

# Multi-Channel Post-Filtering in Non-Stationary Noise Environments

Israel Cohen

*Department of Electrical Engineering, Technion — Israel Institute of Technology,  
Technion City, Haifa 32000, Israel  
E-mail: icohen@ee.technion.ac.il; Tel.: +972 4 8294731; Fax: +972 4 8323041.*

## Abstract

In this paper, we present a multi-channel post-filtering approach for minimizing the log-spectral amplitude distortion in non-stationary noise environments. The beamformer is realistically assumed to have a steering error, a blocking matrix that is unable to block all of the desired signal components, and a noise canceller that is adapted to the pseudo-stationary noise, but not modified during transient interferences. A mild assumption is made, that a desired signal component is stronger at the beamformer output than at any reference noise signal, and a noise component is strongest at one of the reference signals. The ratio between the transient power at the beamformer output and the transient power at the reference noise signals is used for indicating whether such a transient is desired or interfering. Based on a Gaussian statistical model and combined with an appropriate spectral enhancement technique, we derive estimators for the signal presence probability, the noise power spectral density, and the clean signal. The proposed method is tested in various non-stationary noise environments. Compared to single-channel post-filtering, a significantly reduced level of non-stationary noise is achieved without further distorting the desired signal components.

## Keywords

Array signal processing, signal detection, acoustic noise measurement, speech enhancement, spectral analysis, adaptive signal processing.

## I. INTRODUCTION

MULTI channel systems are often used for high quality hands-free communication in reverberant and noisy environments [1]. Compared to single channel systems, a substantial gain in performance is obtainable due to the spatial filtering capability to suppress interfering signals coming from undesired directions. However, in cases of spatially incoherent noise fields, beamforming alone does not provide sufficient noise reduction, and post-filtering is normally required [2], [3].

Multi-channel post-filtering, generalized to an arbitrary number of sensors, was first introduced by Zelinski [4], [5]. Accordingly, the output of a delay-and-sum beamformer is post-filtered using an adaptive Wiener filtering in the time domain, based on the auto and cross spectral densities of the sensor signals. However, Zelinski's approach overestimates the noise power density, and therefore is not optimal in the Wiener sense [6]. A modified post-filtering version was suggested by Simmer and Wasiljeff, which employs the power spectral density of the beamformer output, rather than the average of the power spectral densities of individual sensor signals [6]. The underlying assumption is that noise components at different sensors are mutually uncorrelated. Unfortunately, in a diffuse noise field, where the low-frequency noise components are coherent, the noise reduction performance severely deteriorates.

To overcome this problem, Fischer *et al.* [7], [8], [9] proposed a noise reduction system, which is based on the generalized sidelobe canceller (GSC). The GSC reasonably suppresses the coherent noise components, while a Wiener filter in the look direction is designed to suppress the spatially incoherent noise components. Bitzer *et al.* analyzed the performance of the GSC and adaptive post-filtering techniques in various noise fields [10], [11]. They showed that in a diffuse noise field, neither the GSC nor the adaptive post-filtering performs well at low frequencies. Therefore, at the output of a GSC with standard Wiener post-filtering they used a second post-filter to reduce the spatially correlated noise components [12], [13]. Le Bouquin-Jeannès *et al.* suggested to modify the cross power spectrum estimation and the Wiener post-filtering to take the presence of some correlated noise components into account [14]. The cross power spectrum of the noise signals is averaged during pauses in the desired signal. Subsequently, it is subtracted from the cross power

spectrum of the sensor signals, calculated during signal presence. Meyer and Simmer [15] proposed to combine a delay-and-sum beamformer with Wiener filtering and spectral subtraction. The Wiener filtering is applied in the high-frequency band for the suppression of low-coherence noise components, while the spectral subtraction is used in the low-frequency band for high-coherence noise reduction. Mamhoudi [16] and Mamhoudi and Drygajlo [17] considered a nonlinear coherence filtering in the wavelet domain to improve the performance of the Wiener post-filtering. Instead of the conventional coherence between the individual sensor signals, they used the coherence between the output and the input of the beamformer sensor signals, which is assumed to be low even for correlated noise components. Fischer and Kameyer [18] suggested to apply Wiener filtering to the output of a broadband beamformer, that is built up by several harmonically nested subarrays. They showed that the resulting noise reduction system performance is nearly independent of the correlation properties of the noise field. This structure has been further analyzed by Marro *et al.* [2]. McCowan *et al.* used a near-field super-directive beamforming and investigated the effect of a Wiener post-filter on speech recognition performance [19]. They showed that in the case of nearfield sources and diffuse noise conditions, improved recognition performance can be achieved compared to conventional adaptive beamformers. A theoretical analysis of Wiener multi-channel post-filtering is presented in [3].

A major drawback of existing multi-channel post-filtering techniques is that highly non-stationary noise components are not dealt with. The time variation of the interfering signals is assumed to be sufficiently slow, such that the post-filter can track and adapt to the changes in the noise statistics. Unfortunately, transient interferences are often much too brief and abrupt for the above post-filtering methods. Furthermore, Wiener filtering minimizes the mean-square error (MSE) distortion of the signal estimate, which is essentially not the optimal criterion for enhancing noisy speech. A more appropriate distortion measure for speech enhancement systems is based on the MSE of the spectral, or log-spectral, amplitude [20], [21].

In this paper, we present a multi-channel post-filtering approach for minimizing the log-spectral amplitude distortion in non-stationary noise environments. Presumably, a desired signal component is stronger at the beamformer output than at any reference noise signal, and a noise component is strongest at one of the reference signals. Hence, the ratio between the transient power at beam-

former output and the transient power at the reference signals indicates whether such a transient is desired or interfering. Based on a Gaussian statistical model [20], and an appropriate decision-directed *a priori* SNR estimate [22], we derive an estimator for the signal presence probability. This estimator controls the rate of recursive averaging for obtaining a noise spectrum estimate by the *Minima Controlled Recursive Averaging* (MCRA) approach [22], [23]. Subsequently, spectral enhancement of the beamformer output is achieved by applying an optimal gain function, which minimizes the MSE of the log-spectra. The performance of the proposed post-filtering approach is evaluated under non-stationary noise conditions using objective quality measures, a subjective study of speech spectrograms and informal listening tests. We show that single-channel post-filtering is inefficient at attenuating highly non-stationary noise components, since it lacks the ability to differentiate such components from the desired source components. By contrast, the proposed multi-channel post-filtering approach achieves a significantly reduced level of background noise, whether stationary or not, without further distorting the signal components.

The paper is organized as follows. In Section II, we review the linearly constrained adaptive beamformer, and derive relations in the power-spectral domain between the beamformer output, the reference noise signals, the desired source signal, and the input transient interferences. In Section III, the problem of signal detection in the time-frequency plane is addressed. Signal components are discriminated from transient noise components based on the transient power ratio between the beamformer output and the reference signals. In Section IV, we introduce an estimator for the time-varying spectrum of the beamformer output noise, and describe the multi-channel post-filtering approach. Finally, in Section V, we evaluate the proposed method, and present experimental results, which validate its effectiveness.

## II. LINEARLY CONSTRAINED ADAPTIVE BEAMFORMING

Let  $x(t)$  denote a desired source signal, and let signal vectors  $\mathbf{d}_s(t)$  and  $\mathbf{d}_t(t)$  denote multi-channel uncorrelated interfering signals at the output of  $M$  sensors. The vector  $\mathbf{d}_s(t)$  represents pseudo-stationary interferences, and  $\mathbf{d}_t(t)$  represents undesired transient components. The observed signal at the  $i$ -th sensor is given by

$$z_i(t) = a_i(t) * x(t) + d_{is}(t) + d_{it}(t), \quad i = 1, \dots, M \quad (1)$$

where  $a_i(t)$  is the transfer function from the desired source to the  $i$ -th sensor,  $*$  denotes convolution, and  $d_{is}$  and  $d_{it}$  are the interference signals corresponding to the  $i$ -th sensor. The observed signals are divided in time into overlapping frames by the application of a window function and analyzed using the short-time Fourier transform (STFT). Assuming time-invariant transfer functions [24], we have in the time-frequency domain

$$\mathbf{Z}(k, \ell) = \mathbf{A}(k)X(k, \ell) + \mathbf{D}_s(k, \ell) + \mathbf{D}_t(k, \ell) \quad (2)$$

where  $k$  represents the frequency bin index,  $\ell$  the frame index, and

$$\begin{aligned} \mathbf{Z}(k, \ell) &\triangleq [Z_1(k, \ell) \quad Z_2(k, \ell) \quad \cdots \quad Z_M(k, \ell)]^T \\ \mathbf{A}(k) &\triangleq [A_1(k) \quad A_2(k) \quad \cdots \quad A_M(k)]^T \\ \mathbf{D}_s(k, \ell) &\triangleq [D_{1s}(k, \ell) \quad D_{2s}(k, \ell) \quad \cdots \quad D_{Ms}(k, \ell)]^T \\ \mathbf{D}_t(k, \ell) &\triangleq [D_{1t}(k, \ell) \quad D_{2t}(k, \ell) \quad \cdots \quad D_{Mt}(k, \ell)]^T. \end{aligned}$$

We note that in [24], transient interferences are not dealt with; The interfering signals are assumed to be stationary, and signal enhancement is based on the non-stationarity of the desired source signal. In our case, we have to include a mechanism that discriminates interfering transients from desired signal components.

Fig. 1 shows a generalized sidelobe canceller structure for a linearly constrained adaptive beamformer [25], [26], which is also utilizable in the case of arbitrary transfer functions [24]. The beamformer comprises three parts: 1) A fixed beamformer  $\mathbf{W}$ , proportional to the transfer function ratios  $A_1^{-1}\mathbf{A}$ ; 2) A blocking matrix  $\mathbf{B}$ , which takes into account the assumed propagation path and constructs the reference noise signals  $\{U_i : 2 \leq i \leq M\}$ ; 3) A multi-channel adaptive noise canceller  $\{H_i : 2 \leq i \leq M\}$ , which eliminates the stationary noise that leaks through the sidelobes of the fixed beamformer. We assume that the noise canceller is adapted only to the stationary noise. It is not modified during transient interferences, which are characterized by brief and abrupt variations. Furthermore, we assume that the source is distributed and that steering error might occur. Accordingly, some desired signal components may pass through the blocking matrix.

