

DELAY-PERFORMANCE TRADEOFFS IN CAUSAL MICROPHONE ARRAY PROCESSING

Ryan M. Corey, Naoki Tsuda, and Andrew C. Singer

University of Illinois at Urbana-Champaign

ABSTRACT

In real-time listening enhancement applications, such as hearing aid signal processing, sounds must be processed with no more than a few milliseconds of delay to sound natural to the listener. Listening devices can achieve better performance with lower delay by using microphone arrays to filter acoustic signals in both space and time. Here, we analyze the tradeoff between delay and squared-error performance of causal multichannel Wiener filters for microphone array noise reduction. We compute exact expressions for the delay-error curves in two special cases and present experimental results from real-world microphone array recordings. We find that delay-performance characteristics are determined by both the spatial and temporal correlation structures of the signals.

Index Terms— Microphone arrays, audio enhancement, audio source separation, hearing aids, noise reduction, beamforming

1. INTRODUCTION

Listening enhancement applications, such as hearing aid processing [1] and audio augmented reality [2], differ from other audio enhancement applications, like teleconferencing and speech recognition, in part because of their strict delay constraints. Since users hear both live and processed signals simultaneously, these systems must process sound with no more than a few milliseconds of delay. Discerning listeners can notice delays as low as 3 ms and are disturbed by delays greater than 10 ms [3]. Listeners with hearing loss can tolerate greater delay, around 20 ms for closed-fitting hearing aids [4] and 6 ms for open-fitting hearing aids [5]. Delays longer than about 30 ms can impair the user’s ability to speak [6].

This delay requirement limits the performance of audio enhancement systems. In single-channel systems, the frequency resolution of a frequency-selective filter generally improves with longer delay. Modern single-microphone audio enhancement algorithms [7], such as those employing time-frequency masks [8] and non-negative matrix factorization [9], often process speech using short-time Fourier transform (STFT) frames of 60 ms or longer to maximize time-frequency sparsity [8]. These algorithms are effective in many applications, but their delay is too large for listening enhancement.

Multichannel audio enhancement systems use microphone arrays to spatially separate signals [10–12]. Many multichannel methods are also applied in the STFT domain to more easily model reverberation [12, 13]. In principle, however, spatial processing should require minimal delay: for example, a linear array can enhance a source at broadside with zero delay by simply summing its inputs. Whereas the frequency resolution of a temporal filter depends on its duration, the spatial resolution of an array is determined by its spatial extent. Multichannel listening systems can use both spatial and

spectral diversity to separate signals. It is natural to ask, therefore, whether devices with large arrays can enhance audio with lower delay than those with small arrays. That is, *can we use array processing to trade space for time?*

There is a large body of literature on array processing for listening devices, e.g. [14, 15], and causal multichannel filters have been studied in the contexts of dereverberation [16–19] and noise and echo control [20]. In [21], the authors considered the minimum filter delay required to cover the full aperture of an array. There have also been several proposed low-delay single-microphone filtering and source separation techniques [22–24]. However, to the best of our knowledge, there has been no prior study of delay-performance tradeoffs in array processing.

Here we approach audio enhancement as a stationary linear estimation problem: given an observed signal from the infinite past to time t , what is the linear minimum mean square error (MSE) estimate of a desired signal at time $t - \alpha$? Positive values of α correspond to delay and negative values to prediction. Such problems are well understood in the scalar case: for certain signals, we can use spectral factorization to compute exact expressions for the MSE as a function of α [25–27]. For example, Figure 1 shows delay-error curves for separating several spectrally distinct speechlike sounds, which will be described in Section 3. As α increases, the MSE decreases from the variance of the target signal to the MSE of a noncausal Wiener filter. We can apply similar theoretical tools in the multivariate case [28, 29] to analyze delay-performance tradeoffs for causal multichannel Wiener filters (CMWF) in terms of the spatial and temporal correlation structures of the source signals. In this work, we will derive a general expression for the MSE performance of a CMWF as a function of α , find exact expressions for idealized mixing models, and present experimental results from wearable and distributed microphone arrays in a real room.

