

Comparison of Parameter Estimation Methods for Single-Microphone Multi-Frame Wiener Filtering

Dörte Fischer, Klaus Brümman, Simon Doclo

Department of Medical Physics and Acoustics and Cluster of Excellence Hearing4All, University of Oldenburg, Germany

{doerte.fischer,klaus.bruemann,simon.doclo}@uni-oldenburg.de

Abstract—The multi-frame Wiener filter (MFWF) for single-microphone speech enhancement is able to exploit speech correlation across consecutive time-frames in the short-time Fourier transform (STFT) domain. To achieve a high speech correlation, typically an STFT with a high time-resolution but a low frequency-resolution is applied. The MFWF can be decomposed into a multi-frame minimum power distortionless response (MFMPDR) filter and a single-frame Wiener postfilter. To implement the MFWF using this decomposition, estimates of several parameters are required, namely the speech correlation vector, the noisy speech correlation matrix, and the power spectral densities at the output of the MFMPDR filter. Correlations can be estimated either *directly* in the low frequency-resolution STFT filterbank, *indirectly* by estimating periodograms in a high frequency-resolution filterbank and applying the Wiener-Khinchin theorem, or in a *combined* way. In this paper, we compare the performance of different estimators for the required parameters. Experimental results for different speech material, noise conditions, and signal-to-noise ratios show that using a combined estimator for the speech correlation vector yields the best results in terms of speech quality compared to existing direct and indirect estimators.

I. INTRODUCTION

Speech signals recorded in communication devices are frequently corrupted by undesired additive noise. To improve the speech quality, single-microphone noise reduction is often applied in the short-time Fourier transform (STFT) domain [1]. In contrast to single-frame approaches, where a (real-valued) gain is applied to each noisy STFT coefficient independently, multi-frame approaches aim to exploit speech correlation across consecutive time-frames [1]–[7]. To achieve a high speech correlation, typically an STFT with a high time-resolution but a low frequency-resolution is applied. In [1], [3], a complex-valued multi-frame Wiener filter (MFWF) was proposed which minimizes the mean square error (MSE) between the desired and estimated speech coefficients. In [4], [6], [8], it was reported that in practice the MFWF is very sensitive to estimation errors and more robust results can be obtained by decomposing the MFWF into a multi-frame minimum power distortionless response (MFMPDR) filter [3] and a single-frame Wiener postfilter [8].

The MFMPDR filter requires estimates of the speech correlation vector and noisy speech correlation matrix, where it was shown in [5] that the MFMPDR filter is more sensitive

to estimation errors in the speech correlation vector. In [4], the noisy speech correlation matrix and the speech correlation vector were estimated *directly* in the low frequency-resolution STFT filterbank, where for the estimation of the speech correlation vector a fixed (i.e., time- and frequency-independent) average noise correlation vector was assumed. In [6], it was proposed to estimate the noisy speech correlation matrix and speech correlation vector *indirectly* by first estimating the noisy speech and speech periodograms in a high frequency-resolution filterbank and then applying the Wiener-Khinchin theorem. In this paper, we compare the performance of different estimators for the required parameters of the MFMPDR filter with single-frame Wiener postfilter, i.e., the speech correlation vector, the noisy speech correlation matrix, and the power spectral densities (PSDs) at the output of the MFMPDR filter. In addition to existing direct and indirect estimators, we also propose *combined* estimators using both the low and high frequency-resolution filterbank.

Experimental results for different speech signals, noise types, and signal-to-noise ratios (SNRs) show that using a combined estimator for the speech correlation vector can improve speech quality compared to using a direct or indirect estimator. Furthermore, the results show that a Wiener post-filter improves the speech quality compared to the MFMPDR filter, but independently estimating the PSDs at the output of the MFMPDR does not yield an improvement.

II. MULTI-FRAME SIGNAL MODEL

By applying an STFT with frame length K and analysis window h_K to the noisy microphone signal $y[n]$ at time n , the noisy speech coefficient $Y[k, l]$ at time-frame l and frequency-bin k is obtained, i.e.,

$$Y[k, l] = S[k, l] + N[k, l], \quad k \in \left\{ -\frac{K}{2} + 1, -\frac{K}{2} + 2, \dots, \frac{K}{2} \right\}, \quad (1)$$

where $S[k, l]$ and $N[k, l]$ denote the speech and the noise coefficients, respectively. The noisy speech vector $\mathbf{y}[k, l]$ is defined by considering M consecutive time-frames, i.e.,

$$\mathbf{y}[k, l] = [Y[k, l], Y[k, l-1], \dots, Y[k, l-M+1]]^T, \quad (2)$$

where T denotes the transpose operator. Using (1), this vector can be written as

$$\mathbf{y}[k, l] = \mathbf{s}[k, l] + \mathbf{n}[k, l], \quad (3)$$

This work was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) - Projektnummer 352015383 – SFB 1330 B2 and Cluster of Excellence 1077 Hearing4all.