The reference noise signals  $\mathbf{U}(k, \ell) = [U_2(k, \ell) \quad U_3(k, \ell) \quad \cdots \quad U_M(k, \ell)]^T$  are generated by ap-

plying the blocking matrix to the observed signal vector:

$$\begin{aligned}\mathbf{U}(k, \ell) &= \mathbf{B}^H(k) \mathbf{Z}(k, \ell) \\ &= \mathbf{B}^H(k) [\mathbf{A}(k) X(k, \ell) + \mathbf{D}_s(k, \ell) + \mathbf{D}_t(k, \ell)] .\end{aligned}\quad (3)$$

The reference signals are emphasized by the adaptive noise canceller and subtracted from the output of the fixed beamformer, yielding

$$Y(k, \ell) = [\mathbf{W}^H(k) - \mathbf{H}^H(k, \ell) \mathbf{B}^H(k)] \mathbf{Z}(k, \ell), \quad (4)$$

where  $\mathbf{H}(k, \ell) = [H_2(k, \ell) \ H_3(k, \ell) \ \cdots \ H_M(k, \ell)]^T$ . The optimal solution for the filters  $\mathbf{H}(k, \ell)$  is obtained by minimizing the output power of the stationary noise [27]. Let  $\Phi_{\mathbf{D}_s \mathbf{D}_s}(k, \ell) = E \{ \mathbf{D}_s(k, \ell) \mathbf{D}_s^H(k, \ell) \}$  denote the power-spectral density (PSD) matrix of the input stationary noise. Then, the power of the stationary noise at the beamformer output is minimized by solving the unconstrained optimization problem:

$$\min_{\mathbf{H}} \left\{ [\mathbf{W}(k) - \mathbf{B}(k) \mathbf{H}(k, \ell)]^H \Phi_{\mathbf{D}_s \mathbf{D}_s}(k, \ell) [\mathbf{W}(k) - \mathbf{B}(k) \mathbf{H}(k, \ell)] \right\}. \quad (5)$$

The multi-channel Wiener solution is given by [28]

$$\mathbf{H}(k, \ell) = [\mathbf{B}^H(k) \Phi_{\mathbf{D}_s \mathbf{D}_s}(k, \ell) \mathbf{B}(k)]^{-1} \mathbf{B}^H(k) \Phi_{\mathbf{D}_s \mathbf{D}_s}(k, \ell) \mathbf{W}(k). \quad (6)$$

If we assume that the stationary, as well as transient, noise fields are homogeneous, then the PSD-matrices of the input noise signals are related to the corresponding spatial coherence matrices,  $\Gamma_s(k, \ell)$  and  $\Gamma_t(k, \ell)$ , by

$$\begin{aligned}\Phi_{\mathbf{D}_s \mathbf{D}_s}(k, \ell) &= \lambda_s(k, \ell) \Gamma_s(k, \ell) \\ \Phi_{\mathbf{D}_t \mathbf{D}_t}(k, \ell) &= \lambda_t(k, \ell) \Gamma_t(k, \ell)\end{aligned}$$

where  $\lambda_s(k, \ell)$  and  $\lambda_t(k, \ell)$  represent the input noise power at a single sensor. The input PSD-matrix is therefore given by

$$\Phi_{\mathbf{Z} \mathbf{Z}}(k, \ell) = \lambda_x(k, \ell) \mathbf{A}(k) \mathbf{A}^H(k) + \lambda_s(k, \ell) \Gamma_s(k, \ell) + \lambda_t(k, \ell) \Gamma_t(k, \ell) \quad (7)$$

where  $\lambda_x(k, \ell) \triangleq E\{|X(k, \ell)|^2\}$  is the PSD of the desired source signal. Using (3) and (4), the PSD-matrix of the reference signals and the PSD of the beamformer output are obtained by

$$\Phi_{\mathbf{U}\mathbf{U}}(k, \ell) = \mathbf{B}^H(k) \Phi_{\mathbf{Z}\mathbf{Z}}(k, \ell) \mathbf{B}(k) \quad (8)$$

$$\phi_{YY}(k, \ell) = [\mathbf{W}(k) - \mathbf{B}(k)\mathbf{H}(k, \ell)]^H \Phi_{\mathbf{Z}\mathbf{Z}}(k, \ell) [\mathbf{W}(k) - \mathbf{B}(k)\mathbf{H}(k, \ell)]. \quad (9)$$

Substituting (7) into (8) and (9), we have the following linear relation between the PSD's of the beamformer output, the reference signals, the desired source signal, and the input interferences:

$$\begin{bmatrix} \phi_{YY}(k, \ell) \\ \phi_{U_2 U_2}(k, \ell) \\ \vdots \\ \phi_{U_M U_M}(k, \ell) \end{bmatrix} = \begin{bmatrix} C_{11}(k, \ell) & C_{12}(k, \ell) & C_{13}(k, \ell) \\ \vdots & \vdots & \vdots \\ C_{M1}(k, \ell) & C_{M2}(k, \ell) & C_{M3}(k, \ell) \end{bmatrix} \begin{bmatrix} \lambda_x(k, \ell) \\ \lambda_s(k, \ell) \\ \lambda_t(k, \ell) \end{bmatrix} \quad (10)$$

where

$$[C_{11} \ C_{12} \ C_{13}] = [\mathbf{W} - \mathbf{B}\mathbf{H}]^H [\mathbf{A}\mathbf{A}^H \ \Gamma_s \ \Gamma_t] (\mathbf{I}_3 \otimes [\mathbf{W} - \mathbf{B}\mathbf{H}]) \quad (11)$$

$$[C_{21} \ \cdots \ C_{M1}] = \text{diag}\{\mathbf{B}^H \mathbf{A} \mathbf{A}^H \mathbf{B}\} \quad (12)$$

$$[C_{22} \ \cdots \ C_{M2}] = \text{diag}\{\mathbf{B}^H \Gamma_s \mathbf{B}\} \quad (13)$$

$$[C_{23} \ \cdots \ C_{M3}] = \text{diag}\{\mathbf{B}^H \Gamma_t \mathbf{B}\}, \quad (14)$$

$\mathbf{I}_3$  is a 3-by-3 identity matrix,  $\otimes$  denotes Kronecker product, and  $\text{diag}\{\cdot\}$  represents a row vector constructed from the diagonal of a square matrix.

### III. DETECTION OF SOURCE SIGNALS IN NON-STATIONARY NOISE

Generally, the beamformer output comprises three components. Substituting (2) into (4), we have a non-stationary desired source component, a pseudo-stationary noise component, and a transient interference. Our objective is to detect the desired signal components at the beamformer output, and to differentiate them from the transient interference components.

We assume that the beamformer output and reference noise signals are obtained by adaptively aiming the beamformer at the desired source. Presumably, a desired signal component is stronger at the beamformer output than at any reference noise signal, and a noise component is stronger at one of the reference signals than at the beamformer output. Hence the *transient beam-to-reference*

*ratio* (TBRR), defined by the ratio between the transient power at beamformer output and the transient power at the reference signals, indicates whether such a transient is desired or interfering.

First, we detect transients at the beamformer output. Then, if there are no simultaneous transients at the reference signals, we determine that these transients are likely source components. In that case, a cautious enhancement would be involved. On the other hand, if a simultaneous transient at one of the reference signals is detected, then the TBRR would determine the extent to which such a transient is suppressed or preserved.

#### A. Detection of transients at the beamformer output

Let  $\mathcal{S}$  be a smoothing operator in the power spectral domain, and let  $\mathcal{M}$  denote a single-channel estimator for the PSD of the background pseudo-stationary noise. For example, a causal  $\mathcal{S}$  may be defined by recursively averaging past spectral power values of the noisy measurement:

$$\mathcal{S}Y(k, \ell) = \alpha_s \cdot \mathcal{S}Y(k, \ell - 1) + (1 - \alpha_s) \sum_{i=-w}^w b_i |Y(k - i, \ell)|^2 \quad (15)$$

where  $\alpha_s$  ( $0 \leq \alpha_s \leq 1$ ) is a forgetting factor for the smoothing in time, and  $b$  is a normalized window function ( $\sum_{i=-w}^w b_i = 1$ ) that determines the order of smoothing in frequency. A useful estimator  $\mathcal{M}$ , particularly under low SNR and non-stationary noise conditions, can be obtained by the *Minima Controlled Recursive Averaging* approach [22], [23].

For a given signal, we define its local non-stationarity (LNS) by the local ratio between the total and pseudo-stationary spectral power:

$$\Lambda(Y(k, \ell)) = \frac{\mathcal{S}Y(k, \ell)}{\mathcal{M}Y(k, \ell)}. \quad (16)$$

The LNS is a statistic of  $Y$ , fluctuating about one in the absence of transients, and expectedly well above one in the neighborhood of time-frequency bins that contain transients.

Let three hypotheses  $H_{0s}$ ,  $H_{0t}$ , and  $H_1$  indicate respectively absence of transients, presence of an interfering transient, and presence of a source transient at the beamformer output. Let  $\Lambda_0$  denote a threshold value of the LNS for the detection of transients at the beamformer output (*i.e.*, accept  $H_1 \cup H_{0t}$  if  $\Lambda(Y) > \Lambda_0$ , and accept  $H_{0s}$  otherwise). Then, the false alarm probability is defined by

$$P_{f,Y} = \mathcal{P}(\Lambda(Y) > \Lambda_0 \mid H_{0s})$$



$$= \mathcal{P}(\mathcal{S}Y(k, \ell) > \Lambda_0 \cdot \mathcal{M}Y(k, \ell) \mid H_{0s}) . \quad (17)$$

Since  $\mathcal{S}Y(k, \ell)$  is approximately chi-square distributed with  $\mu$  degrees of freedom (Appendix A),

$$F_{\mathcal{S}Y(k, \ell)}(x) \approx F_{\chi^2; \mu} \left( \frac{\mu x}{\phi_{YY}(k, \ell)} \right) ,$$

and since  $\mathcal{M}Y(k, \ell)$  approximates the PSD of  $Y$  when  $H_{0s}$  is true, we have

$$\begin{aligned} P_{f,Y} &\approx 1 - F_{\chi^2; \mu} \left( \frac{\mu \Lambda_0 \cdot \mathcal{M}Y(k, \ell)}{\phi_{YY}(k, \ell)} \right) \Big|_{H_{0s}} \\ &\approx 1 - F_{\chi^2; \mu}(\mu \Lambda_0) . \end{aligned} \quad (18)$$

From this equation, the required threshold value for a specified  $P_{f,Y}$  is

$$\Lambda_0 = \frac{1}{\mu} F_{\chi^2; \mu}^{-1}(1 - P_{f,Y}) . \quad (19)$$

The probability of detection is given by

$$\begin{aligned} P_{d,Y} &= \mathcal{P}(\Lambda(Y) > \Lambda_0 \mid H_1 \cup H_{0t}) \\ &= \mathcal{P}(\mathcal{S}Y(k, \ell) > \Lambda_0 \cdot \mathcal{M}Y(k, \ell) \mid H_1 \cup H_{0t}) \\ &\approx 1 - F_{\chi^2; \mu} \left( \frac{\mu \Lambda_0 \cdot \mathcal{M}Y(k, \ell)}{\phi_{YY}(k, \ell)} \right) \Big|_{H_1 \cup H_{0t}} . \end{aligned} \quad (20)$$

Using Eq. (10) and the approximation  $\mathcal{M}Y \approx \phi_{YY}|_{H_{0s}}$  yields

$$P_{d,Y} \approx 1 - F_{\chi^2; \mu} \left( \frac{\mu \Lambda_0 C_{12} \lambda_s}{C_{11} \lambda_x + C_{12} \lambda_s + C_{13} \lambda_t} \right) . \quad (21)$$

Substituting (19) into (21) we obtain that for a specified false alarm probability, the detection probability is

$$P_{d,Y} = 1 - F_{\chi^2; \mu} \left[ \frac{1}{1 + \xi_Y} F_{\chi^2; \mu}^{-1}(1 - P_{f,Y}) \right] \quad (22)$$

where

$$\xi_Y \triangleq \frac{C_{11} \lambda_x + C_{13} \lambda_t}{C_{12} \lambda_s} \quad (23)$$

represents the ratio between the transient and pseudo-stationary power at the beamformer output. Fig. 2 shows the receiver operating characteristic (ROC) curve for detection of transients at the beamformer output, with the false alarm probability as parameter, and  $\mu$  set to 32.2 (this value of  $\mu$  is obtained for a smoothing  $\mathcal{S}$  of the form (15), with  $\alpha_s = 0.9$ , and  $b = [0.25 \ 0.5 \ 0.25]$ ). Suppose

that we require a false alarm probability no larger than  $P_{f,Y} = 10^{-2}$ , and suppose that transients at the beamformer output are defined by  $\xi_Y \geq 2$ . Then, the detection probability obtained using a detector  $\Lambda(Y) > \Lambda_0 = 1.67$  is  $P_{d,Y} = 0.98$ .

