

A NOVEL PERSPECTIVE ON STEREOPHONIC ACOUSTIC ECHO CANCELLATION

Cristian Stanciu[†], Jacob Benesty^{*}, Constantin Paleologu[†], Tomas Gaensler[‡], and Silviu Ciochină[†]

[†] University Politehnica of Bucharest, Romania, e-mail: {cristian, pale, silviu}@comm.pub.ro

^{*} INRS-EMT, University of Quebec, Montreal, Canada, e-mail: benesty@emt.inrs.ca

[‡] mh Acoustics, Summit, NJ, USA, e-mail: tfg@mhacoustics.com

ABSTRACT

The stereophonic acoustic echo is due to the coupling between two loudspeakers and two microphones. In the classical approach, this configuration is modelled by a two-input/two-output system with real random variables. In this paper, we propose to redesign this scheme as a single-input/single-output system with complex random variables. In this framework, we illustrate the behavior of some basic adaptive algorithms and present a distortion method which is more suitable for this model.

Index Terms— Stereophonic acoustic echo cancellation (SAEC), widely linear (WL) model, nonlinear distortion, adaptive filters.

1. INTRODUCTION

In hands-free teleconferencing systems, stereo transmission provides telepresence thanks to our binaural hearing system. These stereophonic systems give a realistic presence that actual single-channel systems cannot offer [1], [2]. In this context, stereophonic acoustic echo cancellation (SAEC) is necessary for full-duplex quality communication. For each microphone in the receiving (i.e., near-end) location, the SAEC consists of the identification of a two-input unknown system, consisting of the parallel combination of two acoustic echo paths (from the two loudspeakers to the microphone). Therefore, in the usual approach, an SAEC system consists of four adaptive filters aiming at identifying four echo paths from two loudspeakers to two microphones.

Despite the inherent similarities, SAEC is fundamentally different (and also more difficult) as compared to single-channel acoustic echo cancellation. The main challenge of SAEC is that the two channels may carry linearly related signals, which in turn may make the normal equation to be solved by the adaptive algorithm singular. This implies that there is no unique solution to the equation (as in the single-channel case) but an infinite number of solutions [3]. It was demonstrated that the only practical solution to the nonuniqueness problem is to reduce the coherence between the input (loudspeaker) signals [4]. Consequently, we need to distort these signals but without affecting too much the stereo perception and the sound quality.

In this paper, we propose a different approach for SAEC by recasting the classical two-input/two-output scheme with real random variables as a single-input/single-output system with complex random variables. In this framework, we present some basic adaptive algorithms and also a nonlinear distortion method which could be more suitable in this context.

This work was supported under the Grant POSDRU/107/1.5/S/76903 and Grant UEFISCDI PN-II-RU-TE no. 7/5.08.2010.

2. A NOVEL MODEL FOR SAEC

The stereophonic acoustic echo is usually modelled by a two-input/two-output system. In this classical setup, we have two input or loudspeaker signals denoted by $x_L(n)$ and $x_R(n)$ (i.e., “left” and “right”), and two output or microphone signals denoted by $d_L(n)$ and $d_R(n)$, where n is the time index. The microphone signals can be expressed as

$$d_L(n) = y_L(n) + v_L(n), \quad (1)$$

$$d_R(n) = y_R(n) + v_R(n), \quad (2)$$

where $y_L(n)$ and $y_R(n)$ are the stereo echo signals, and $v_L(n)$ and $v_R(n)$ are the near-end signals. The echo signals can be obtained as [3], [4]

$$y_L(n) = \mathbf{h}_{t,LL}^T \mathbf{x}_L(n) + \mathbf{h}_{t,RL}^T \mathbf{x}_R(n), \quad (3)$$

$$y_R(n) = \mathbf{h}_{t,LR}^T \mathbf{x}_L(n) + \mathbf{h}_{t,RR}^T \mathbf{x}_R(n), \quad (4)$$