where the speech vector $\mathbf{s}[k, l]$ and the noise vector $\mathbf{n}[k, l]$ are defined similarly as in (2). In multi-frame approaches the speech coefficient $S[k, l]$ is estimated by applying an M -dimensional (complex-valued) finite impulse response filter $\mathbf{h}[k, l]$ to the noisy speech vector, i.e.,

$$\hat{S}[k, l] = \mathbf{h}^H[k, l] \mathbf{y}[k, l], \quad (4)$$

where H denotes the Hermitian operator. For conciseness, in the remainder of this paper the indices k and l will be omitted wherever possible.

Assuming that the speech and noise signals are uncorrelated, the $M \times M$ -dimensional noisy speech correlation matrix $\mathbf{R}_y = \mathbb{E}[\mathbf{y}\mathbf{y}^H]$, with $\mathbb{E}[\cdot]$ the expectation operator, is given by

$$\mathbf{R}_y = \mathbf{R}_s + \mathbf{R}_n, \quad (5)$$

where $\mathbf{R}_s = \mathbb{E}[\mathbf{s}\mathbf{s}^H]$ and $\mathbf{R}_n = \mathbb{E}[\mathbf{n}\mathbf{n}^H]$ denote the speech and noise correlation matrices, respectively.

Considering the speech correlation across time-frames, it was proposed in [3] to decompose the speech vector \mathbf{s} into a temporally correlated speech component \mathbf{x} and a temporally uncorrelated speech component \mathbf{x}' with respect to the speech coefficient S , i.e.,

$$\mathbf{s} = \mathbf{x} + \mathbf{x}' = \boldsymbol{\gamma}_s S + \mathbf{x}', \quad (6)$$

where $\boldsymbol{\gamma}_s$ denotes the normalized speech correlation vector, which is defined as

$$\boldsymbol{\gamma}_s = \frac{\mathbb{E}[\mathbf{s}S^*]}{\mathbb{E}[|S|^2]} = \frac{\mathbf{r}_s}{\phi_S}, \quad (7)$$

where $*$ denotes the complex-conjugate operator and \mathbf{r}_s denotes the speech correlation vector. Due to the normalization with the speech PSD $\phi_S = \mathbb{E}[|S|^2] = \mathbf{r}_s(1)$, the first element of $\boldsymbol{\gamma}_s$ is equal to 1.

Using (5), (6) and (7), the speech correlation matrix \mathbf{R}_s can be decomposed into the rank-1 correlation matrix $\mathbf{R}_x = \phi_S \boldsymbol{\gamma}_s \boldsymbol{\gamma}_s^H$ and the correlation matrix $\mathbf{R}_{x'} = \mathbb{E}[\mathbf{x}'\mathbf{x}'^H]$. Hence, the speech correlation vector \mathbf{r}_s and the speech PSD ϕ_S in (7) can be computed as

$$\mathbf{r}_s = \mathbf{R}_x \mathbf{e}, \quad \phi_S = \mathbf{e}^T \mathbf{R}_x \mathbf{e}, \quad (8)$$

with $\mathbf{e} = [1, 0, \dots, 0]^T$ an M -dimensional selection vector. Considering the uncorrelated speech component \mathbf{x}' in (6) as an interference, the *multi-frame signal model* is given by

$$\mathbf{y} = \boldsymbol{\gamma}_s S + \mathbf{u}, \quad (9)$$

where the undesired signal vector $\mathbf{u} = \mathbf{x}' + \mathbf{n}$ and the correlated speech vector $\boldsymbol{\gamma}_s S$ are uncorrelated. Using (9), the noisy speech correlation matrix can therefore be written as

$$\mathbf{R}_y = \phi_S \boldsymbol{\gamma}_s \boldsymbol{\gamma}_s^H + \mathbf{R}_u, \quad (10)$$

with the undesired correlation matrix $\mathbf{R}_u = \mathbf{R}_{x'} + \mathbf{R}_n$. Since $\mathbf{e}^T \mathbf{R}_{x'} \mathbf{e} = 0$, the undesired PSD $\phi_U = \mathbf{e}^T \mathbf{R}_u \mathbf{e} = \phi_N$.