### B. Detection of transients at the reference noise signals

Given that a transient was detected at the beamformer output, its modification rule depends on the presence of a simultaneous transient at one of the reference signals. Let

$$\Lambda(\mathbf{U}(k, \ell)) = \max_{2 \leq i \leq M} \left\{ \frac{SU_i(k, \ell)}{\mathcal{M}U_i(k, \ell)} \right\} \quad (24)$$

denote the LNS of the reference signals, and let  $\Lambda_1$  be a corresponding threshold value for detecting transients. Then the false alarm probability is defined by

$$\begin{aligned} P_{f,\mathbf{U}} &= \mathcal{P}(\Lambda(\mathbf{U}(k, \ell)) > \Lambda_1 \mid H_{0s}) \\ &= \mathcal{P}\left(\max_{2 \leq i \leq M} \left\{ \frac{SU_i(k, \ell)}{\mathcal{M}U_i(k, \ell)} \right\} > \Lambda_1 \mid H_{0s}\right). \end{aligned} \quad (25)$$

Assuming that  $\left\{ \frac{SU_i(k, \ell)}{\mathcal{M}U_i(k, \ell)} \right\}_{i=2}^M$  are statistically independent, we have

$$\begin{aligned} P_{f,\mathbf{U}} &= 1 - \prod_{i=2}^M \mathcal{P}\left(\frac{SU_i(k, \ell)}{\mathcal{M}U_i(k, \ell)} \leq \Lambda_1 \mid H_{0s}\right) \\ &\approx 1 - \prod_{i=2}^M F_{\chi^2; \mu}\left(\frac{\mu \Lambda_1 \cdot \mathcal{M}U_i(k, \ell)}{\phi_{U_i U_i}(k, \ell)}\right) \Big|_{H_{0s}} \\ &\approx 1 - F_{\chi^2; \mu}^{M-1}(\mu \Lambda_1) \end{aligned} \quad (26)$$

where the last approximation was obtained by using  $\mathcal{M}U_i \approx \phi_{U_i U_i} \mid_{H_{0s}}$ . Thus, for a specified false alarm probability,  $P_{f,\mathbf{U}}$ , the threshold value is

$$\Lambda_1 = \frac{1}{\mu} F_{\chi^2; \mu}^{-1} \left[ (1 - P_{f,\mathbf{U}})^{\frac{1}{M-1}} \right]. \quad (27)$$

The detection probability of a transient at one of the reference signals is given by

$$P_{d,\mathbf{U}} = \mathcal{P}(\Lambda(\mathbf{U}(k, \ell)) > \Lambda_1 \mid H_1 \cup H_{0t})$$

$$\begin{aligned}
&= 1 - \prod_{i=2}^M \mathcal{P} \left( \frac{SU_i(k, \ell)}{\mathcal{M}U_i(k, \ell)} \leq \Lambda_1 \mid H_1 \cup H_{0t} \right) \\
&\approx 1 - \prod_{i=2}^M F_{\chi^2; \mu} \left( \frac{\mu \Lambda_1 \cdot C_{i2} \lambda_s}{C_{i1} \lambda_x + C_{i2} \lambda_s + C_{i3} \lambda_t} \right). \tag{28}
\end{aligned}$$

Substituting (27) into (28), and denoting by  $\xi_{U_i} = \frac{C_{i1} \lambda_x + C_{i3} \lambda_t}{C_{i2} \lambda_s}$  the ratio of transient to pseudo-stationary power at the  $i$ -th reference signal, we have

$$\begin{aligned}
P_{d, \mathbf{U}} &\approx 1 - \prod_{i=2}^M F_{\chi^2; \mu} \left( \frac{1}{1 + \xi_{U_i}} F_{\chi^2; \mu}^{-1} \left[ (1 - P_{f, \mathbf{U}})^{\frac{1}{M-1}} \right] \right) \\
&\geq 1 - (1 - P_{f, \mathbf{U}})^{\frac{M-2}{M-1}} \cdot F_{\chi^2; \mu} \left( \frac{1}{1 + \xi_{\mathbf{U}}} F_{\chi^2; \mu}^{-1} \left[ (1 - P_{f, \mathbf{U}})^{\frac{1}{M-1}} \right] \right) \tag{29}
\end{aligned}$$

where  $\xi_{\mathbf{U}} \triangleq \max \{ \xi_{U_i} \mid 2 \leq i \leq M \}$ . Equality in (29) is derived when all  $\xi_{U_i, t}$  but one are identically zero. Fig. 3 shows the receiver operating characteristic (ROC) curve for detection of transients at the reference noise signals, with the false alarm probability as parameter. Four sensors are used, and  $\mu$  is set to 32.2. Suppose that we require a false alarm probability no larger than  $P_{f, \mathbf{U}} = 10^{-2}$ , and suppose that transients at the reference outputs are defined by  $\xi_{\mathbf{U}} \geq 2$ . Then, the detection probability obtained using a detector  $\Lambda(\mathbf{U}) > \Lambda_1 = 1.81$  is  $P_{d, \mathbf{U}} = 0.96$ .

### C. The transient beam-to-reference ratio

The TBRR is a useful statistic to determine the origin of a transient, once detected simultaneously at the beamformer output and at one of the reference noise signals [29]. Since the operator  $\mathcal{S}$  gives a measure of local spectral power, and  $\mathcal{M}$  estimates the background pseudo-stationary power, then their difference yields a measure of the local transient power<sup>1</sup>. We define the TBRR by

$$\Omega(Y, \mathbf{U}) = \frac{\mathcal{S}Y - \mathcal{M}Y}{\max_{2 \leq i \leq M} \{ \mathcal{S}U_i - \mathcal{M}U_i \}}. \tag{30}$$

Transient signal components are relatively strong at the beamformer output, whereas transient noise components are relatively strong at one of the reference signals. Hence, we expect  $\Omega(Y, \mathbf{U})$  to be large for signal transients, and small for noise transients. Let  $\Omega_0$  denote a threshold value of

<sup>1</sup>Recall that transient components are assumed to be uncorrelated with pseudo-stationary noise components

the TBRR for the decision between signal and noise (*i.e.*, accept  $H_1$  only if  $\Omega(Y, \mathbf{U}) > \Omega_0$ ). The conditional false alarm probability is defined by the probability of accepting  $H_1$  in the absence of a source signal, given that a transient was simultaneously detected at the beamformer output and at one of the reference signals:

$$P_{f,\Omega} = \mathcal{P} \{ \Omega(Y, \mathbf{U}) > \Omega_0 \mid (H_{0s} \cup H_{0t}) \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \} . \quad (31)$$

Assuming that during absence of transients ( $H_{0s}$ ), simultaneous transients at the beamformer output and at the reference signals are improbable (*i.e.*, the threshold  $\Lambda_0$  and  $\Lambda_1$  are chosen such that  $\mathcal{P} \{ H_{0s} \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \}$  is negligible), and assuming that transient beam-to-reference ratios for individual reference signals ( $\frac{\mathcal{S}Y - \mathcal{M}Y}{\mathcal{S}U_i - \mathcal{M}U_i}$ ) are statistically independent, we have

$$\begin{aligned} P_{f,\Omega} &\approx \mathcal{P} \{ \Omega(Y, \mathbf{U}) > \Omega_0 \mid H_{0t} \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \} \\ &\approx \prod_{i=2}^M \mathcal{P} \left\{ \frac{\mathcal{S}Y - \mathcal{M}Y}{\mathcal{S}U_i - \mathcal{M}U_i} > \Omega_0 \mid H_{0t} \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \right\} \\ &\approx \prod_{i=2}^M \mathcal{P} \left\{ \frac{C_{12}(\hat{\lambda}_s - \lambda_s) + C_{13}\hat{\lambda}_t}{C_{i2}(\hat{\lambda}_s - \lambda_s) + C_{i3}\hat{\lambda}_t} > \Omega_0 \mid H_{0t} \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \right\} , \end{aligned} \quad (32)$$

where we used Eq. (10) and the fact that  $\mathcal{S}$  is an estimator for the PSD:

$$\begin{bmatrix} \mathcal{S}Y(k, \ell) \\ \mathcal{S}U_2(k, \ell) \\ \vdots \\ \mathcal{S}U_M(k, \ell) \end{bmatrix} = \begin{bmatrix} C_{11}(k, \ell) & C_{12}(k, \ell) & C_{13}(k, \ell) \\ \vdots & \vdots & \vdots \\ C_{M1}(k, \ell) & C_{M2}(k, \ell) & C_{M3}(k, \ell) \end{bmatrix} \begin{bmatrix} \hat{\lambda}_x(k, \ell) \\ \hat{\lambda}_s(k, \ell) \\ \hat{\lambda}_t(k, \ell) \end{bmatrix} . \quad (33)$$

Given that  $H_{0t}$  is true, detection of a transient at the beamformer output implies  $C_{12}\hat{\lambda}_s + C_{13}\hat{\lambda}_t > \Lambda_0 C_{12}\lambda_s$ . Detection of a transient at one of the reference signals implies that there exists  $i \in [2, M]$  such that  $C_{i2}\hat{\lambda}_s + C_{i3}\hat{\lambda}_t > \Lambda_1 C_{i2}\lambda_s$ . Furthermore, since we assume that the pseudo-stationary noise at the beamformer output is weak compared to that associated with any reference noise signal ( $C_{12}/C_{13} \leq C_{i2}/C_{i3}$  for all  $i \in [2, M]$ ), then with probability one there exists  $i \in [2, M]$  such that

$$\frac{C_{12}(\hat{\lambda}_s - \lambda_s) + C_{13}\hat{\lambda}_t}{C_{i2}(\hat{\lambda}_s - \lambda_s) + C_{i3}\hat{\lambda}_t} \leq \frac{C_{13}}{C_{i3}} . \quad (34)$$

Accordingly, by choosing

$$\Omega_0 \geq \frac{C_{13}}{\min_{2 \leq i \leq M} \{C_{i3}\}} \quad (35)$$

we have that  $P_{f,\Omega} = 0$  with probability one.