where $\mathbf{h}_{t,LL}$, $\mathbf{h}_{t,RL}$, $\mathbf{h}_{t,LR}$, $\mathbf{h}_{t,RR}$ are L -dimensional vectors of the loudspeaker-to-microphone (“true”) acoustic impulse responses, the superscript T denotes transpose of a vector or a matrix, and the vectors

$$\mathbf{x}_L(n) = [x_L(n) \quad x_L(n-1) \quad \cdots \quad x_L(n-L+1)]^T$$

$$\mathbf{x}_R(n) = [x_R(n) \quad x_R(n-1) \quad \cdots \quad x_R(n-L+1)]^T$$

contain the most recent L samples of the loudspeaker signals. Consequently, the main goal of this application is to estimate the four acoustic impulse responses (i.e., $\mathbf{h}_{t,LL}$, $\mathbf{h}_{t,RL}$, $\mathbf{h}_{t,LR}$, $\mathbf{h}_{t,RR}$) from the microphone signals in order to cancel the echo due to the coupling between the loudspeakers and the microphones.

In the context of acoustic echo cancellation, the loudspeaker and microphone signals are all real random variables. In order to introduce the proposed model, let us use the complex notation

$$d(n) = d_L(n) + jd_R(n) = y(n) + v(n), \quad (5)$$

where $j = \sqrt{-1}$, $y(n) = y_L(n) + jy_R(n)$, and $v(n) = v_L(n) + jv_R(n)$. Furthermore, let us define the complex random vector

$$\begin{aligned} \mathbf{x}(n) &= [x(n) \quad x(n-1) \quad \cdots \quad x(n-L+1)]^T \\ &= \mathbf{x}_L(n) + j\mathbf{x}_R(n), \end{aligned} \quad (6)$$

where $x(n) = x_L(n) + jx_R(n)$, so that the complex echo signal can be expressed as

$$y(n) = \mathbf{h}_t^H \mathbf{x}(n) + \mathbf{h}_t'^H \mathbf{x}^*(n), \quad (7)$$

where the superscripts H and $*$ denote transpose-conjugate and conjugate, respectively, and

$$\mathbf{h}_t = \mathbf{h}_{t,1} + j\mathbf{h}_{t,2}, \quad (8)$$

$$\mathbf{h}'_t = \mathbf{h}'_{t,1} + j\mathbf{h}'_{t,2}, \quad (9)$$

with

$$\begin{aligned} \mathbf{h}_{t,1} &= \frac{\mathbf{h}_{t,LL} + \mathbf{h}_{t,RR}}{2}, & \mathbf{h}_{t,2} &= \frac{\mathbf{h}_{t,RL} - \mathbf{h}_{t,LR}}{2}, \\ \mathbf{h}'_{t,1} &= \frac{\mathbf{h}_{t,LL} - \mathbf{h}_{t,RR}}{2}, & \mathbf{h}'_{t,2} &= -\frac{\mathbf{h}_{t,RL} + \mathbf{h}_{t,LR}}{2}. \end{aligned}$$

Consequently, (7) can be rewritten as

$$y(n) = \tilde{\mathbf{h}}_t^H \tilde{\mathbf{x}}(n), \quad (10)$$

where

$$\tilde{\mathbf{h}}_t = \begin{bmatrix} \mathbf{h}_t \\ \mathbf{h}'_t \end{bmatrix}, \quad \tilde{\mathbf{x}}(n) = \begin{bmatrix} \mathbf{x}(n) \\ \mathbf{x}^*(n) \end{bmatrix}.$$

Finally, the complex reference signal (7) becomes

$$d(n) = \tilde{\mathbf{h}}_t^H \tilde{\mathbf{x}}(n) + v(n). \quad (11)$$

In this context, our new goal is to estimate the complex acoustic impulse response $\tilde{\mathbf{h}}_t$ (of length $2L$) from the complex microphone signal, $d(n)$, and the complex loudspeaker signal, $x(n)$. In fact, the classical two-input/two-output system with real random variables has been converted to a single-input/single-output system with complex random variables. Looking of (7) or (10), we can recognize the widely linear (WL) model for complex random variables proposed in [5]; also, this approach is in consistence with the duality principle explained in [6].

