and the GSC outperform the MF by up to  $32.8$  dB and  $2.1$  dB, respectively. At low  $SNR_y^{coh}$  values, MCLP performs only somewhat better than the MF, while the GSC in contrast outperforms the MF by up to  $5.1$  dB. This difference at low  $SNR_y^{coh}$  values is expected as MCLP suppresses the late coherent-noise component  $\hat{q}_{s_2}(l)$  only, while the GSC suppresses the entire coherent-noise component  $q_{s_2}(l)$ , cf. (47)–(49), (55)–(57), Table I, and the time domain simulation results in Sec. VII-B1. Audio examples of the STFT domain simulations are available online [40].

## VIII. CONCLUSION

In this paper, we formally analyzed and compared the MCLP and GSC frameworks in terms of blind dereverberation and noise reduction performance. Both frameworks are theoretically able to perform complete dereverberation if incoherent noise is absent. Due to the use of a blocking matrix, the GSC is theoretically able to completely cancel coherent noise in the absence of incoherent noise, while MCLP cancels the late coherent-noise component only. For complete cancellation, the GSC requires one additional microphone as compared to MCLP. Furthermore, the blocking matrix design requires spatial information in form of the early speech-source RIR or the RETF, which needs to be acquired in practice. In order to confirm the theory and to assess the practical relevance of the theoretical findings, we carried out time domain simulations using oracle knowledge on the early RIRs, resulting in *complete blocking* of the early reverberant-speech component, and STFT domain simulations using estimated RETFs, resulting in *incomplete blocking*.

The simulation results confirm that in terms of noise reduction, as opposed to the GSC performance, the performance of MCLP is limited. In terms of dereverberation, the GSC performs equally well if *complete blocking* is achieved, as expected from the theoretical analysis, but performs inferior for *incomplete blocking*. Both MCLP and the GSC exhibit strong sensitivity to incoherent noise. For both frameworks, dereverberation and noise reduction performance reach their maximum at a relative filter length of about one, while moderate undermodeling of the filter length does not appear to be extremely critical. The simulations further confirm that for one coherent-noise component, the GSC requires four microphones, while MCLP requires three microphones only.

In summary, we can state that if sufficiently many microphones are available and complete blocking is achieved, the GSC performs superior to MCLP in terms of noise reduction and equally well in terms of dereverberation, but inferior in terms of dereverberation for incomplete blocking. In practice therefore, in acoustic conditions with only mild noise but predominantly late-reverberant-speech interference, MCLP is to be preferred, while in case of predominantly noise but mild to moderate late-reverberant-speech interference, the GSC is to be preferred. In acoustic conditions with both strong reverberation and strong noise, combined schemes may be most appropriate.

## APPENDIX A

1): Analogously to (16)–(17), let  $\dot{\mathbf{C}} \in \mathbb{C}^{(NL_s \times NL_s)}$  and its counterpart  $\ddot{\mathbf{C}}$  be defined by

$$\dot{\mathbf{C}} = \begin{pmatrix} \mathbf{I}^{d \times d} & \mathbf{0}^{d \times (NL_s - d)} \\ \mathbf{0}^{(NL_s - d) \times d} & \mathbf{0}^{d \times (NL_s - d)} \end{pmatrix}, \quad (64)$$

$$\ddot{\mathbf{C}} = \mathbf{I}^{NL_s \times NL_s} - \dot{\mathbf{C}}, \quad (65)$$

and let  $\underline{\Psi}_{\bar{s}} \in \mathbb{C}^{NL_s - d \times NL_s - d}$  be the submatrix of  $\Psi_{\bar{s}}$  spanning its last  $NL_s - d$  rows and columns, such that

$$\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}} = \begin{pmatrix} \mathbf{0}^{d \times d} & \mathbf{0}^{d \times (NL_s - d)} \\ \mathbf{0}^{(NL_s - d) \times d} & \underline{\Psi}_{\bar{s}} \end{pmatrix}. \quad (66)$$