The conditional detection probability is defined by the probability of accepting  $H_1$  in the presence of a desired signal, given that a transient was simultaneously detected at the beamformer output and at one of the reference signals:

$$P_{d,\Omega} = \mathcal{P} \{ \Omega(Y, \mathbf{U}) > \Omega_0 \mid H_1 \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \} . \quad (36)$$

Assuming that desired and interfering transients do not overlap in the time-frequency domain, we have

$$\begin{aligned} P_{d,\Omega} &\approx \prod_{i=2}^M \mathcal{P} \left\{ \frac{\mathcal{S}Y - \mathcal{M}Y}{\mathcal{S}U_i - \mathcal{M}U_i} > \Omega_0 \mid H_1 \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \right\} \\ &\approx \prod_{i=2}^M \mathcal{P} \left\{ \frac{C_{11}\hat{\lambda}_x + C_{12}(\hat{\lambda}_s - \lambda_s)}{C_{i1}\hat{\lambda}_x + C_{i2}(\hat{\lambda}_s - \lambda_s)} > \Omega_0 \mid H_1 \cap (\Lambda(Y) > \Lambda_0) \cap (\Lambda(\mathbf{U}) > \Lambda_1) \right\} . \end{aligned} \quad (37)$$

Since there is no correlation between the desired signal, transient noise, and pseudo-stationary noise components, the distributions of  $\hat{\lambda}_x$  and  $\hat{\lambda}_s$  are the same as the distribution of  $\mathcal{S}Y$  (chi-square with  $\mu$  degrees of freedom). Accordingly,  $E\{\hat{\lambda}_x\} = \lambda_x$ ,  $Var\{\hat{\lambda}_x\} = \frac{2}{\mu}\lambda_x^2$ ,  $E\{\hat{\lambda}_s\} = \lambda_s$ , and  $Var\{\hat{\lambda}_s\} = \frac{2}{\mu}\lambda_s^2$ . For  $\mu \gg 1$ , the transient power at the beamformer output is relative to  $C_{11}\hat{\lambda}_x$ , and at a reference signal is relative to  $C_{i1}\hat{\lambda}_x$ . Therefore, to retain a high detection probability,  $P_{d,\Omega}$ , we require

$$\Omega_0 \leq \min_{2 \leq i \leq M} \left\{ \frac{C_{11}}{C_{i1}} \right\} = \frac{C_{11}}{\max_{2 \leq i \leq M} \{C_{i1}\}} . \quad (38)$$

Let

$$Q \triangleq \frac{C_{11}}{C_{13}} \cdot \frac{\min_{2 \leq i \leq M} \{C_{i3}\}}{\max_{2 \leq i \leq M} \{C_{i1}\}} \quad (39)$$

define a *transient discrimination quality* (TDQ) of a beamformer. Then from Eqs. (35) and (38) it follows that discrimination between transient noise and desired signal components is possible when  $Q \geq 1$  (in practice, we obtained good performance,  $P_{f,\Omega} \rightarrow 0$ ,  $P_{d,\Omega} \rightarrow 1$ , for  $Q \geq 3$ ).

Fig. 4 summarizes a block diagram for the detection of desired source components at the beamformer output. The detection is carried out in the time-frequency plane for each frame and frequency bin. Case 1 is reached when no transients have been detected at the beamformer output, or when the TBRR is lower than the threshold  $\Omega_0$ . In this case, presumably no desirable transients are present at the beamformer output, and consequently strong noise suppression is applicable. Considering Case 2, a transient has been detected at the beamformer output, but not at any reference signal. This case indicates that the transient is likely a desirable source component, and a cautious noise suppression would therefore be involved. Finally, Case 3 is determined when transients are simultaneously detected at the beamformer output and at a reference signal, and conjunctionally the value of the TBRR is above  $\Omega_0$ . In this case, the larger the TBRR is, the higher the likelihood that a transient comes from a desired source.

#### IV. MULTI-CHANNEL POST-FILTERING

In this section, we address the problem of estimating the time-varying spectrum of the beamformer output noise, and present the multi-channel post-filtering approach. Fig. 5 describes the block diagram of the proposed multi-channel post-filtering. Desired source components are detected at the beamformer output, and an estimate  $\hat{q}(k, \ell)$  for the *a priori* signal absence probability is produced. Based on a Gaussian statistical model [20], and a decision-directed estimator for the *a priori* SNR under signal presence uncertainty [22], we derive an estimator  $p(k, \ell) \triangleq \mathcal{P}(H_1 | Y, \mathbf{U})$  for the signal presence probability. This estimator controls the components that are introduced as noise into the PSD estimator. Finally, spectral enhancement of the beamformer output is achieved by applying an *optimally-modified log-spectral amplitude* (OM-LSA) gain function [22]. This gain minimizes the mean-square error of the log-spectra under signal presence uncertainty.

Referring to Fig. 4, Cases 1 and 2 imply presumable signal absence and presence, respectively. Therefore, we set  $\hat{q}(k, \ell)$  to 1 in Case 1, and to 0 in Case 2. However, when transients are simultaneously detected in both the beamformer output and one of the reference signals, and the TBRR is larger than  $\Omega_0$  (Case 3), then the value of the *a priori* signal absence probability is determined

according to

$$\hat{q}(k, \ell) = \begin{cases} 1, & \text{if } \gamma_s(k, \ell) \leq 1 \\ 0, & \text{if } \gamma_s(k, \ell) > \gamma_0 \text{ and } \Omega(Y, \mathbf{U}) \geq 3\Omega_0 \\ \max \left\{ \frac{\gamma_0 - \gamma_s(k, \ell)}{\gamma_0 - 1}, \frac{3\Omega_0 - \Omega(Y, \mathbf{U})}{2\Omega_0} \right\}, & \text{otherwise,} \end{cases} \quad (40)$$

where  $\gamma_s(k, \ell) \triangleq |Y(k, \ell)|^2 / \mathcal{M}Y(k, \ell)$  represents the *a posteriori* SNR at the beamformer output with respect to the pseudo-stationary noise, and  $\gamma_0$  denotes a constant satisfying

$$\mathcal{P}(\gamma_s(k, \ell) \geq \gamma_0 \mid H_{0s}) < \epsilon \quad (41)$$

for a certain significance level  $\epsilon$ . Eq. (40) suggests that the likelihood of signal presence increases with the values of  $\gamma_s$  and  $\Omega(Y, \mathbf{U})$ . Indeed, from (41) we have that when the *a posteriori* SNR is larger than  $\gamma_0$ , either  $H_1$  or  $H_{0t}$  is true ( $H_{0s}$  is very unlikely). On the other hand,  $\Omega(Y, \mathbf{U})$  discriminates between desired source components ( $H_1$ ) and noise transients ( $H_{0t}$ ). Therefore, Eq. (40) is obtained by combining conditions on  $\gamma_s$  and  $\Omega(Y, \mathbf{U})$ , and assuming smooth bilinear transition from signal absence to presence in the regions  $\gamma_s \in [E\{\gamma_s \mid H_{0s}\}, \gamma_0]$  and  $\Omega(Y, \mathbf{U}) \in [\Omega_0, 3\Omega_0]$ .

Under the assumed statistical model, the distribution of  $\gamma_s(k, \ell)$ , in the absence of transients, is exponential [23]:

$$f(\gamma_s(k, \ell) \mid H_{0s}) = e^{-\gamma_s(k, \ell)} u(\gamma_s(k, \ell)) \quad (42)$$

where  $u(\cdot)$  is the unit step function (*i.e.*,  $u(\gamma) = 1$  for  $\gamma \geq 0$  and  $u(\gamma) = 0$  otherwise). Accordingly,  $\gamma_0 = -\log(\epsilon)$  (typically, we use  $\epsilon = 0.01$ , so  $\gamma_0 = 4.6$ ). Furthermore, the signal presence probability is given by

$$p(k, \ell) = \left\{ 1 + \frac{q(k, \ell)}{1 - q(k, \ell)} (1 + \xi(k, \ell)) \exp(-v(k, \ell)) \right\}^{-1} \quad (43)$$

where  $\xi(k, \ell) \triangleq E\{|X(k, \ell)|^2\} / \lambda_d(k, \ell)$  is the *a priori* SNR,  $\lambda_d(k, \ell)$  is the noise PSD at the beamformer output,  $v(k, \ell) \triangleq \gamma(k, \ell) \xi(k, \ell) / (1 + \xi(k, \ell))$ , and  $\gamma(k, \ell) \triangleq |Y(k, \ell)|^2 / \lambda_d(k, \ell)$  is the *a posteriori* SNR. The *a priori* SNR is estimated by [22]

$$\hat{\xi}(k, \ell) = \alpha G_{H_1}^2(k, \ell - 1) \gamma(k, \ell - 1) + (1 - \alpha) \max\{\gamma(k, \ell) - 1, 0\} \quad (44)$$

where  $\alpha$  is a weighting factor that controls the trade-off between noise reduction and signal distor-

tion, and

$$G_{H_1}(k, \ell) \triangleq \frac{\xi(k, \ell)}{1 + \xi(k, \ell)} \exp\left(\frac{1}{2} \int_{v(k, \ell)}^{\infty} \frac{e^{-t}}{t} dt\right) \quad (45)$$

is the spectral gain function of the *Log-Spectral Amplitude* (LSA) estimator when signal is surely present<sup>2</sup> [21]. The MCRA approach for noise spectrum estimation [23] is to recursively average past spectral power values of the noisy measurement, using a smoothing parameter that is controlled by the minima values of a smoothed periodogram. The recursive averaging is given by

$$\hat{\lambda}_d(k, \ell + 1) = \tilde{\alpha}_d(k, \ell) \hat{\lambda}_d(k, \ell) + \beta \cdot [1 - \tilde{\alpha}_d(k, \ell)] |Y(k, \ell)|^2 \quad (46)$$

where  $\tilde{\alpha}_d(k, \ell)$  is a time-varying frequency-dependent smoothing parameter, and  $\beta$  is a factor that compensates the bias when signal is absent. The smoothing parameter is determined by the signal presence probability,  $p(k, \ell)$ , and a constant  $\alpha_d$  ( $0 < \alpha_d < 1$ ) that represents its minimal value:

$$\tilde{\alpha}_d(k, \ell) \triangleq \alpha_d + (1 - \alpha_d) p(k, \ell). \quad (47)$$

When signal is present,  $\tilde{\alpha}_d$  is close to one, thus preventing the noise estimate from increasing as a result of signal components. As the probability of signal presence decreases, the smoothing parameter gets smaller, facilitating a faster update of the noise estimate.

The estimate of the clean signal STFT is finally given by

$$\hat{X}(k, \ell) = G(k, \ell) Y(k, \ell), \quad (48)$$

where

$$G(k, \ell) = \{G_{H_1}(k, \ell)\}^{p(k, \ell)} \cdot G_{min}^{1-p(k, \ell)} \quad (49)$$

is the OM-LSA gain function and  $G_{min}$  denotes a lower bound constraint for the gain when signal is absent. The implementation of the multi-channel post-filtering algorithm is summarized in Fig. 6. Typical values of the respective parameters, for a sampling rate of 8 kHz, are given in Table I.

## V. EXPERIMENTAL RESULTS

To validate the usefulness of the proposed post-filtering approach under non-stationary noise conditions, we compare its performance to a single-channel post-filtering in various car environments. Specifically, multi-channel speech signals are degraded by interfering speakers and various

<sup>2</sup>The advantage of  $\hat{\xi}(k, \ell)$  over the “decision-directed” estimator of Ephraim and Malah [20], particularly for weak signal components and low input SNR, is discussed in [22].



car noise types. Then, beamforming is applied to the noisy signals, followed by either single-channel or multi-channel post-filtering. The performance evaluation includes objective quality measures, as well as a subjective study of speech spectrograms and informal listening tests.

A linear array, consisting of four microphones with 5 cm spacing, is mounted in a car on the visor. Clean speech signals are recorded at a sampling rate of 8 kHz in the absence of background noise (standing car, silent environment). An interfering speaker and car noise signals are recorded while the car speed is about 60 km/h, and windows are either closed, or the window next to the driver is slightly open (about 5 cm). The input microphone signals are generated by mixing the speech and noise signals at various SNR levels in the range  $[-5, 10]$  dB.