3. SOME BASIC ADAPTIVE ALGORITHMS

Let $\tilde{\mathbf{h}}(n) = [\tilde{h}_0(n) \ \tilde{h}_1(n) \ \cdots \ \tilde{h}_{2L-1}(n)]^T$, be an adaptive filter of length $2L$, which is an estimate of $\tilde{\mathbf{h}}_t$, and let

$$\hat{y}(n) = \tilde{\mathbf{h}}^H(n-1)\tilde{\mathbf{x}}(n) \quad (12)$$

be the output of the adaptive filter at time n . Thus, the error signal is

$$e(n) = d(n) - \hat{y}(n). \quad (13)$$

Based on (12) and (13), we can write the update of the normalized least-mean-square (NLMS) algorithm as

$$\tilde{\mathbf{h}}(n) = \tilde{\mathbf{h}}(n-1) + \frac{\alpha \tilde{\mathbf{x}}(n)e^*(n)}{\delta + \tilde{\mathbf{x}}^H(n)\tilde{\mathbf{x}}(n)}, \quad (14)$$

where α is the normalized stepsize parameter ($0 < \alpha < 2$) and $\delta \geq 0$ is the regularization constant. The NLMS algorithm could be useful in practice mainly due to its simplicity. However, it converges slowly for long length adaptive filters and highly correlated inputs.

In order to improve the convergence rate, we can take advantage of the sparseness character of the echo paths, which inspired the idea to ‘‘proportionate’’ the algorithm behavior [7]. In other words, we can update each coefficient of the filter independently of the others, by adjusting the adaptation stepsize in proportion to the magnitude of the estimated filter coefficient. Hence, the adaptation gain is ‘‘proportionately’’ redistributed among all the coefficients to emphasize the large ones in order to speed up their convergence and,

consequently, to increase the overall convergence rate. Among the many proportionate-type algorithms developed for echo cancellation (e.g., see [8] and the references therein), the improved proportionate NLMS (IPNLMS) algorithm [9] is one of the most attractive choice. The good features of this algorithm include its simplicity and the robustness to the sparseness degree of the echo path.

In the context of the proposed model for SAEC, the update of the IPNLMS algorithm can be expressed as

$$\tilde{\mathbf{h}}(n) = \tilde{\mathbf{h}}(n-1) + \frac{\alpha \mathbf{G}(n-1)\tilde{\mathbf{x}}(n)e^*(n)}{\delta + \tilde{\mathbf{x}}^H(n)\mathbf{G}(n-1)\tilde{\mathbf{x}}(n)}, \quad (15)$$

where

$$\mathbf{G}(n-1) = \text{diag}[g_0(n-1), g_1(n-1), \dots, g_{2L-1}(n-1)], \quad (16)$$

is a diagonal matrix (of size $2L \times 2L$) containing the proportionate (or gain) factors, which are evaluated as

$$g_l(n-1) = \frac{1-\kappa}{4L} + (1+\kappa) \frac{|\tilde{h}_l(n-1)|}{2 \sum_{i=0}^{2L-1} |\tilde{h}_i(n-1)|}, \quad 0 \leq l \leq 2L-1, \quad (17)$$

where κ ($-1 \leq \kappa < 1$) is a parameter that controls the amount of proportionality in the IPNLMS algorithm [9].

Another very good candidate for echo cancellation is the affine projection algorithm (APA) [10], since it converges and tracks faster than the NLMS algorithm. Besides, it can be efficient from an arithmetic complexity viewpoint as compared to more complex algorithms from the recursive least-squares (RLS) family. In order to derive the APA in the context of the proposed model, let us write the $2L \times P$ input matrix

$$\tilde{\mathbf{X}}(n) = [\tilde{\mathbf{x}}(n) \ \tilde{\mathbf{x}}(n-1) \ \cdots \ \tilde{\mathbf{x}}(n-P+1)],$$

where P is the projection order. Also, we can define the $P \times 1$ a priori error vector as