Similarly to (7), the normalized noisy speech correlation vector $\boldsymbol{\gamma}_y$ and the normalized noise correlation vector $\boldsymbol{\gamma}_n$ are defined as

$$\boldsymbol{\gamma}_y = \frac{\mathbf{r}_y}{\phi_Y} = \frac{\mathbf{R}_y \mathbf{e}}{\mathbf{e}^T \mathbf{R}_y \mathbf{e}}, \quad \boldsymbol{\gamma}_n = \frac{\mathbf{r}_n}{\phi_N} = \frac{\mathbf{R}_n \mathbf{e}}{\mathbf{e}^T \mathbf{R}_n \mathbf{e}}. \quad (11)$$

Using (5) and (11), it can be easily shown that

$$\phi_Y \boldsymbol{\gamma}_y = \phi_S \boldsymbol{\gamma}_s + \phi_N \boldsymbol{\gamma}_n, \quad (12)$$

such that the normalized speech correlation vector can be obtained as

$$\boldsymbol{\gamma}_s = \frac{\phi_S + \phi_N}{\phi_S} \boldsymbol{\gamma}_y - \frac{\phi_N}{\phi_S} \boldsymbol{\gamma}_n. \quad (13)$$

III. MULTI-FRAME WIENER FILTER

The aim of the MFWF is to minimize the MSE between the desired speech coefficient S and the estimated speech coefficient \hat{S} in (4), i.e.,

$$\hat{\mathbf{h}} = \underset{\mathbf{h}}{\operatorname{argmin}} E\{|\mathbf{h}^H \mathbf{y} - S|^2\}. \quad (14)$$

Solving this optimization problem using the multi-frame signal model in (9) leads to the MFWF [3]

$$\mathbf{h}_{\text{MFWF}} = \mathbf{R}_y^{-1} \boldsymbol{\gamma}_s \phi_S. \quad (15)$$

In [4], [6], [8] it was reported that in practice the MFWF is very sensitive to estimation errors and more robust results can be obtained by decomposing the MFWF into an MFMPDR filter [3] and a single-frame Wiener postfilter [8], i.e.,

$$\mathbf{h}_{\text{MFMPDR-WG}} = \underbrace{\mathbf{R}_y^{-1} \boldsymbol{\gamma}_s}_{\mathbf{h}_{\text{MFMPDR}}} \underbrace{\frac{\phi_S}{\phi_Y^{\text{out}}}}_{G_{\text{WG}}}, \quad (16)$$

where $\phi_Y^{\text{out}} = (\boldsymbol{\gamma}_s^H \mathbf{R}_y^{-1} \boldsymbol{\gamma}_s)^{-1}$ denotes the signal PSD at the output of the MFMPDR filter.

As can be observed from (16), the MFWF requires estimates of the noisy speech correlation matrix \mathbf{R}_y , the normalized speech correlation vector $\boldsymbol{\gamma}_s$, the speech PSD ϕ_S and the signal output PSD ϕ_Y^{out} . In Section IV, we present an indirect method to estimate the correlation vectors \mathbf{r}_y , \mathbf{r}_s , and \mathbf{r}_n . In Section V, we present different estimators (direct, indirect, and combined) for the noisy speech correlation matrix, the normalized speech correlation vector, and the signal output PSD.

IV. INDIRECT CORRELATION VECTOR ESTIMATION

The Wiener-Khinchin theorem states that the correlation of a wide-sense stationary process is given by the inverse discrete Fourier transform (IDFT) of the PSD. To accurately estimate the correlation based on the Wiener-Khinchin theorem, it was proposed in [6] to estimate the periodogram in a filterbank with a $2M/O$ higher frequency-resolution than the (processing) STFT filterbank, where O denotes the oversampling factor. In the following, the parameters in the high frequency-resolution filterbank will be denoted with a superscript $\{\}^F$. An overview of the indirect estimation of the correlation vectors is depicted in Fig. 1.

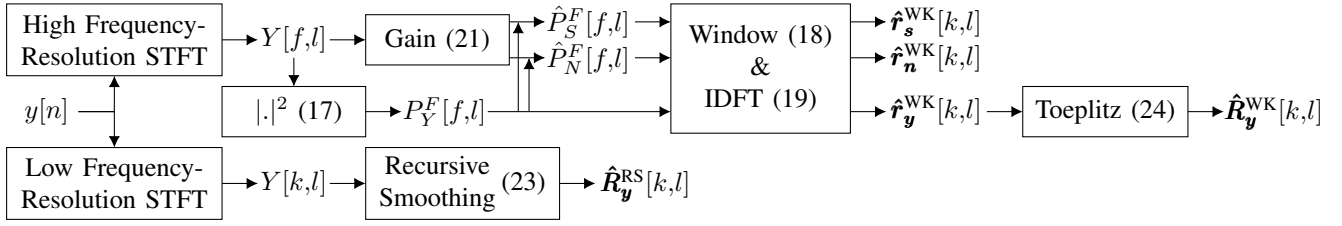


Fig. 1. Block diagram of parameter estimation from Sections IV and V-A.