An adaptive beamformer (specifically, the TF-GSC, proposed by Gannot *et al.* [24]) is applied to the noisy multi-channel signals. The beamformer output is enhanced using the OM-LSA estimator [22], and is referred to as the single-channel post-filtering output. Alternatively, the beamformer output, enhanced using the procedure described in the previous section, is referred to as the multi-channel post-filtering output. Three different objective quality measures are used in our evaluation. The first is segmental SNR defined by [30]

$$\text{SegSNR} = \frac{1}{L} \sum_{\ell=0}^{L-1} 10 \cdot \log \frac{\sum_{n=0}^{N-1} x^2(n + \ell N/2)}{\sum_{n=0}^{N-1} [x(n + \ell N/2) - \hat{x}(n + \ell N/2)]^2} \quad [\text{dB}] \quad (50)$$

where  $L$  represents the number of frames in the signal, and  $N = 256$  is the number of samples per frame (corresponding to 32 ms frames, and 50% overlap). The segmental SNR at each frame is limited to perceptually meaningful range between 35 dB and  $-10$  dB [31], [32]. This measure takes into account both residual noise and speech distortion. The second quality measure is noise reduction (NR), which is defined by

$$\text{NR} = \frac{1}{|\mathcal{L}'|} \sum_{\ell \in \mathcal{L}'} 10 \cdot \log \frac{\sum_{n=0}^{N-1} z_1^2(n + \ell N/2)}{\sum_{n=0}^{N-1} \hat{x}^2(n + \ell N/2)} \quad [\text{dB}] \quad (51)$$

where  $\mathcal{L}'$  represents the set of frames that contain only noise, and  $|\mathcal{L}'|$  its cardinality. The NR measure compares the noise level in the enhanced signal to the noise level recorded by the first

microphone. The third quality measure is log-spectral distance (LSD), which is defined by

$$\text{LSD} = \frac{1}{L} \sum_{\ell=0}^{L-1} \left\{ \frac{1}{N/2+1} \sum_{k=0}^{N/2} \left[ 10 \cdot \log \mathcal{A}X(k, \ell) - 10 \cdot \log \mathcal{A}\hat{X}(k, \ell) \right]^2 \right\}^{\frac{1}{2}} \quad [\text{dB}] \quad (52)$$

where  $\mathcal{A}X(k, \ell) \triangleq \max \{ |X(k, \ell)|^2, \delta \}$  is the spectral power, clipped such that the log-spectrum dynamic range is confined to about 50 dB (that is,  $\delta = 10^{-50/10} \cdot \max_{k, \ell} \{ |X(k, \ell)|^2 \}$ ).

Fig. 7 shows experimental results of the average segmental SNR, obtained for various noise types and at various noise levels. The segmental SNR is evaluated at the first microphone, the beamformer output, and the post-filtering outputs. A theoretical limit post-filtering, achievable by calculating the noise spectrum from the noise itself, is also considered. Results of the NR and LSD measures are presented in Figs. 8 and 9, respectively. It can be readily seen that beamforming alone does not provide sufficient noise reduction in a car environment, owing to its limited ability to reduce diffuse noise [24]. Furthermore, multi-channel post-filtering is consistently better than single-channel post-filtering under all noise conditions. The improvement in performance of the former over the latter is expectedly high in non-stationary noise environments (specifically, open windows or interfering speaker), but is insignificant otherwise, since multi-channel post-filtering reduces to single-channel in pseudo-stationary noise environments.

A subjective comparison between multi-channel and single-channel post-filtering was conducted using speech spectrograms and validated by informal listening tests. Typical examples of speech spectrograms are presented in Fig. 10 for the case of non-stationary noise (interfering speaker, open window) at  $\text{SNR} = -0.9$  dB. The beamformer output (Fig. 10(c)) is clearly characterized by a high level of noise. Its enhancement using single-channel post-filtering well suppresses the pseudo-stationary noise, but adversely retains the transient noise components. By contrast, the enhancement using multi-channel post-filtering results in superior noise attenuation, while preserving the desired source components.

Fig. 11 shows traces of the improvement in segmental SNR and LSD measures, gained by the multi-channel post-filtering and theoretical limit, in comparison with a single-channel post-filtering. The traces are averaged out over a period of about 400 ms (25 frames of 32 ms each, with 50% over-

lap). The noise PSD at the beamformer output varies substantially due to the residual interfering components of speech, wind blows, and passing cars. The improvement in performance over the single-channel post-filtering is obtained when the noise spectrum fluctuates. In some instances the increase in segmental SNR surpasses as much as 8 dB, and the decrease in LSD is greater than 6 dB. Clearly, a single-channel post-filter is inefficient at attenuating highly non-stationary noise components, since it lacks the ability to differentiate such components from the speech components. On the other hand, the proposed multi-channel post-filtering approach achieves a significantly reduced level of background noise, whether stationary or not, without further distorting speech components. This is verified by subjective informal listening tests.

## VI. CONCLUSION

We have described a multi-channel post-filtering approach for arbitrary beamformers, that is particularly advantageous in non-stationary noise environments. The beamformer is realistically assumed to have a steering error, a blocking matrix that is unable to block all of the desired signal components, and a noise canceller that is adapted to the pseudo-stationary noise, but not modified during transient interferences. Accordingly, the reference noise signals may include some desired signal components. Furthermore, transient noise components that leak through the sidelobes of the fixed beamformer may proceed to the beamformer primary output. A mild assumption is made with regard to the beamformer, that a desired signal component is stronger at the beamformer output than at any reference noise signal, and a noise component is strongest at one of the reference signals. Consequently, transients are detected at the beamformer output and either suppressed or preserved based on the transient beam-to-reference ratio.

We derived an estimator for the signal presence probability, that controls the rate of recursive averaging for obtaining a noise spectrum estimate. It also modifies the spectral gain function for obtaining an estimate of the clean signal spectral amplitude. The proposed method was tested in various non-stationary car noise environments, and its performance was compared to a single-channel post-filtering approach. We showed that multi-channel post-filtering is better than single-channel post-filtering particularly under highly non-stationary noise conditions (such as noise resulting from wind blows, passing cars, interfering speakers, etc.). While transient noise components are indistin-

guishable from desired source components if using a state-of-the-art single-channel post-filtering, the enhancement of the beamformer output by multi-channel post-filtering produces a significantly reduced level of residual transient noise without further distorting the desired signal components.

## APPENDIX

### I. STATISTICS OF $\mathcal{S}Y(k, \ell)$

Successive spectral power values of the beamformer output  $|Y(k, \ell)|^2$  are generally correlated, and there is no closed form solution for the probability density function of  $\mathcal{S}Y(k, \ell)$ . However, Eq. (15) can be written as

$$\mathcal{S}Y(k, \ell) = (1 - \alpha_s) \sum_{i=-w}^w \sum_{j=0}^{\infty} b_i \alpha_s^j |Y(k - i, \ell - j)|^2. \quad (53)$$

Approximating  $\mathcal{S}Y(k, \ell)$  as the sum of  $\mu$  squared mutually independent normal variables [33], [23], its distribution function is given by

$$F_{\mathcal{S}Y(k, \ell)}(x) \approx F_{\chi^2; \mu} \left( \frac{\mu x}{\phi_{YY}(k, \ell)} \right) \quad (54)$$

where  $F_{\chi^2; \mu}(x)$  denotes the standard chi-square distribution function, with  $\mu$  degrees of freedom. Specifically,  $F_{\mathcal{S}Y(k, \ell)}(x) = \Gamma\left(\frac{\mu}{2}, \frac{\mu x}{2\phi_{YY}(k, \ell)}\right) u(x) / \Gamma\left(\frac{\mu}{2}\right)$ , where  $\Gamma(\cdot)$  is the gamma function, and  $\Gamma(a, x) \triangleq \int_0^\infty e^{-t} t^{a-1} dt$  is the incomplete gamma function. The equivalent degrees of freedom,  $\mu$ , is determined by the smoothing parameter  $\alpha_s$ , the window function  $b$ , and the spectral analysis parameters of the STFT (size and shape of the analysis window, and frame-update step). The value of  $\mu$  can be estimated by generating a stationary white Gaussian noise  $d(t)$ , transforming it to the time-frequency domain, and substituting the sample mean and variance (over the entire time-frequency plane) into the expression  $\hat{\mu} \approx 2 E^2 \{\mathcal{S}D(k, \ell)\} / \text{var} \{\mathcal{S}D(k, \ell)\}$ .

## ACKNOWLEDGEMENT

The author thanks Dr. Baruch Berdugo for helpful discussions, and Dr. Sharon Gannot for making his adaptive beamforming code (TF-GSC) available.

## REFERENCES

- [1] M. S. Brandstein and D. B. Ward, Eds., *Microphone Arrays: Signal Processing Techniques and Applications*, Springer-Verlag, Berlin, 2001.

- [2] C. Marro, Y. Mahieux, and K. U. Simmer, "Analysis of noise reduction and dereverberation techniques based on microphone arrays with postfiltering," *IEEE Trans. Speech and Audio Processing*, vol. 6, no. 3, pp. 240–259, May 1998.
- [3] K. U. Simmer, J. Bitzer, and C. Marro, *Post-Filtering Techniques*, chapter 3, pp. 39–60, In Brandstein and Ward [1], 2001.
- [4] R. Zelinski, "A microphone array with adaptive post-filtering for noise reduction in reverberant rooms," in *Proc. 13th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-88*, New York, USA, 11–14 April 1988, pp. 2578–2581.
- [5] R. Zelinski, "Noise reduction based on microphone array with LMS adaptive post-filtering," November 1990, vol. 26, pp. 2036–2581.
- [6] K. U. Simmer and A. Wasiljeff, "Adaptive microphone arrays for noise suppression in the frequency domain," in *Proc. 2nd Cost-229 Workshop on Adaptive Algorithms in Communications*, Bordeaux, France, 30 September–2 October 1992, pp. 185–194.
- [7] S. Fischer and K. U. Simmer, "An adaptive microphone array for hands-free communication," in *Proc. 4th International Workshop on Acoustic Echo and Noise Control, IWAENC-95*, Røros, Norway, 21–23 June 1995, pp. 44–47.
- [8] S. Fischer and K. U. Simmer, "Beamforming microphone arrays for speech acquisition in noisy environments," *Speech Communication*, vol. 20, no. 3–4, pp. 215–227, December 1996.
- [9] K. U. Simmer, S. Fischer, and A. Wasiljeff, "Suppression of coherent and incoherent noise using a microphone array," *Annales des Télécommunications*, vol. 49, no. 7–8, pp. 439–446, July 1994.
- [10] J. Bitzer, K. U. Simmer, and K.-D. Kammeyer, "Multichannel noise reduction - algorithms and theoretical limits," in *Proc. European Signal Processing Conference, EUSIPCO-98*, Rhodes, Greece, 8–11 September 1998, pp. 105–108.
- [11] J. Bitzer, K. U. Simmer, and K.-D. Kammeyer, "Theoretical noise reduction limits of the generalized sidelobe canceller (GSC) for speech enhancement," in *Proc. 24th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-99*, Phoenix, Arizona, 15–19 March 1999, pp. 2965–2968.
- [12] J. Bitzer, K. U. Simmer, and K.-D. Kammeyer, "Multi-microphone noise reduction by post-filter and superdirective beamformer," in *Proc. 6th International Workshop on Acoustic Echo and Noise Control, IWAENC-99*, Pocono Manor, Pennsylvania, 27–30 September 1999, pp. 100–103.
- [13] J. Bitzer, K. U. Simmer, and K.-D. Kammeyer, "Multi-microphone noise reduction techniques as front-end devices for speech recognition," *Speech Communication*, vol. 34, no. 1, pp. 3–12, April 2001.
- [14] R. Le Bouquin-Jeannès, A. A. Azirani, and G. Faucon, "Enhancement of speech degraded by coherent and incoherent noise using a cross-spectral estimator," *IEEE Trans. Speech and Audio Processing*, vol. 5, no. 5, pp. 484–487, September 1997.
- [15] J. Meyer and K. U. Simmer, "Multi-channel speech enhancement in a car environment using Wiener filtering and spectral subtraction," in *Proc. 22th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-97*, Munich, Germany, 20–24 April 1997, pp. 21–24.
- [16] D. Mahmoudi, "A microphone array for speech enhancement using multiresolution wavelet transform," in *Proc. 5th European Conf. Speech, Communication and Technology, EUROSPEECH'97*, Rhodes, Greece, 22–25 September 1997, pp. 339–342.
- [17] D. Mahmoudi and A. Drygajlo, "Combined Wiener and coherence filtering in wavelet domain for microphone array speech enhancement," in *Proc. 23th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-98*, Seattle, Washington, 12–15 May 1998, pp. 385–388.
- [18] S. Fischer and K.-D. Kammeyer, "Broadband beamforming with adaptive postfiltering for speech acquisition in noisy environments," in *Proc. 22th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-97*, Munich, Germany, 20–24 April 1997, pp. 359–362.
- [19] I. A. McCowan, C. Marro, and L. Mauuary, "Robust speech recognition using near-field superdirective beamforming with post-filtering," in *Proc. 25th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-2000*, Istanbul, Turkey, 5–9 June 2000, pp. 1723–1726.
- [20] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator," *IEEE Trans. Acoustics, Speech and Signal Processing*, vol. ASSP-32, no. 6, pp. 1109–1121, December 1984.
- [21] Y. Ephraim and D. Malah, "Speech enhancement using a minimum mean-square error log-spectral amplitude estimator," *IEEE Trans. Acoustics, Speech and Signal Processing*, vol. ASSP-33, no. 2, pp. 443–445, April 1985.