$$\mathbf{e}(n) = \mathbf{d}(n) - \tilde{\mathbf{X}}^T(n)\tilde{\mathbf{h}}^*(n-1), \quad (18)$$

where $\mathbf{d}(n) = [d(n) \ d(n-1) \ \cdots \ d(n-P+1)]^T$. Using this notation, the update of the APA is

$$\begin{aligned} \tilde{\mathbf{h}}(n) &= \tilde{\mathbf{h}}(n-1) + \alpha \tilde{\mathbf{X}}(n) [\delta \mathbf{I}_P + \tilde{\mathbf{X}}^H(n)\tilde{\mathbf{X}}(n)]^{-1} \mathbf{e}^*(n), \end{aligned} \quad (19)$$

where \mathbf{I}_P is the $P \times P$ identity matrix. It can be noticed that, by taking $P = 1$, we obtain the update of the NLMS algorithm (14).

The ‘‘proportionate’’ idea can be also extended in the case of APA, in order to further increase its performance when identifying sparse impulse responses. For example, using the gain factors of the IPNLMS algorithm, we can derive the improved proportionate APA (IPAPA):

$$\begin{aligned} \tilde{\mathbf{h}}(n) &= \tilde{\mathbf{h}}(n-1) + \alpha \mathbf{G}(n-1)\tilde{\mathbf{X}}(n) \times \\ &\quad [\delta \mathbf{I}_P + \tilde{\mathbf{X}}^H(n)\mathbf{G}(n-1)\tilde{\mathbf{X}}(n)]^{-1} \mathbf{e}^*(n), \end{aligned} \quad (20)$$

where $\mathbf{G}(n-1)$ is defined in (16) and (17). Clearly, for $P = 1$ we find the IPNLMS algorithm [see (15)].

Of course, many other adaptive algorithms can be derived in the context of the proposed model for SAEC. However, due to the lack of space, we limit our presentation to these four basic algorithms, i.e., NLMS, IPNLMS, APA, and IPAPA.

4. SOLUTIONS TO THE NONUNIQUENESS PROBLEM

It is well known that in the SAEC problem, most of the time, the two input signals [i.e., $x_L(n)$ and $x_R(n)$] are obtained by filtering a common source, so that a problem of nonuniqueness is expected [3]. Also, it was found that preprocessing of these far-end loudspeaker signals that actually are transmitted to the near-end room is the only way to achieve a unique solution [4]. In other words, it may be required to distort the input signals $x_L(n)$ and $x_R(n)$, in order to reduce the coherence between these two signals, which can lead to the estimation of the true acoustic impulse responses. However, this distortion must be performed in such a way that the quality of the signals and the stereo effect are not degraded.

A simple but efficient method uses positive and negative half-wave rectifiers on each channel respectively [4], i.e.,

$$x'_L(n) = x_L(n) + \alpha_r \frac{x_L(n) + |x_L(n)|}{2}, \quad (21)$$

$$x'_R(n) = x_R(n) + \alpha_r \frac{x_R(n) - |x_R(n)|}{2}, \quad (22)$$

where α_r is a parameter used to control the amount of nonlinearity. Experiments show that stereo perception is not affected by this method even with α_r as large as 0.5.

In the context of the proposed model, the complex input signal can be expressed as

$$x(n) = x_L(n) + jx_R(n) = e^{j\theta_r(n)} |x(n)|, \quad (23)$$

where $\theta_r(n)$ [with $\tan \theta_r(n) = x_R(n)/x_L(n)$] and $|x(n)| = \sqrt{x_L^2(n) + x_R^2(n)}$ are the phase and module of $x(n)$, respectively. In this formulation, we represent the stereo perception with $\theta_r(n)$ and the quality of the stereo signals with $|x(n)|$. A modification of $\theta_r(n)$ only, will mostly affect the stereo effect of $x(n)$; while a modification of $|x(n)|$ will mostly affect the quality of the stereo signals.