In the high frequency-resolution filterbank, the noisy speech periodogram is given by

$$P_Y^F[f,l] = |Y^F[f,l]|^2, \quad f \in \left\{ -\frac{F}{2} + 1, -\frac{F}{2} + 2, \dots, \frac{F}{2} \right\}, \quad (17)$$

with \$Y^F[f,l]\$ the noisy speech coefficient at time-frame \$l\$ and frequency-bin \$f\$. An estimate of the \$2M\$-dimensional noisy speech PSD vector \$\hat{\phi}_Y[k,l]\$ in the low frequency-resolution filterbank is obtained by applying a windowing operation, i.e.,

$$\hat{\phi}_Y[k,l](\tau) = \frac{1}{O} |H_F[\tau]|^2 P_Y^F \left[\frac{2Mk}{O} + \tau, l \right], \quad \tau = -M+1, -M+2, \dots, M, \quad (18)$$

where \$\hat{\phi}_Y[k,l](\tau)\$ denotes the \$\tau\$-th element of \$\hat{\phi}_Y[k,l]\$ and \$H^F\$ denotes the Fourier transform of the zero-padded analysis window \$h_K\$ of length \$F = 2MK/O\$. Using the Wiener-Khinchin theorem, an estimate of the \$M\$-dimensional noisy speech correlation vector \$\hat{\mathbf{r}}_y^{\text{WK}}[k,l]\$ is obtained by applying the IDFT to \$\hat{\phi}_Y[k,l]\$ in (18), i.e.,

$$\hat{\mathbf{r}}_y^{\text{WK}}[k,l](m) = \frac{1}{2M} \sum_{\tau=-M+1}^M \hat{\phi}_Y[k,l](\tau) e^{-j2\pi\tau m/2M}, \quad m=0,1,\dots,M-1. \quad (19)$$

Similarly, it was proposed in [6] to estimate the speech correlation vector \$\mathbf{r}_s\$ using an estimate of the speech periodogram \$P_S^F\$ in the high frequency-resolution filterbank. The speech periodogram can be estimated by applying power subtraction to the noisy speech periodogram, i.e.,

$$\hat{P}_S^F = \hat{G}^F P_Y^F, \quad (20)$$

where the Wiener gain \$\hat{G}^F\$ is computed as

$$\hat{G}^F = \frac{\hat{\xi}^F}{\hat{\xi}^F + 1}, \quad (21)$$

with \$\hat{\xi}^F\$ denoting an estimate of the a-priori SNR. The a-priori SNR is estimated using the decision-directed approach (DDA) [9] with the noise PSD estimator in [10]. Replacing \$P_Y^F\$ with \$\hat{P}_S^F\$ in (18), followed by (19), yields an estimate of the speech correlation vector \$\hat{\mathbf{r}}_s^{\text{WK}}\$.

To estimate the noise correlation vector \$\mathbf{r}_n\$, we propose a similar approach, namely by estimating the noise periodogram \$P_N^F\$ in the high frequency-resolution filterbank as

$$\hat{P}_N^F = (1 - \hat{G}^F) P_Y^F. \quad (22)$$

Replacing \$P_Y^F\$ with \$\hat{P}_N^F\$ in (18), followed by (19), yields an estimate of the noise correlation vector \$\hat{\mathbf{r}}_n^{\text{WK}}\$.

V. PARAMETER ESTIMATION

In this section, we present several existing and novel estimators for the required parameters of the MFWF, specifically, estimating the parameters *directly* in the low frequency-resolution filterbank, *indirectly* using the estimated correlation vectors from Section IV or in a *combined* approach.

A. Noisy Speech Correlation Matrix

The noisy speech correlation matrix \$\mathbf{R}_y\$ can be estimated *directly* using first-order recursive smoothing as

$$\hat{\mathbf{R}}_y^{\text{RS}}[k,l] = \lambda \hat{\mathbf{R}}_y^{\text{RS}}[k,l-1] + (1-\lambda) \mathbf{y}[k,l] \mathbf{y}^H[k,l] \quad (23)$$

with \$\lambda\$ a forgetting factor.