- [22] I. Cohen and B. Berdugo, "Speech enhancement for non-stationary noise environments," *Signal Processing*, vol. 81, no. 11, pp. 2403–2418, October 2001.
- [23] I. Cohen, "Noise spectrum estimation in adverse environments: Improved minima controlled recursive averaging," Technical Report, EE PUB 1291, Technion - Israel Institute of Technology, Haifa, Israel, October 2001.
- [24] S. Gannot, D. Burshtein, and E. Weinstein, "Signal enhancement using beamforming and nonstationarity with applications to speech," *IEEE Trans. Signal Processing*, vol. 49, no. 8, pp. 1614–1626, August 2001.
- [25] L. J. Griffiths and C. W. Jim, "An alternative approach to linearly constrained adaptive beamforming," *IEEE Trans. Antennas and Propagation*, vol. AP-30, no. 1, pp. 27–34, January 1982.
- [26] C. W. Jim, "A comparison of two LMS constrained optimal array structures," *Proceedings of the IEEE*, vol. 65, no. 12, pp. 1730–1731, December 1977.
- [27] B. Widrow and S. D. Stearns, *Adaptive Signal Processing*, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1985.
- [28] S. Nordholm, I. Claesson, and P. Eriksson, "The broadband Wiener solution for Griffiths-Jim beamformers," *IEEE Trans. Signal Processing*, vol. 40, no. 9, pp. 474–478, February 1992.
- [29] I. Cohen and B. Berdugo, "Microphone array post-filtering for non-stationary noise suppression," in *Proc. 27th IEEE Internat. Conf. Acoust. Speech Signal Process., ICASSP-2002*, Orlando, Florida, 13-17 May 2002.
- [30] S. R. Quackenbush, T. P. Barnwell, and M. A. Clements, *Objective Measures of Speech Quality*, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1988.
- [31] J. R. Deller, J. H. L. Hansen, and J. G. Proakis, *Discrete-Time Processing of Speech Signals*, IEEE Press, New York, 2nd edition, 2000.
- [32] P. E. Papamichalis, *Practical Approaches to Speech Coding*, Prentice-Hall, Inc., Englewood Cliffs, New Jersey, 1987.
- [33] R. Martin, "Noise power spectral density estimation based on optimal smoothing and minimum statistics," *IEEE Trans. Speech and Audio Processing*, vol. 9, no. 5, pp. 504–512, July 2001.

#### TABLE CAPTIONS

Table I: Values of Parameters Used in the Implementation of the Proposed Multi-Channel Post-Filtering, For a Sampling Rate of 8 kHz.

#### FIGURE CAPTIONS

- Fig. 1: Block diagram of the Griffiths-Jim adaptive beamformer.
- Fig. 2: Receiver operating characteristic curve for detection of transients at the beamformer output ( $\mu = 32.2$ ).
- Fig. 3: Receiver operating characteristic curve for detection of transients at the reference noise signals, using  $M = 4$  sensors ( $\mu = 32.2$ ).
- Fig. 4: Block diagram for detection of desired source components at the beamformer output.
- Fig. 5: Block diagram of the multi-channel post-filtering.
- Fig. 6: The multi-channel post-filtering algorithm.
- Fig. 7: Average segmental SNR at ( $\triangle$ ) microphone #1, ( $\circ$ ) beamformer output, ( $\times$ ) single-channel post-filtering output, (solid line) multi-channel post-filtering output, and (\*) theoretical limit post-filtering output, for various car noise conditions: (a) Closed windows; (b) Open

window; (c) Interfering speaker.

Fig. 8: Average noise reduction at (o) beamformer output, ( $\times$ ) single-channel post-filtering output, (solid line) multi-channel post-filtering output, and (\*) theoretical limit post-filtering output, for various car noise conditions: (a) Closed windows; (b) Open window; (c) Interfering speaker.

Fig. 9: Average log-spectral distance at ( $\triangle$ ) microphone #1, (o) beamformer output, ( $\times$ ) single-channel post-filtering output, (solid line) multi-channel post-filtering output, and (\*) theoretical limit post-filtering output, for various car noise conditions: (a) Closed windows; (b) Open window; (c) Interfering speaker.

Fig. 10: Speech spectrograms. (a) Original clean speech signal at microphone #1: “Five six seven eight nine.”; (b) Noisy signal at microphone #1 (car noise, open window, interfering speaker. SNR =  $-0.9$  dB, SegSNR =  $-6.2$  dB, LSD =  $15.4$  dB); (c) Beamformer output (SegSNR =  $-5.3$  dB, NR =  $5.2$  dB, LSD =  $12.2$  dB); (d) Single-channel post-filtering output (SegSNR =  $-3.8$  dB, NR =  $12.1$  dB, LSD =  $7.4$  dB); (e) Multi-channel post-filtering output (SegSNR =  $-1.3$  dB, NR =  $23.2$  dB, LSD =  $4.6$  dB); (f) Theoretical limit (SegSNR =  $-0.4$  dB, NR =  $24.0$  dB, LSD =  $4.0$  dB).

Fig. 11: Trace of the improvement over a single-channel post-filtering gained by the proposed multi-channel post-filtering (solid) and theoretical limit (dashed): (a) Increase in segmental SNR; (b) Decrease in Log-Spectral Distance.

TABLE I  
VALUES OF PARAMETERS USED IN THE IMPLEMENTATION OF THE PROPOSED MULTI-CHANNEL  
POST-FILTERING, FOR A SAMPLING RATE OF 8 kHz

$\Lambda_0 = 1.67$	$\Lambda_1 = 1.81$	$\Omega_0 = 1$	$\gamma_0 = 4.6$
$\alpha = 0.92$	$\alpha_s = 0.9$	$\alpha_d = 0.85$	$\beta = 1.47$
$b = [0.25 \quad 0.5 \quad 0.25]$		$\mu = 32.2$	$G_{min} = -20 \text{ dB}$

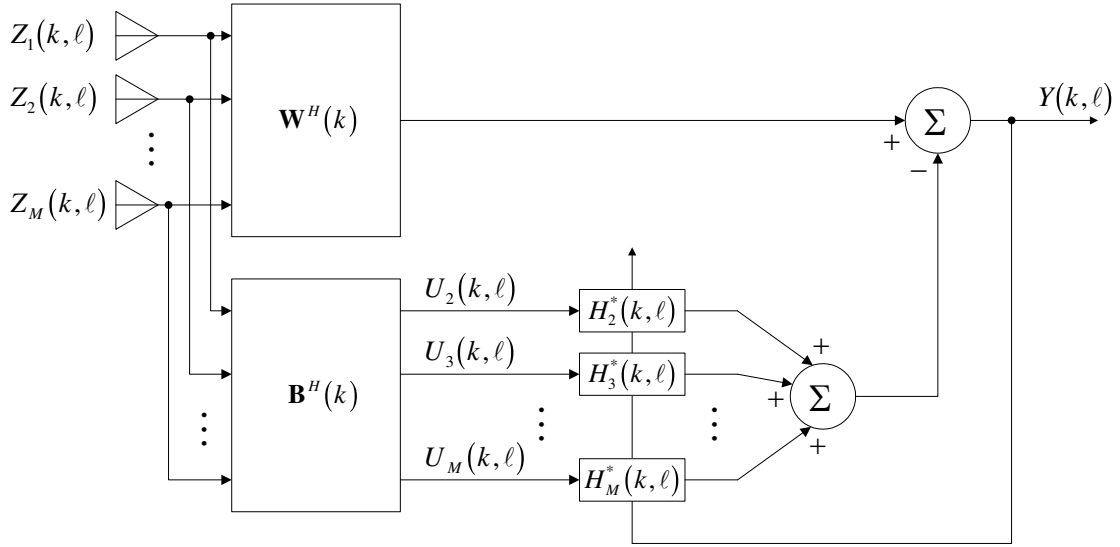


Fig. 1. Block diagram of the Griffiths-Jim adaptive beamformer.



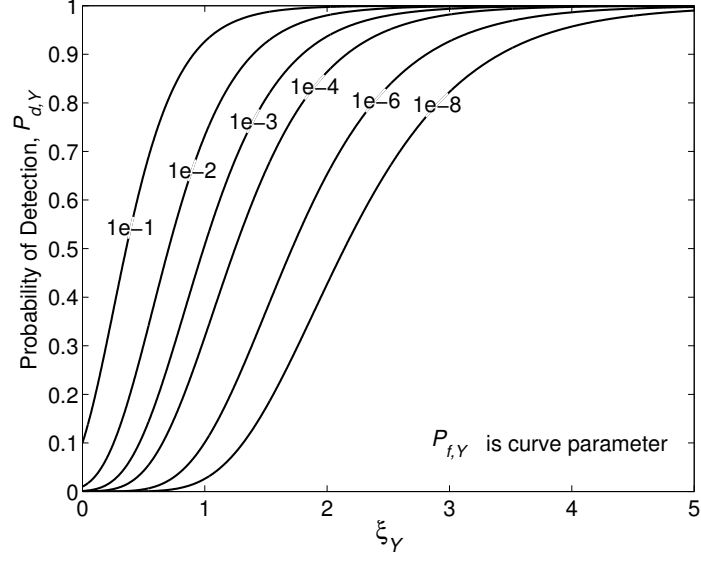


Fig. 2. Receiver operating characteristic curve for detection of transients at the beamformer output ( $\mu = 32.2$ ).