With  $\mathbf{H}\mathbf{B}$  as given in (27), the expression  $((\mathbf{H}\mathbf{B})^H \Psi_{\bar{s}} \mathbf{H}\mathbf{B})^+$  in (54) can then be written as

$$\begin{aligned} ((\mathbf{H}\mathbf{B})^H \Psi_{\bar{s}} \mathbf{H}\mathbf{B})^+ &= (\mathbf{H}_B^H \underline{\Psi}_{\bar{s}} \mathbf{H}_B)^+ \\ &= \mathbf{H}_B^+ \underline{\Psi}_{\bar{s}}^{-1} \mathbf{H}_B^{+T}. \end{aligned} \quad (67)$$

Inserting (27) and (67) in (54) while noting that  $\mathbf{H}_B \mathbf{H}_B^+ = \mathbf{I}$  since  $\mathbf{H}_B$  is assumed to have full row rank, and further using (66) and Lemma 1 from Appendix B, we obtain

$$\begin{aligned} z(l) &= \left( \mathbf{H}\mathbf{B} ((\mathbf{H}\mathbf{B})^H \Psi_{\bar{s}} \mathbf{H}\mathbf{B})^+ (\mathbf{H}\mathbf{B})^H \Psi_{\bar{s}} \mathbf{H}\mathbf{g} \right)^H \mathbf{s}(l) \\ &= \left( \left( \begin{pmatrix} \mathbf{0}^{d \times d} & \mathbf{0}^{d \times (NL_s - d)} \\ \mathbf{0}^{(NL_s - d) \times d} & \underline{\Psi}_{\bar{s}}^{-1} \end{pmatrix} \Psi_{\bar{s}} \mathbf{H}\mathbf{g} \right)^H \right) \mathbf{s}(l) \\ &= ((\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+ \Psi_{\bar{s}} \mathbf{H}\mathbf{g})^H \mathbf{s}(l). \end{aligned} \quad (68)$$

With  $\dot{\mathbf{C}} + \ddot{\mathbf{C}} = \mathbf{I}$ , the matrix  $\Psi_{\bar{s}}$  may be written as

$$\begin{aligned} \Psi_{\bar{s}} &= \dot{\mathbf{C}}\Psi_{\bar{s}} + \ddot{\mathbf{C}}\Psi_{\bar{s}} \\ &= \dot{\mathbf{C}}\Psi_{\bar{s}} + \ddot{\mathbf{C}}\Psi_{\bar{s}}\dot{\mathbf{C}} + \ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}}. \end{aligned} \quad (69)$$

Substituting (69) in (68), we find  $(\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+(\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}}) = \ddot{\mathbf{C}}$  from (66), while  $(\dot{\mathbf{C}}\Psi_{\bar{s}}\dot{\mathbf{C}})^+\dot{\mathbf{C}}\Psi_{\bar{s}} = \mathbf{0}$ , such that  $(\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+\Psi_{\bar{s}}$  in (68) becomes

$$(\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+\Psi_{\bar{s}} = \ddot{\mathbf{C}} + (\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+\ddot{\mathbf{C}}\Psi_{\bar{s}}\dot{\mathbf{C}}. \quad (70)$$

Using (64)–(65), (16)–(17), and Lemma 1 from Appendix B, the term  $(\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+\ddot{\mathbf{C}}\Psi_{\bar{s}}\dot{\mathbf{C}}$  in (70) takes the form

$$(\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+\ddot{\mathbf{C}}\Psi_{\bar{s}}\dot{\mathbf{C}}$$

$$= \begin{pmatrix} (\ddot{\mathbf{C}}\Psi_{\bar{s}_1}\ddot{\mathbf{C}})^+ + \ddot{\mathbf{C}}\Psi_{\bar{s}_1}\dot{\mathbf{C}} & \mathbf{0}_{L_s \times (N-1)L_s} \\ \mathbf{0}^{(N-1)L_s \times L_s} & \mathbf{0}^{(N-1)L_s \times (N-1)L_s} \end{pmatrix}. \quad (71)$$