Alternatively, since \$\mathbf{R}_y\$ can be assumed to be a Hermitian Toeplitz matrix, this matrix can be fully defined by the noisy speech correlation vector \$\mathbf{r}_y\$. Hence, using the estimate in (19), the noisy speech correlation matrix can be estimated *indirectly* as in [6]

$$\hat{\mathbf{R}}_y^{\text{WK}} = \text{Toeplitz}(\hat{\mathbf{r}}_y^{\text{WK}}) \quad (24)$$

B. Normalized Speech Correlation Vector

Based on (13), a maximum-likelihood estimator has been proposed in [4] to *directly* estimate the normalized speech correlation vector \$\gamma_s\$ by replacing the normalized noise correlation vector \$\gamma_n\$ with a long-term estimate \$\hat{\mu}_{\gamma_n}\$, i.e.,

$$\hat{\gamma}_s^{\diamond\text{-ML}} = \frac{\hat{\phi}_S + \hat{\phi}_N}{\hat{\phi}_S} \hat{\gamma}_y^{\diamond} - \frac{\hat{\phi}_N}{\hat{\phi}_S} \hat{\mu}_{\gamma_n} \quad (25)$$

where \$\hat{\mu}_{\gamma_n}\$ is defined by the STFT frame overlap and the analysis window \$h_K\$ [4]. The parameter \$\hat{\gamma}_y^{\diamond}\$ denotes an estimate of \$\gamma_y\$, computed similarly to (11) using \$\hat{\mathbf{R}}_y^{\diamond}\$, where \$\diamond\$ denotes either (23) or (24). The parameters \$\hat{\phi}_N\$ and \$\hat{\phi}_S\$ are estimates of the noise and speech PSDs in the low frequency-resolution filterbank, respectively. The noise PSD is estimated using [10] and the speech PSD is estimated similarly as the speech periodogram in the high frequency-resolution filterbank in (20), by applying a Wiener gain \$G\$ to the noisy PSD, i.e., \$\hat{\phi}_S = G e^T \hat{\mathbf{R}}_y^{\diamond} e\$.

Alternatively, it has been proposed in [6] to *indirectly* estimate the normalized speech correlation vector based on

the estimated speech correlation vector $\hat{\mathbf{r}}_s^{\text{WK}}$ from Section IV, i.e.,

$$\hat{\boldsymbol{\gamma}}_s^{\text{WK}} = \frac{\hat{\mathbf{r}}_s^{\text{WK}}}{\hat{\mathbf{r}}_s^{\text{WK}}(1)} \quad (26)$$

In addition, in this paper we propose a *combined* estimator by replacing the (time- and frequency-independent) long-term estimate $\hat{\mu}_{\gamma_n}$ by the (time- and frequency-dependent) estimated noise correlation vector $\hat{\mathbf{r}}_n^{\text{WK}}$ from Section IV, normalized similarly as in (26), i.e.,

$$\hat{\boldsymbol{\gamma}}_s^{\circ\text{-WK}} = \frac{\hat{\phi}_S + \hat{\phi}_N}{\hat{\phi}_S} \hat{\boldsymbol{\gamma}}_y^{\circ} - \frac{\hat{\phi}_N}{\hat{\phi}_S} \frac{\hat{\mathbf{r}}_n^{\text{WK}}}{\hat{\mathbf{r}}_n^{\text{WK}}(1)} \quad (27)$$

C. Signal Output PSD

The MFMPDR filter is designed to prevent speech distortion such that in theory, the speech PSD at the output of the MFMPDR filter is equal to the input speech PSD. Using this assumption, the speech PSD estimate $\hat{\phi}_S$ can be applied at the output of the MFMPDR filter.

The signal PSD ϕ_Y^{out} at the output of the MFMPDR filter can be estimated in several ways. The most straightforward way is to estimate this PSD as

$$\hat{\phi}_Y^{\text{out}, \mathbf{R}_y} = \left(\hat{\boldsymbol{\gamma}}_s^H \hat{\mathbf{R}}_y^{-1} \hat{\boldsymbol{\gamma}}_s \right)^{-1} \quad (28)$$

using the previously presented estimators for $\hat{\boldsymbol{\gamma}}_s$ and $\hat{\mathbf{R}}_y$.