Similarly, using the complex notation, (21) and (22) can be expressed as

$$x'(n) = x'_L(n) + jx'_R(n) = e^{j\theta'_r(n)} |x'(n)|, \quad (24)$$

where $\tan \theta'_r(n) = x'_R(n)/x'_L(n)$ and $|x'(n)| = \sqrt{x'^2_L(n) + x'^2_R(n)}$. In order to preserve the quality of the stereo signals, we propose not to modify the module of the complex input signal $x(n)$, but only to change its phase. Therefore, we can use the new following transformations [2]:

$$x''_L(n) = \cos \theta'_r(n) |x(n)|, \quad (25)$$

$$x''_R(n) = \sin \theta'_r(n) |x(n)|, \quad (26)$$

where the phase $\theta'_r(n)$ is computed from the half-wave rectifiers [see (24)] while the module corresponds to the module of the original signals.

5. SIMULATION RESULTS

Simulations are performed in the context of the proposed model for SAEC. The acoustic impulse responses in the far-end location have 2048 coefficients, while the length of the impulse responses in the near-end location [i.e., $h_{t,LL}(n)$, $h_{t,RL}(n)$, $h_{t,LR}(n)$, and $h_{t,RR}(n)$] is $L = 512$. The length of the adaptive filter $\tilde{\mathbf{h}}(n)$ is $2L = 1024$ and the sampling rate is 8 kHz.

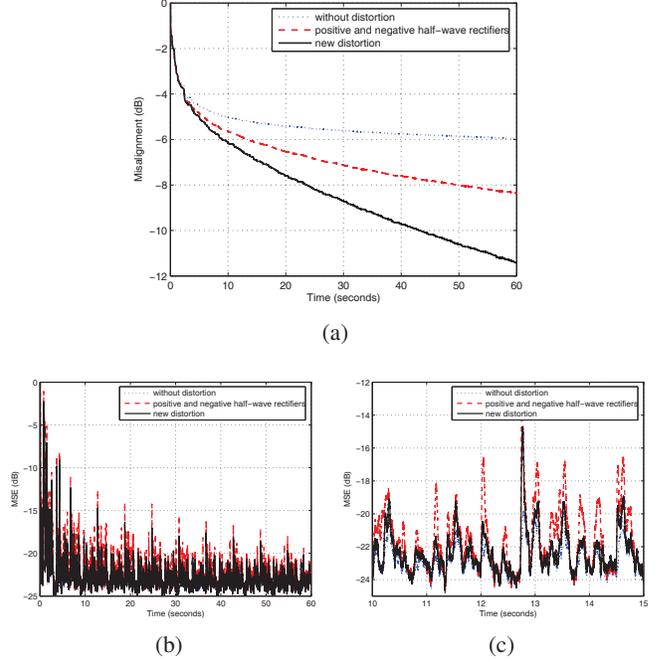


Fig. 1. Results of the NLMS algorithm for different types of distortion with $\alpha_r = 0.3$. (a) Misalignment; (b) MSE; (c) MSE detail.

The source signal in the far-end location is a speech sequence. All simulations are performed in the single-talk scenario, i.e., absence of a near-end talker. In this case, the near-end signal $v(n)$ consists only of the background noise. We can define the stereo echo-to-noise ratio (SENr) [which is equivalent to the signal-to-noise ratio (SNR)] as $\text{SENr} = \sigma_y^2 / \sigma_v^2$, where $\sigma_y^2 = E[|y(n)|^2]$ and $\sigma_v^2 = E[|v(n)|^2]$ are the variances of $y(n)$ and $v(n)$, respectively. In our simulations, the background noise in the near-end is an independent white Gaussian signal and its level is set such that $\text{SENr} = 30$ dB.