Inserting (70) in (68), multiplying out, using (64)–(65) and (14)–(17), it can be shown that

$$\begin{aligned} z(l) &= (\ddot{\mathbf{C}}\mathbf{H}\mathbf{g})^H \mathbf{s}(l) + ((\ddot{\mathbf{C}}\Psi_{\bar{s}}\ddot{\mathbf{C}})^+ + \ddot{\mathbf{C}}\Psi_{\bar{s}}\dot{\mathbf{C}}\mathbf{H}\mathbf{g})^H \mathbf{s}(l) \\ &= \sum_{n=2}^N q_{s_n}(l) + \ddot{q}_{s_1}(l) \\ &\quad + ((\ddot{\mathbf{C}}\Psi_{\bar{s}_1}\ddot{\mathbf{C}})^+ + \ddot{\mathbf{C}}\Psi_{\bar{s}_1}\dot{\mathbf{C}}\mathbf{H}_1\mathbf{g})^H \mathbf{s}_0(l). \end{aligned} \quad (72)$$

With Lemma 1 from Appendix B and  $\ddot{\mathbf{C}}_d$  defined in (19),  $(\ddot{\mathbf{C}}_0\Psi_{\bar{s}_0}\ddot{\mathbf{C}}_0)^+$  can be written as

$$(\ddot{\mathbf{C}}_1\Psi_{\bar{s}_1}\ddot{\mathbf{C}}_1)^+ = \ddot{\mathbf{C}}_d^H (\ddot{\mathbf{C}}_d\Psi_{\bar{s}_1}\ddot{\mathbf{C}}_d^H)^+ \ddot{\mathbf{C}}_d, \quad (73)$$

Substituting (73) into (72) and extracting  $\ddot{\mathbf{C}}_d^H$  on the left, we obtain (55)–(56).

2): With Lemma 1 in Appendix B, the first term in (56),  $(\ddot{\mathbf{C}}_d\Psi_{\bar{s}_0}\ddot{\mathbf{C}}_d^H)^+$ , can be written as

$$(\ddot{\mathbf{C}}_d\Psi_{\bar{s}_1}\ddot{\mathbf{C}}_d^H)^+ = \begin{pmatrix} \Psi_{\bar{s}_1}^{-1} & \mathbf{0}_{(L_s-d) \times d} \\ \mathbf{0}^{d \times (L_s-d)} & \mathbf{0}^{d \times d} \end{pmatrix}, \quad (74)$$

where  $\Psi_{\bar{s}_1} \in \mathbb{C}^{L_s-d \times L_s-d}$  is the submatrix of  $\Psi_{\bar{s}_1}$  spanning the last  $L_s - d$  rows and columns, matching the first term in (48) for  $n = 1$ . Note that for  $d = L_b$ ,  $\Psi_{\bar{s}_1}^{-1} \in \mathbb{C}^{L_s-d \times L_s-d}$  in (74) and  $\Psi_{\bar{s}_1}^{-1} \in \mathbb{C}^{L_s \times L_s}$  in (48) also correspond in terms of dimensions: inserting (26) into (9) yields  $L_s - d = L_h + L_w - 1$  in the GSC, while inserting (21) into (9) yields  $L_s = L_h + L_w - 1$  in MCLP. Finally, since the second term in (56),  $\ddot{\mathbf{C}}_d\Psi_{\bar{s}_1}\dot{\mathbf{C}}\mathbf{H}_1\mathbf{g}$ , is equivalent to the second term in (48), both expressions (56) and (48) yield the same bias component.