Alternatively, by substituting (10) into (28) and applying the matrix inversion lemma, the signal output PSD can be estimated as

$$\hat{\phi}_Y^{\text{out}, \mathbf{R}_u} = \hat{\phi}_S + \left(\hat{\boldsymbol{\gamma}}_s^H \hat{\mathbf{R}}_u^{-1} \hat{\boldsymbol{\gamma}}_s \right)^{-1} \quad (29)$$

with $\left(\hat{\boldsymbol{\gamma}}_s^H \hat{\mathbf{R}}_u^{-1} \hat{\boldsymbol{\gamma}}_s \right)^{-1}$ an estimate of the undesired output PSD ϕ_U^{out} . According to (10), the correlation matrix $\hat{\mathbf{R}}_u$ can be determined as $\hat{\mathbf{R}}_u = \hat{\mathbf{R}}_y - \hat{\phi}_S \hat{\boldsymbol{\gamma}}_s \hat{\boldsymbol{\gamma}}_s^H$. Since $\hat{\mathbf{R}}_u$ may not be positive semi-definite due to estimation errors, we set the negative eigenvalues of $\hat{\mathbf{R}}_u$ to zero.

The undesired output PSD can also be estimated independently of the previously estimated parameters. Since $\phi_U = \phi_N$, we can replace $\hat{\phi}_U^{\text{out}}$ in (29) with an estimate of the noise output PSD $\hat{\phi}_N^{\text{out}}$, i.e.,

$$\hat{\phi}_Y^{\text{out}, \phi_N} = \hat{\phi}_S + \hat{\phi}_N^{\text{out}} \quad (30)$$

where $\hat{\phi}_N^{\text{out}}$ is determined at the output of the MFMPDR filter using [10].

VI. EXPERIMENTAL RESULTS

In this section, we compare the performance of the different presented estimators for the parameters of the MFWF. In Section VI-A we describe the used speech and noise material and the algorithmic implementation details. In Section VI-B we compare the performance of different MFMPDR filters. Using the best MFMPDR filter, in Section VI-C we compare the performance of different Wiener postfilters.

TABLE I
OVERVIEW OF THE EVALUATED PARAMETER ESTIMATORS IN THE MFMPDR FILTER.

| Label | Estimation of \mathbf{R}_y | Estimation of $\boldsymbol{\gamma}_s$ |
|---|------------------------------|---|
| $\hat{\mathbf{R}}_y^{\text{RS}} - \hat{\boldsymbol{\gamma}}_s^{\text{RS-ML}}$ | (23) | $\hat{\boldsymbol{\gamma}}_s^{\circ\text{-ML}}$: (25) using $\hat{\boldsymbol{\gamma}}_y^{\text{RS}}$: (23) |
| $\hat{\mathbf{R}}_y^{\text{RS}} - \hat{\boldsymbol{\gamma}}_s^{\text{WK}}$ | (23) | $\hat{\boldsymbol{\gamma}}_s^{\text{WK}}$: (26) |
| $\hat{\mathbf{R}}_y^{\text{RS}} - \hat{\boldsymbol{\gamma}}_s^{\text{RS-WK}}$ | (23) | $\hat{\boldsymbol{\gamma}}_s^{\circ\text{-WK}}$: (27) using $\hat{\boldsymbol{\gamma}}_y^{\text{RS}}$: (23) |
| $\hat{\mathbf{R}}_y^{\text{WK}} - \hat{\boldsymbol{\gamma}}_s^{\text{WK-ML}}$ | (24) | $\hat{\boldsymbol{\gamma}}_s^{\circ\text{-ML}}$: (25) using $\hat{\boldsymbol{\gamma}}_y^{\text{WK}}$: (24) |
| $\hat{\mathbf{R}}_y^{\text{WK}} - \hat{\boldsymbol{\gamma}}_s^{\text{WK}}$ | (24) | $\hat{\boldsymbol{\gamma}}_s^{\text{WK}}$: (26) |
| $\hat{\mathbf{R}}_y^{\text{WK}} - \hat{\boldsymbol{\gamma}}_s^{\text{WK-WK}}$ | (24) | $\hat{\boldsymbol{\gamma}}_s^{\circ\text{-WK}}$: (27) using $\hat{\boldsymbol{\gamma}}_y^{\text{WK}}$: (24) |

A. Algorithm Implementation and Performance Measures

The performance is evaluated in terms of the perceptual evaluation of speech quality (PESQ) [11] improvement over the noisy speech signal, using the clean speech signal as the reference signal. We used audio material from [12] sampled at 16 kHz and we evaluated the average performance over 176 s of speech material (92 s female, 84 s male) under five different noise conditions (babble, white Gaussian noise (WGN), traffic, modulated WGN, crossroad) at three different input SNRs (0, 5, and 10 dB).