We choose for comparisons the four algorithms presented in Section 3, i.e., NLMS, IPNLMS, APA, and IPAPA. The stepsize for all the algorithms is $\alpha = 0.25$ and the regularization constants are $\delta_{\text{NLMS}} = \delta_{\text{APA}} = 20\sigma_x^2$ and $\delta_{\text{IPNLMS}} = \delta_{\text{IPAPA}} = 20\sigma_x^2 / (2L)$ [11], where $\sigma_x^2 = E[|x(n)|^2]$ is the variance of $x(n)$. The proportionate-type algorithms (i.e., IPNLMS and IPAPA) use $\kappa = 0$. The performance of the algorithms is evaluated in terms of two measures, i.e., (a) the normalized misalignment (in dB), defined as $20 \log_{10} \left\| \tilde{\mathbf{h}}_t - \tilde{\mathbf{h}}(n) \right\|_2 / \left\| \tilde{\mathbf{h}}_t \right\|_2$ (with $\|\cdot\|_2$ denoting the ℓ_2 norm) and (b) the mean-square error (MSE) averaged over 256 points for the purpose of smoothing the results.

In all the experiments, we compare the performance of the algorithms using positive and negative half-wave rectifiers [see (21) and (22)] versus the new proposed distortion [see (25) and (26)]; the distortion parameter is set to $\alpha_r = 0.3$. Also, the case without distortion is shown as a reference.

Figure 1 presents the performance of the NLMS algorithm. It can be noticed from Fig. 1(a) that the misalignment is greatly reduced by the new distortion. Also, as we can see in Fig. 1(b) and in the detail presented in Fig. 1(c), the new distortion leads to a better performance in terms of the MSE as compared to the positive and

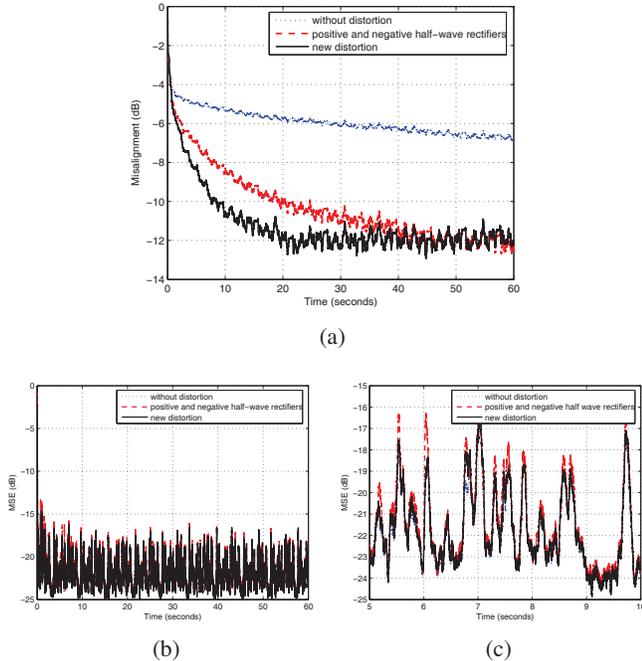


Fig. 2. Results of the APA (using $P = 8$) for different types of distortion with $\alpha_r = 0.3$. (a) Misalignment; (b) MSE; (c) MSE detail.

negative half-wave rectifiers method.

Figure 2 shows the performance of the APA with $P = 8$; it was found that this value of the projection order offers a proper compromise between the performance and complexity. It can be noticed from Fig. 2(a) that the APA converges faster with the new distortion, outperforming by far the NLMS algorithm [see for comparison Fig. 1(a)]. Also, as we can see in Fig. 2(b) and in the detail presented in Fig. 2(c), the new distortion leads to a slightly better performance in terms of the MSE as compared to the positive and negative half-wave rectifiers.

Since the IPAPA has resulted as a combination between the IPNLMS algorithm and the APA, it is expected that the IPAPA should outperform both its predecessors. The last experiment outlines this aspect, by comparing these three algorithms in a tracking situation (the impulse responses in the near-end location are shifted to the right by 12 samples). The new distortion is used with $\alpha_r = 0.3$. The projection order is $P = 8$ for the APA and IPAPA. The results are shown in Fig. 3. According to these plots, it is clear that IPAPA outperforms both the IPNLMS and APA.

6. CONCLUSIONS

In this paper, we proposed to recast the SAEC problem as a single-input/single-output system with complex random variables. As a consequence, the four real-valued acoustic impulse responses are converted to one complex-valued impulse response. The main advantage of this approach is that instead of handling two (real) output signals separately, we only handle one (complex) output signal. In this framework, we have presented some typical adaptive algorithms and a new distortion method suitable for this model.