## APPENDIX B

*Lemma 1: The pseudoinverse  $\mathbf{A}^+$  of a block-diagonal matrix  $\mathbf{A}$  defined by the blocks  $\mathbf{A}_n$ ,  $n = 1 \dots N$ , on its diagonal is given by a block-diagonal matrix composed of the pseudoinverses  $\mathbf{A}_n^+$  of the individual blocks, i.e.*

$$\begin{aligned} \text{if} \quad & \mathbf{A} = \text{blkdiag} [\mathbf{A}_1, \dots, \mathbf{A}_N], \\ \text{then} \quad & \mathbf{A}^+ = \text{blkdiag} [\mathbf{A}_1^+, \dots, \mathbf{A}_N^+]. \end{aligned}$$

This lemma can be proven easily by verifying the four criteria defining the pseudoinverse  $\mathbf{A}^+$  of the matrix  $\mathbf{A}$ , i.e.  $\mathbf{A}\mathbf{A}^+\mathbf{A} = \mathbf{A}$ ,  $\mathbf{A}^+\mathbf{A}\mathbf{A}^+ = \mathbf{A}^+$ ,  $(\mathbf{A}\mathbf{A}^+)^H = \mathbf{A}\mathbf{A}^+$ , and  $(\mathbf{A}^+\mathbf{A})^H = \mathbf{A}^+\mathbf{A}$ . It is further important to note that the pseudoinverse of a zero matrix is equal to its transpose.

## ACKNOWLEDGMENTS

This work was carried out at the ESAT Laboratory of KU Leuven, in the frame of KU Leuven internal funds C2-16-00449 "Distributed Digital Signal Processing for Ad-hoc Wireless Local Area Audio Networking," VES/16/032, VES/19/004, and Impulse Fund IMP/14/037; IWT O&O Project 150611 "Proof-of-concept of a Rationed Architecture

for Vehicle Entertainment and NVH Next-generation Acoustics (RAVENNA)"; VLAIO TETRA Project HBC.2016.0085 "m-sense: innovative use of sensors in mobile platforms" and O&O Project HBC.2017.0358 "SPOTTTomorrow's Scalable and Personalised advertising Technology, Today"; and EU FP7-PEOPLE Marie Curie Initial Training Network "Dereverberation and Reverberation of Audio, Music, and Speech (DREAMS)," funded by the European Commission under Grant 316969. The research leading to these results has received funding from the European Research Council under the European Union's Horizon 2020 research and innovation program / ERC Consolidator Grant: SONORA (773268).