To achieve a high speech correlation, we use a frame length of $K = 64$ samples (4 ms) and a frame shift of 16 samples (1 ms) in the low frequency-resolution STFT filterbank. As analysis and synthesis window h_K we use a Hann window. In the high frequency-resolution STFT filterbank, we use a four-times higher frequency-resolution, i.e. a frame length of $F = 256$ samples (16 ms), a frame shift of 16 samples (1 ms), and apply an asymmetric analysis window similarly to [6]. The number of the consecutive time-frames is $M = 8$ as in [6], resulting in 11 ms of analysis data in the low frequency-resolution filterbank. In the high and low frequency-resolution filterbanks, the weighting parameter for the DDA [9] is set to 0.97 and to reduce the amount of musical noise the Wiener gain is limited to -17 dB. The forgetting factor in (23) is experimentally set to $\lambda = 0.9$, resulting in a smoothing window of 10 ms. Before computing the inverse of $\hat{\mathbf{R}}_y$ estimated using (23) or (24) regularization based on diagonal loading is performed with a regularization parameter of 0.04 as in [4].

B. Comparison of MFMPDR Filters

In this section, we compare the performance of the MFMPDR filter for different direct, indirect, and combined estimators of the noisy speech correlation matrix \mathbf{R}_y and the normalized speech correlation vector $\boldsymbol{\gamma}_s$ (see Table I).

For different SNRs, Fig. 2 depicts the average PESQ improvements of the considered MFMPDR filters. Using the direct estimate $\hat{\mathbf{R}}_y^{\text{RS}}$ in (23) leads to a higher PESQ improvement than using the indirect estimate $\hat{\mathbf{R}}_y^{\text{WK}}$ in (24). In combination with either $\hat{\mathbf{R}}_y^{\text{RS}}$ or $\hat{\mathbf{R}}_y^{\text{WK}}$, using the indirect estimate of the normalized speech correlation vector $\hat{\boldsymbol{\gamma}}_s^{\circ\text{-WK}}$ in (26) leads to a degradation (or minor improvement) of the PESQ scores. The best PESQ improvements are obtained by combining the direct estimate $\hat{\mathbf{R}}_y^{\text{RS}}$ in (23) with the proposed combined estimate $\hat{\boldsymbol{\gamma}}_s^{\text{RS-WK}}$ in (27). This suggest that using the (time- and

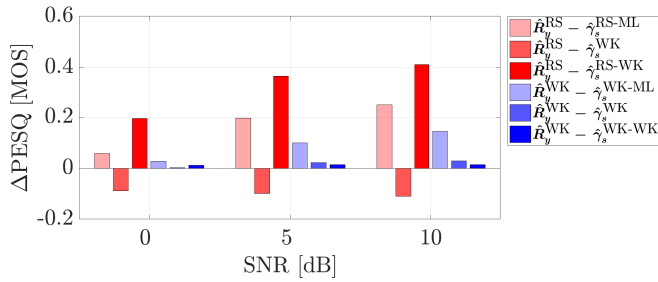


Fig. 2. Average PESQ improvement of the MFMPDR filters using the parameters in Table I at different SNRs.

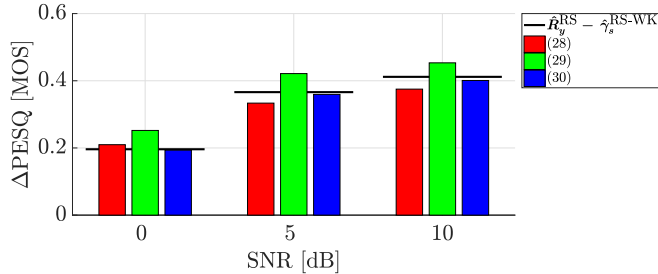


Fig. 3. Average PESQ improvement of the postfilters using the parameters from Section V-C at different SNRs, in reference to the MFMPDR filter $\hat{R}_y^{RS} - \hat{\gamma}_s^{RS-WK}$.

frequency-dependent) indirect estimate \hat{r}_n^{WK} is more effective to estimate the normalized speech correlation vector than using the (time- and frequency-independent) long-term estimate $\hat{\mu}_{\gamma_n}$.