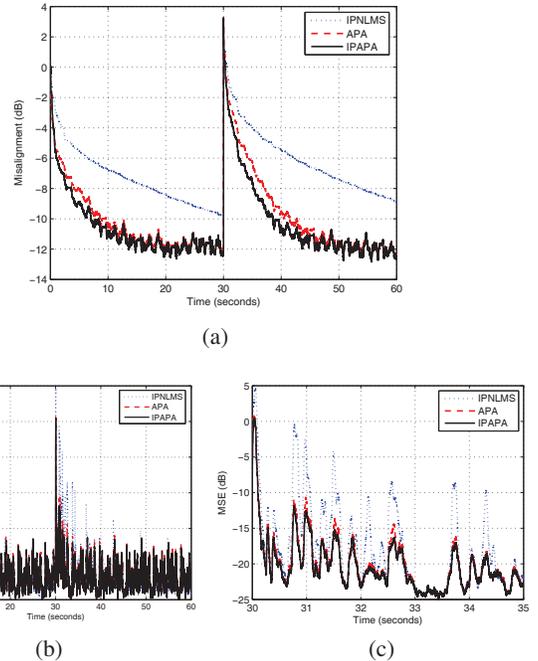


Fig. 3. Results of the IPNLMS, APA, and IPAPA (using $P = 8$) for different types of distortion with $\alpha_r = 0.3$. Echo paths changes at time 30 seconds. (a) Misalignment; (b) MSE; (c) MSE detail.

7. REFERENCES

- [1] J. Benesty, T. Gaensler, D. R. Morgan, M. M. Sondhi, and S. L. Gay, *Advances in Network and Acoustic Echo Cancellation*. Berlin, Germany: Springer-Verlag, 2001.
- [2] J. Benesty, C. Paleologu, T. Gänsler, and S. Ciochină, *A Perspective on Stereophonic Acoustic Echo Cancellation*. Springer-Verlag, Berlin, Germany, 2011.
- [3] M. M. Sondhi, D. R. Morgan, and J. L. Hall, "Stereophonic acoustic echo cancellation—An overview of the fundamental problem," *IEEE Signal Process. Lett.*, vol. 2, pp. 148–151, Aug. 1995.
- [4] J. Benesty, D. R. Morgan, and M. M. Sondhi, "A better understanding and an improved solution to the specific problems of stereophonic acoustic echo cancellation," *IEEE Trans. Speech, Audio Process.*, vol. 6, pp. 156–165, Mar. 1998.
- [5] B. Picinbono and P. Chevalier, "Widely linear estimation with complex data," *IEEE Trans. Signal Process.*, vol. 43, pp. 2030–2033, Aug. 1995.
- [6] D. P. Mandic, S. Still, and S. C. Douglas, "Duality between widely linear and dual channel adaptive filtering," in *Proc. IEEE ICASSP*, 2009, pp. 1729–1732.
- [7] D. L. Duttweiler, "Proportionate normalized least-mean-squares adaptation in echo cancelers," *IEEE Trans. Speech, Audio Process.*, vol. 8, pp. 508–518, Sept. 2000.
- [8] C. Paleologu, J. Benesty, and S. Ciochină, *Sparse Adaptive Filters for Echo Cancellation*. Morgan & Claypool Publishers, 2010.
- [9] J. Benesty and S. L. Gay, "An improved PNLMS algorithm," in *Proc. IEEE ICASSP*, 2002, pp. 1881–1884.
- [10] K. Ozeki and T. Umeda, "An adaptive filtering algorithm using an orthogonal projection to an affine subspace and its properties," *Electron. Commun. Jpn.*, vol. 67-A, pp. 19–27, May 1984.
- [11] J. Benesty, C. Paleologu, and S. Ciochină, "On regularization in adaptive filtering," *IEEE Trans. Audio, Speech, Language Process.*, vol. 19, pp. 1734–1742, Aug. 2011.