## REFERENCES

- [1] R. Beutelmänn and T. Brand, "Prediction of speech intelligibility in spatial noise and reverberation for normal-hearing and hearing-impaired listeners," *J. Acoustic. Society America*, vol. 120, no. 1, pp. 331–342, Apr. 2006.
- [2] M. Miyoshi and Y. Kaneda, "Inverse filtering of room acoustics," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 36, no. 2, pp. 145–152, Feb. 1988.
- [3] T. Hikichi, M. Delcroix, and M. Miyoshi, "Inverse filtering for speech dereverberation less sensitive to noise and room transfer function fluctuations," *EURASIP J. Adv. Signal Process.*, pp. 1–12, Apr. 2007.
- [4] F. Lim and P. Naylor, "Statistical modelling of multichannel blind system identification errors," in *Proc. Int. Workshop Acoustic Signal Enhancement (IWAENC 2014)*, Juan-les-Pins, France, Sept. 2014, pp. 119–123.
- [5] M. R. P. Thomas, N. D. Gaubitch, E. A. P. Habets, and P. A. Naylor, "An insight into common filtering in noisy speech blind system identification," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP 2012)*, Kyoto, Japan, Mar. 2012, pp. 521–524.
- [6] M. Triki and D. T. M. Slock, "Blind dereverberation of a single source based on multichannel linear prediction," in *Proc. Int. Workshop Acoustic Echo Noise Control (IWAENC 2005)*, Eindhoven, Netherlands, Sep. 2005, pp. 173–176.
- [7] M. Delcroix, T. Hikichi, and M. Miyoshi, "Precise dereverberation using multichannel linear prediction," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, no. 2, pp. 430–440, Jan. 2007.
- [8] T. Nakatani, B. H. Juang, T. Yoshioka, K. Kinoshita, M. Delcroix, and M. Miyoshi, "Speech dereverberation based on maximum-likelihood estimation with time-varying Gaussian source model," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 16, no. 8, pp. 1512–1527, Nov. 2008.
- [9] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B. H. Juang, "Speech dereverberation based on variance-normalized delayed linear prediction," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 7, pp. 1717–1731, Aug. 2010.
- [10] T. Yoshioka, N. T. M. Miyoshi, and H. G. Okuno, "Blind separation and dereverberation of speech mixtures by joint optimization," *IEEE Trans. Audio, Speech, Language Process.*, vol. 19, pp. 69–84, 2011.
- [11] T. Yoshioka and T. Nakatani, "Generalization of multi-channel linear prediction methods for blind MIMO impulse response shortening," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 20, no. 10, pp. 2707–2720, July 2012.
- [12] T. Yoshioka, "Dereverberation for reverberation-robust microphone arrays," in *Proc. 21st European Signal Process. Conf. (EUSIPCO 2013)*, Marrakech, Morocco, Sep. 2013, pp. 1–5.
- [13] A. Jukić, T. van Waterschoot, T. Gerkmann, and S. Doclo, "Multi-channel linear prediction-based speech dereverberation with sparse priors," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 9, pp. 1509–1520, June 2015.
- [14] S. Braun and E. A. P. Habets, "Online dereverberation for dynamic scenarios using a Kalman filter with an autoregressive model," *IEEE Signal Process. Letters*, vol. 23, no. 12, pp. 1741–1745, Dec. 2016.
- [15] A. Jukić, T. van Waterschoot, and S. Doclo, "Adaptive speech dereverberation using constrained sparse multichannel linear prediction," *IEEE Signal Process. Letters*, vol. 24, no. 1, pp. 101–105, Jan. 2017.
- [16] T. Dietzen, S. Doclo, A. Spriet, W. Tirry, M. Moonen, and T. van Waterschoot, "Low complexity Kalman filter for multi-channel linear prediction based blind speech dereverberation," in *IEEE Workshop Appl.*