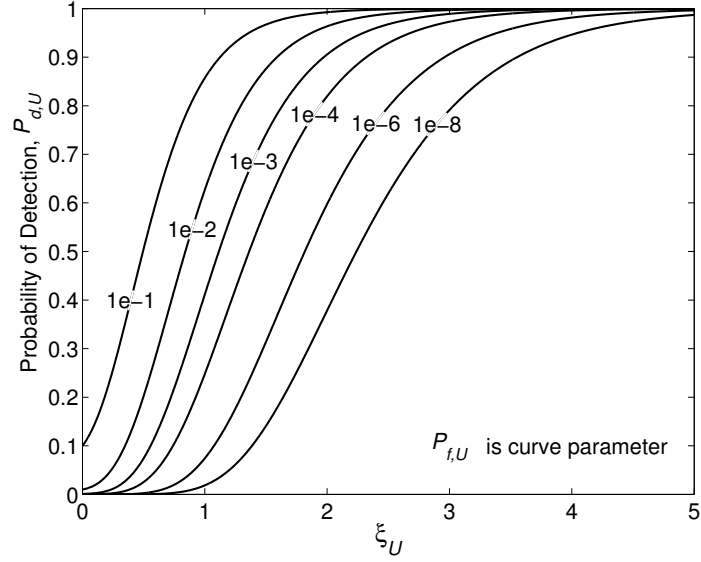


Fig. 3. Receiver operating characteristic curve for detection of transients at the reference noise signals, using  $M = 4$  sensors ( $\mu = 32.2$ ).

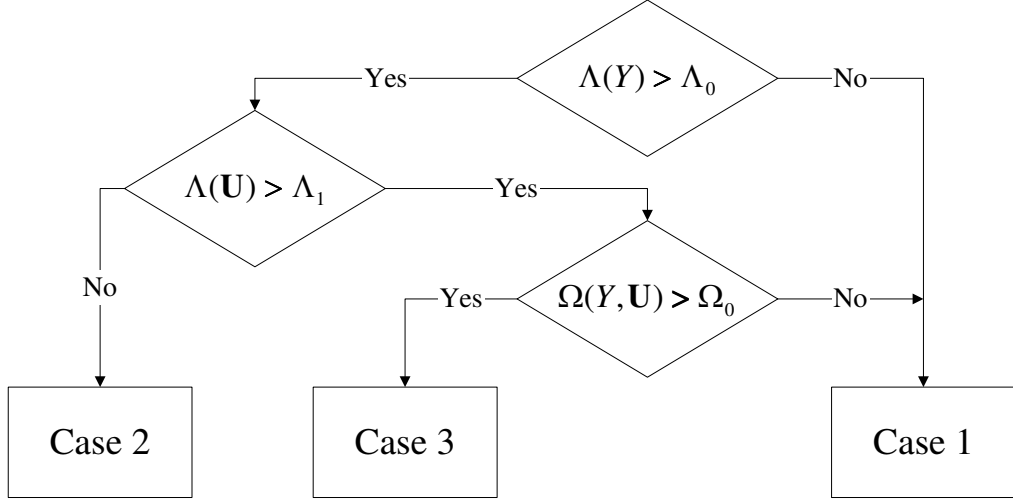


Fig. 4. Block diagram for detection of desired source components at the beamformer output.

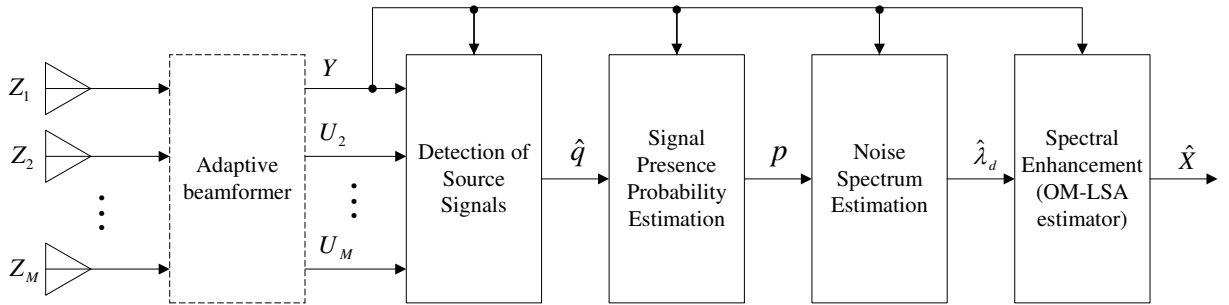


Fig. 5. Block diagram of the multi-channel post-filtering.

Initialize variables at the first frame for all frequency bins  $k$ :

$$\mathcal{S}Y(k, 0) = \mathcal{M}Y(k, 0) = \hat{\lambda}_d(k, 0) = |Y(k, 0)|^2; \quad G_{H_1}(k, 0) = \gamma(k, 0) = 1.$$

For all time frames  $\ell$

For all frequency bins  $k$

Compute the recursively averaged spectrum of the beamformer output  $\mathcal{S}Y(k, \ell)$  using Eq. (15), and update the MCRA estimate of the background pseudo-stationary noise  $\mathcal{M}Y(k, \ell)$  using [23].

Compute the local non-stationarities of the beamformer output and reference signals,  $\Lambda(Y)$  and  $\Lambda(\mathbf{U})$ , using Eqs. (16) and (24), and compute the transient beam-to-reference ratio,  $\Omega(Y, \mathbf{U})$ , using Eq. (30).

Using the block diagram in Fig. 4, determine which case applies to each frequency bin; Set the *a priori* signal absence probability  $\hat{q}(k, \ell)$  to 1 in Case 1, and to 0 in Case 2, and compute its value according to Eq. (40) in Case 3.

Compute the *a priori* SNR  $\hat{\xi}(k, \ell)$  using Eq. (44), the conditional gain  $G_{H_1}(k, \ell)$  using Eq. (45), and the signal presence probability  $p(k, \ell)$  using Eq. (43).

Compute the time-varying smoothing parameter  $\tilde{\alpha}_d(k, \ell)$  using Eq. (47), and update the noise spectrum estimate  $\hat{\lambda}_d(k, \ell + 1)$  using Eq. (46).

Compute the OM-LSA estimate of the clean signal,  $\hat{X}(k, \ell)$ , using Eqs. (48) and (49).

Fig. 6. The multi-channel post-filtering algorithm.

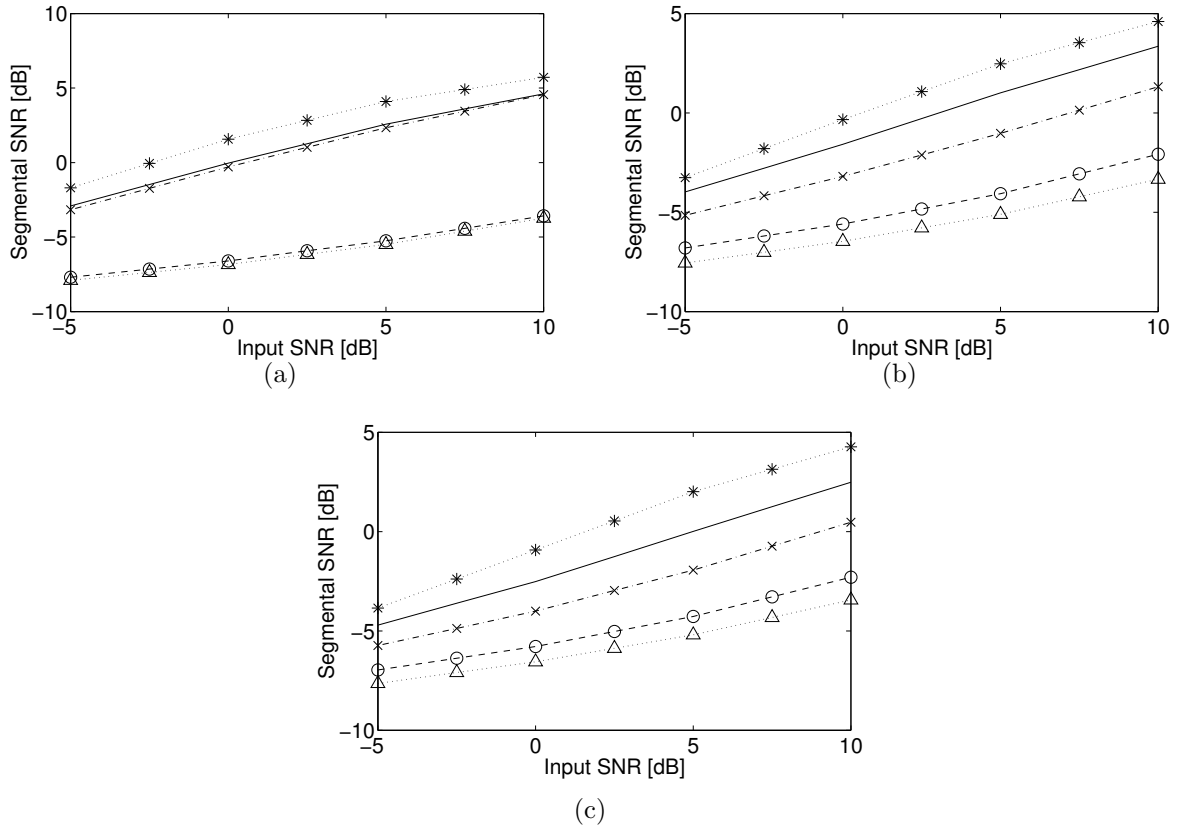


Fig. 7. Average segmental SNR at ( $\triangle$ ) microphone #1, ( $\circ$ ) beamformer output, ( $\times$ ) single-channel post-filtering output, (solid line) multi-channel post-filtering output, and ( $*$ ) theoretical limit post-filtering output, for various car noise conditions: (a) Closed windows; (b) Open window; (c) Interfering speaker.

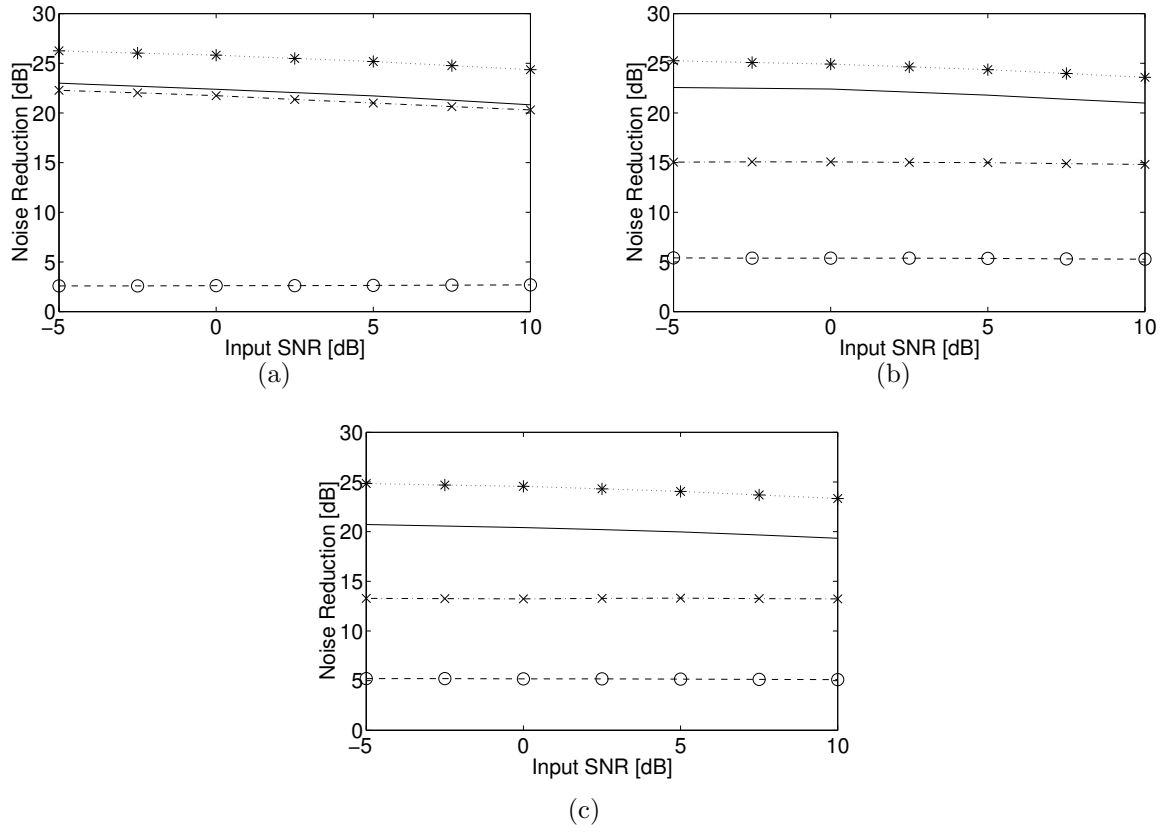


Fig. 8. Average noise reduction at (o) beamformer output, (x) single-channel post-filtering output, (solid line) multi-channel post-filtering output, and (\*) theoretical limit post-filtering output, for various car noise conditions: (a) Closed windows; (b) Open window; (c) Interfering speaker.