C. Comparison of Single-Frame Wiener Postfilters

In this section we compare the performance of different Wiener postfilters using the signal output PSD estimators from Section V-C at the output of the best MFMPDR filter, i.e. using \hat{R}_y^{RS} and $\hat{\gamma}_s^{RS-WK}$. For different SNRs, Fig. 3 depicts the average PESQ improvements of the different parameter estimators in the Wiener postfilter, in reference to the MFMPDR filter. From these results it can be observed that the best PESQ improvements are obtained using the signal output PSD estimate $\hat{\phi}_Y^{out, R_u}$ in (29). It can also be observed that independently estimating the undesired PSD at the output of the MFMPDR filter, i.e., using $\hat{\phi}_Y^{out, \phi_N}$ in (30), does not yield a PESQ improvement compared to the MFMPDR filter. This can presumably be explained by the used noise PSD estimator [10], which assumes that the noise is more stationary than the speech, whereas the output of the MFMPDR filter may contain highly fluctuating residual noise leading to an underestimation of the noise output PSD.

VII. CONCLUSION

In this paper, we compare the noise reduction performance of different estimators for the required parameters of the single-microphone MFMPDR filter with single-frame Wiener postfilter, i.e., the noisy speech correlation matrix, the speech correlation vector and the signal output PSD. Existing estimation methods estimate the required parameters either directly

in the low frequency-resolution STFT filterbank or indirectly by estimating periodograms in a high frequency-resolution filterbank and applying the Wiener-Khinchin theorem. We compare these existing estimators with proposed combined estimators using both the low and the high frequency-resolution filterbank. The results show that combining the direct estimate of the noisy speech correlation matrix with the proposed combined estimate of the speech correlation vector achieves the best speech quality improvement. Furthermore, the results show that with the Wiener postfilter the speech quality can be improved in reference to the MFMPDR filter. However, when estimating the signal output PSD independently to the output of the MFMPDR filter, no speech quality improvement can be predicted.

REFERENCES

- [1] J. Benesty, J. Chen, and E. A. P. Habets, *Speech enhancement in the STFT domain*. Springer Science & Business Media, 2011.
- [2] T. Esch and P. Vary, "Modified Kalman Filter Exploiting Interframe Correlation of Speech and Noise Magnitudes," in *Proc. of Int. Workshop Acoustic Echo, Noise Control (IWAENC)*, Seattle, WA, USA, Sep. 2008.
- [3] Y. Huang and J. Benesty, "A multi-frame approach to the frequency-domain single-channel noise reduction problem," *IEEE Trans. Audio, Speech, Language Process.*, vol. 20, no. 4, pp. 1256–1269, May 2012.
- [4] A. Schasse and R. Martin, "Estimation of subband speech correlations for noise reduction via MVDR processing," *IEEE Trans. Audio, Speech, Language Process.*, vol. 22, no. 9, pp. 1355–1365, Sep. 2014.
- [5] D. Fischer and S. Doclo, "Sensitivity analysis of the multi-frame MVDR filter for single-microphone speech enhancement," in *Proc. of Europ. Signal Process. Conf. (EUSIPCO)*, Kos, Greece, Aug. 2017, pp. 603–607.
- [6] K. T. Andersen and M. Moonen, "Robust speech-distortion weighted interframe Wiener filters for single-channel noise reduction," *IEEE Trans. Audio, Speech, Language Process.*, vol. 26, no. 1, pp. 97–107, Jan. 2018.
- [7] D. Fischer and S. Doclo, "Robust constrained MFMPDR filtering for single-microphone speech enhancement," in *Proc. of Int. Workshop Acoustic Echo, Noise Control (IWAENC)*, Tokyo, Japan, Sep. 2018, pp. 41–45.
- [8] D. Fischer and T. Gerkmann, "Single-microphone speech enhancement using MVDR filtering and Wiener post-filtering," in *Proc. IEEE Int. Conf. Acoust., Speech, Signal Process. (ICASSP)*, Shanghai, China, Mar. 2016, pp. 201–205.
- [9] Y. Ephraim and D. Malah, "Speech enhancement using a minimum-mean square error short-time spectral amplitude estimator," *IEEE Trans. Acoust., Speech, Signal Process.*, vol. 32, no. 6, pp. 1109–1121, Dec. 1984.
- [10] T. Gerkmann and R. C. Hendriks, "Unbiased MMSE-based noise power estimation with low complexity and low tracking delay," *IEEE Trans. Audio, Speech, Language Process.*, vol. 20, no. 4, pp. 1383–1393, May 2012.
- [11] "ITU-T recommendation P.862. Perceptual evaluation of speech quality (PESQ): an objective method for end-to-end speech quality assessment of narrow-band telephone networks and speech codecs," Feb. 2001.
- [12] J. S. Garofolo, "DARPA TIMIT acoustic-phonetic speech database," in *National Institute of Standards and Technology (NIST)*, 1988.