- Signal Process. Audio Acoust. (WASPAA)*, New Paltz, NY, USA, Oct. 2017.
- [17] S. Braun and E. A. P. Habets, "Linear prediction-based online dereverberation and noise reduction using alternating Kalman filters," *IEEE/ACM Trans. Audio, Speech, Language Proc.*, vol. 26, no. 6, pp. 240–251, June 2018.
- [18] J. Benesty, J. Chen, Y. Huang, and J. Dmochowski, "On microphone-array beamforming from a MIMO acoustic signal processing perspective," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 15, no. 3, pp. 1053–1063, Mar. 2007.
- [19] E. A. P. Habets, J. Benesty, I. Cohen, S. Gannot, and J. Dmochowski, "New insights into the MVDR beamformer in room acoustics," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 18, no. 1, pp. 158–170, Jan. 2010.
- [20] F. Lim, M. R. P. Thomas, and P. Naylor, "MINTFormer: A spatially aware channel equalizer," in *IEEE Workshop Appl. Signal Process. Audio Acoust. (WASPAA 2013)*, New Paltz, NY, USA, Oct. 2013, pp. 1–4.
- [21] O. Schwartz, S. Gannot, and E. A. P. Habets, "Multi-microphone speech dereverberation and noise reduction using relative early transfer functions," *IEEE/ACM Trans. Audio, Speech, Language Proc.*, vol. 23, no. 2, pp. 240–251, Feb. 2015.
- [22] —, "Nested generalized sidelobe canceller for joint dereverberation and noise reduction," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP 2015)*, Brisbane, Australia, April 2015, pp. 106–110.
- [23] T. Dietzen, N. Huleihel, A. Spriet, W. Tirry, S. Doclo, M. Moonen, and T. van Waterschoot, "Speech dereverberation by data-dependent beamforming with signal pre-whitening," in *Proc. 23rd European Signal Process. Conf. (EUSIPCO 2015)*, Nice, France, Aug. 2015.
- [24] T. Dietzen, A. Spriet, W. Tirry, S. Doclo, M. Moonen, and T. van Waterschoot, "On the relation between data-dependent beamforming and multichannel linear prediction for dereverberation," in *Proc. AES 60th Int. Conf. Dereverberation and Reverberation of Audio, Music, and Speech (DREAMS)*, Leuven, Belgium, Jan. 2016, pp. 1–8.
- [25] I. Kodrasi and S. Doclo, "Joint dereverberation and noise reduction based on acoustic multi-channel equalization," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 24, no. 4, pp. 680–693, Apr. 2016.
- [26] M. Delcroix, T. Yoshioka, A. Ogawa, Y. Kubo, M. Fujimoto, N. Ito, K. Kinoshita, M. Espi, S. Araki, T. Hori, and T. Nakatani, "Strategies for distant speech recognition in reverberant environments," *EURASIP J. Adv. Signal Process.*, pp. 1–15, Dec. 2015.
- [27] T. Yoshioka, N. Ito, M. Delcroix, A. Ogawa, K. Kinoshita, M. Fujimoto, C. Yu, W. J. Fabian, M. Espi, T. Higuchi, S. Araki, and T. Nakatani, "The NTT CHiME-3 system: Advances in speech enhancement and recognition for mobile multi-microphone devices," in *Proc. 2015 IEEE Workshop Automatic Speech Recognition, Understanding (ASRU)*, Scottsdale, AZ, USA, Dec. 2015, pp. 436–443.
- [28] T. Dietzen, S. Doclo, M. Moonen, and T. van Waterschoot, "Joint multi-microphone speech dereverberation and noise reduction using integrated sidelobe cancellation and linear prediction," in *Proc. Int. Workshop Acoustic Signal Enhancement (IWAENC 2018)*, Tokyo, Japan, Sep. 2018.
- [29] J. Barker, S. Watanabe, E. Vincent, and J. Trmal, "The fifth 'CHiME' speech separation and recognition challenge: Dataset, task and baselines," in *Proc. Interspeech 2018*, Hyderabad, India, Sep. 2018, pp. 1561–1565.
- [30] P. A. Naylor and N. D. Gaubitch, *Speech Dereverberation*. Springer, July 2010.
- [31] L. J. Griffiths and C. W. Jim, "An alternative approach to linearly constrained adaptive beamforming," *IEEE Trans. Antennas Propag.*, vol. 1, no. 30, pp. 27–34, Jan. 1982.
- [32] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, "A consolidated perspective on multimicrophone speech enhancement and source separation," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 25, no. 4, pp. 692–730, Apr. 2017.
- [33] Y. Avargel and I. Cohen, "System identification in the short-time Fourier transform domain with crossband filtering," *IEEE Trans. Audio Speech Lang. Process.*, vol. 15, no. 4, pp. 1305–1319, Apr. 2007.
- [34] C. D. Meyer, Ed., *Matrix Analysis and Applied Linear Algebra*. Philadelphia, PA, USA: Society for Industrial and Applied Mathematics, 2000.
- [35] I. Kodrasi and S. Doclo, "EVD-based multi-channel dereverberation of a moving speaker using different RETF estimation methods," in *Proc. IEEE Hands-free Speech Com. Mic. Arrays (HSCMA 2017)*, San Francisco, CA, USA, Mar. 2017, pp. 116–120.
- [36] Y. Hu and P. C. Loizou, "Evaluation of objective quality measures for speech enhancement," *IEEE Trans. Audio Speech Lang. Process.*, vol. 16, no. 1, pp. 229–238, 2008.
- [37] D. S. E., N. Antonello, M. Moonen, and T. van Waterschoot, "On the modeling of rectangular geometries in room acoustic simulations," *IEEE Trans. Audio, Speech, Lang. Process.*, vol. 23, no. 4, pp. 774–786, Apr. 2015.
- [38] Bang and Olufsen, "Music for Archimedes," Compact Disc B&O, 1992.
- [39] S. Markovich, S. Gannot, and I. Cohen, "Multichannel eigenspace beamforming in a reverberant noisy environment with multiple interfering speech signals," *IEEE Trans. Audio Speech Lang. Process.*, vol. 17, no. 6, pp. 1071–1086, Aug. 2009.
- [40] T. Dietzen, "Audio examples for IEEE/ACM TASLP 2018," <ftp://ftp.esat.kuleuven.be/pub/SISTA/tdietzen/reports/taslp18/audio>, Aug. 2018.