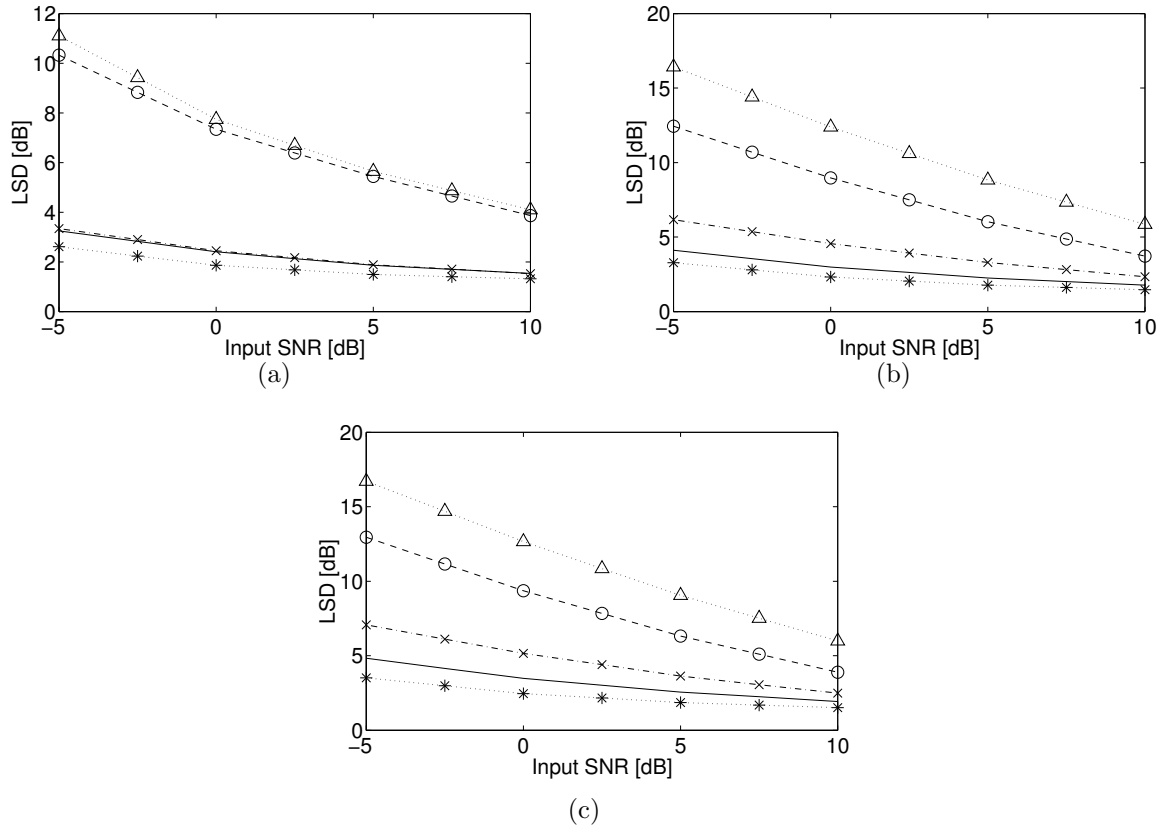


Fig. 9. Average log-spectral distance at ( $\Delta$ ) microphone #1, ( $\circ$ ) beamformer output, ( $\times$ ) single-channel post-filtering output, (solid line) multi-channel post-filtering output, and ( $*$ ) theoretical limit post-filtering output, for various car noise conditions: (a) Closed windows; (b) Open window; (c) Interfering speaker.

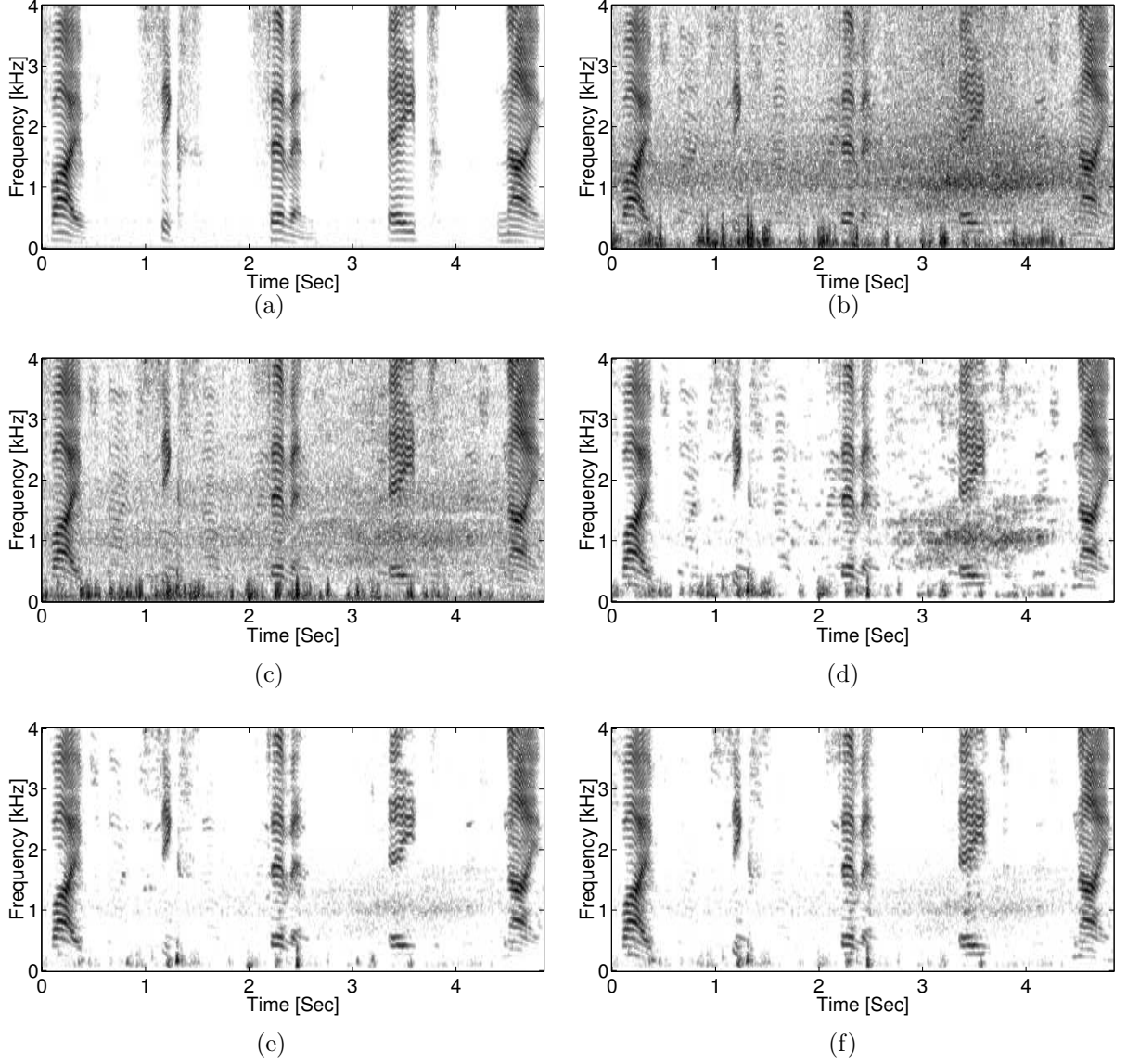


Fig. 10. Speech spectrograms. (a) Original clean speech signal at microphone #1: “Five six seven eight nine.”; (b) Noisy signal at microphone #1 (car noise, open window, interfering speaker.  $\text{SNR} = -0.9$  dB,  $\text{SegSNR} = -6.2$  dB,  $\text{LSD} = 15.4$  dB); (c) Beamformer output ( $\text{SegSNR} = -5.3$  dB,  $\text{NR} = 5.2$  dB,  $\text{LSD} = 12.2$  dB); (d) Single-channel post-filtering output ( $\text{SegSNR} = -3.8$  dB,  $\text{NR} = 12.1$  dB,  $\text{LSD} = 7.4$  dB); (e) Multi-channel post-filtering output ( $\text{SegSNR} = -1.3$  dB,  $\text{NR} = 23.2$  dB,  $\text{LSD} = 4.6$  dB); (f) Theoretical limit ( $\text{SegSNR} = -0.4$  dB,  $\text{NR} = 24.0$  dB,  $\text{LSD} = 4.0$  dB).

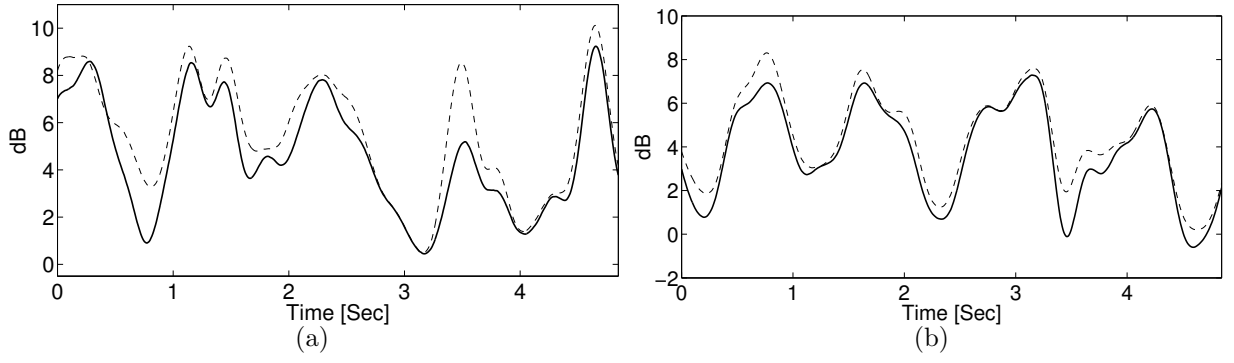


Fig. 11. Trace of the improvement over a single-channel post-filtering gained by the proposed multi-channel post-filtering (solid) and theoretical limit (dashed): (a) Increase in segmental SNR; (b) Decrease in Log-Spectral Distance.