**Thomas Dietzen** received his Dipl.-Ing. degree from Kaiserslautern University, Germany, in 2011. Between 2012 and 2014, he was a research assistant at University of Heidelberg, Germany, and at Fraunhofer Institute for Integrated Circuits IIS, Germany. From 2014 to 2017, he has been a doctoral researcher at NXP Semiconductors Belgium NV, Belgium, in the frame of the FP7-PEOPLE Marie Curie ITN 'Dereverberation and Reverberation of Audio, Music, and Speech (DREAMS)'. Currently, he pursues his PhD at KU Leuven, Belgium. His research is focused on signal processing algorithms for microphone-array based speech enhancement, specifically on dereverberation and noise reduction. He has served as a reviewer for the IEEE Transactions on Audio, Speech, and Language Processing and the IEEE Signal Processing Letters.



**Ann Spriet** received her MSc and PhD degree in Electrical Engineering from KU Leuven, Belgium, in 1999 and 2004, respectively. From 2004 to 2010, she worked as a postdoctoral researcher at KU Leuven, Belgium. During her PhD and postdoctoral research, she worked on signal processing algorithms for speech enhancement in hearing aids and cochlear implants. In 2003 and 2005, she received a best student paper award at the International Workshop on Acoustic Echo and Noise Control and the International Conference on Acoustics, Speech and Signal Processing, respectively. From 2008 to 2010, she was a nominated officer of EURASIP, where she was active as newsletter editor. Since 2010, she is a senior audio research engineer at NXP Semiconductors Belgium N.V. Her main activities are the development of speech enhancement algorithms for mobile devices and technology scouting. Since 2013, she is a member of the IEEE Audio and Acoustic Signal Processing Technical Committee.



ogy development activities.

**Wouter Tirry** received his M.Sc. degree in Physics and his PhD degree in Solar Physics from the University of Leuven, Belgium, in 1994 and 1998, respectively. As a post-doc, he further pursued his research at the National Centre for Atmospheric Research, Boulder, Colorado. Since 1999, he has been building up expertise in the domain of speech enhancement for mobile devices at Philips and NXP as research engineer and system architect. Currently he is Senior Principal at the Product Line Voice and Audio Solutions, NXP leading the speech technology development activities.





**Simon Doclo** (S'95-M'03-SM'13) received the M.Sc. degree in electrical engineering and the Ph.D. degree in applied sciences from the Katholieke Universiteit Leuven, Belgium, in 1997 and 2003. From 2003 to 2007 he was a Postdoctoral Fellow with the Research Foundation Flanders at the Electrical Engineering Department (Katholieke Universiteit Leuven) and the Cognitive Systems Laboratory (McMaster University, Canada). From 2007 to 2009 he was a Principal Scientist with NXP Semiconductors at the Sound and Acoustics Group in Leuven,

Belgium. Since 2009 he is a full professor at the University of Oldenburg, Germany, and scientific advisor for the project group Hearing, Speech and Audio Technology of the Fraunhofer Institute for Digital Media Technology. His research activities center around signal processing for acoustical and biomedical applications, more specifically microphone array processing, speech enhancement, active noise control, acoustic sensor networks and hearing aid processing. Prof. Doclo received the Master Thesis Award of the Royal Flemish Society of Engineers in 1997, the Best Student Paper Award at the International Workshop on Acoustic Echo and Noise Control in 2001, the EURASIP Signal Processing Best Paper Award in 2003 (with Marc Moonen) and the IEEE Signal Processing Society 2008 Best Paper Award (with Jingdong Chen, Jacob Benesty, Arden Huang). He is member of the IEEE Signal Processing Society Technical Committee on Audio and Acoustic Signal Processing, the EURASIP Special Area Team on Acoustic, Speech and Music Signal Processing and the EAA Technical Committee on Audio Signal Processing. Prof. Doclo was and is involved in several large-scale national and European research projects (ITN DREAMS, Cluster of Excellence Hearing4All, CRC Hearing Acoustics). He was Technical Program Chair of the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA) in 2013 and Chair of the ITG Conference on Speech Communication in 2018. In addition, he served as guest editor for several special issues (IEEE Signal Processing Magazine, Elsevier Signal Processing) and is associate editor for IEEE/ACM Transactions on Audio, Speech and Language Processing and EURASIP Journal on Advances in Signal Processing.



**Marc Moonen** (M'94, SM'06, F'07) is a Full Professor at the Electrical Engineering Department of KU Leuven, where he is heading a research team working in the area of numerical algorithms and signal processing for digital communications, wireless communications, DSL and audio signal processing. He is a Fellow of the IEEE (2007) and a Fellow of EURASIP (2018). He received the 1994 KU Leuven Research Council Award, the 1997 Alcatel Bell (Belgium) Award (with Piet Vandaele), the 2004 Alcatel Bell (Belgium) Award (with Raphael

Cendrillon), and was a 1997 Laureate of the Belgium Royal Academy of Science. He received journal best paper awards from the IEEE Transactions on Signal Processing (with Geert Leus and with Daniele Giacobello) and from Elsevier Signal Processing (with Simon Doclo). He was chairman of the IEEE Benelux Signal Processing Chapter (1998-2002), a member of the IEEE Signal Processing Society Technical Committee on Signal Processing for Communications, and President of EURASIP (European Association for Signal Processing, 2007-2008 and 2011-2012). He has served as Editor-in-Chief for the EURASIP Journal on Applied Signal Processing (2003-2005), Area Editor for Feature Articles in IEEE Signal Processing Magazine (2012-2014), and has been a member of the editorial board of Signal Processing, IEEE Transactions on Circuits and Systems II, IEEE Signal Processing Magazine, Integration-the VLSI Journal, EURASIP Journal on Wireless Communications and Networking and EURASIP Journal on Advances in Signal Processing.



**Toon van Waterschoot** (S'04, M'12) received MSc (2001) and PhD (2009) degrees in Electrical Engineering, both from KU Leuven, Belgium, where he is currently an Associate Professor and Consolidator Grantee of the European Research Council (ERC). He has previously also held teaching and research positions at Delft University of Technology in The Netherlands and the University of Lugano in Switzerland. His research interests are in signal processing, machine learning, and numerical optimization, applied to acoustic signal enhancement,

acoustic modeling, audio analysis, and audio reproduction. He has been serving as an Associate Editor for the Journal of the Audio Engineering Society and for the EURASIP Journal on Audio, Music, and Speech Processing, and as a Guest Editor for Elsevier Signal Processing. He is a Director of the European Association for Signal Processing (EURASIP), a Member of the IEEE Audio and Acoustic Signal Processing Technical Committee, and a Founding Member of the EAA Technical Committee in Audio Signal Processing. He was the General Chair of the 60th AES International Conference in Leuven, Belgium (2016), and has been serving on the Organizing Committee of the European Conference on Computational Optimization (EUCCO 2016) and the IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA 2017). He is a member of EURASIP, IEEE, ASA, and AES.