2. DELAY-CONSTRAINED MULTICHANNEL FILTERING

Consider a mixture of N sources captured by M microphones. Let the sources $\mathbf{s}(t) = [s_1(t), \dots, s_N(t)]^T$ and additive noise $\mathbf{z}(t) = [z_1(t), \dots, z_M(t)]^T$ be wide-sense stationary continuous-time random processes that are uncorrelated with each other. Let $a_{m,n}(t)$, $m = 1, \dots, M$, $n = 1, \dots, N$ be known causal impulse responses and let $\mathbf{w}_\alpha^T(t) = [w_{\alpha,1}(t), \dots, w_{\alpha,M}(t)]$ be filter impulse responses. Denote the observed signals by $\mathbf{x}(t) = [x_1(t), \dots, x_M(t)]^T$ and the system output by $y_\alpha(t)$, where

$$x_m(t) = \sum_{n=1}^N (a_{m,n} * s_n)(t) + z_m(t), \quad m = 1, \dots, M, \text{ and } (1)$$

$$y_\alpha(t) = \sum_{m=1}^M (w_{\alpha,m} * x_m)(t), \quad (2)$$

This material is based upon work supported by the National Science Foundation Graduate Research Fellowship Program under Grant Number DGE-1144245.

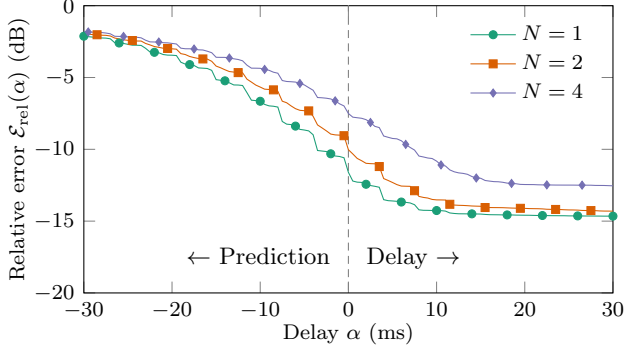


Fig. 1. Relative MSE as a function of delay for isolating one source from a mixture of N synthetic speechlike sounds (see Section 3) and uncorrelated noise using single-channel Wiener filters.

and $*$ denotes linear convolution. We define the desired output signal $d_\alpha(t)$ to be the first source as captured by the first microphone—for example, a target talker reproduced at the microphone nearest the listener’s ear—and delayed by time α :

$$d_\alpha(t) = (a_{11} * s_1)(t - \alpha). \quad (3)$$

To understand fundamental tradeoffs in performance, we restrict our attention to the best-case scenario in which all signals are stationary in both space and time and have known statistics. Let $\mathbf{A}(\omega)$ be the $M \times N$ frequency response matrix corresponding to the $a_{m,n}(t)$ ’s. Let $\mathbf{r}_s(t)$, $\mathbf{r}_z(t)$, $r_d(t)$, and $\mathbf{r}_x(t)$ be the autocorrelation sequences of the corresponding random variables and let $\mathbf{R}_s(\omega)$, $\mathbf{R}_z(\omega)$, $R_d(\omega) = |A_{1,1}(\omega)|^2 R_{s_1}(\omega)$, and $\mathbf{R}_x(\omega) = \mathbf{A}(\omega)\mathbf{R}_s(\omega)\mathbf{A}^H(\omega) + \mathbf{R}_z(\omega)$ be their respective Fourier transforms. To ensure that the CMWF is well defined, we assume that $\mathbf{R}_x(\omega)$ is positive definite for all ω of interest. Let $\mathbf{r}_{xd}(t)$ be the cross-correlation of $\mathbf{x}(t)$ with $d_0(t)$ and let $\mathbf{R}_{xd}(\omega) = \mathbf{A}_1(\omega)R_{s_1}(\omega)\mathbf{A}_{1,1}^*(\omega)$ be its Fourier transform, where $\mathbf{A}_1(\omega)$ is the column of $\mathbf{A}(\omega)$ corresponding to the target source. Let $\mathbf{W}_\alpha^T(\omega)$ be the Fourier transform of $\mathbf{w}_\alpha^T(t)$.

2.1. Causal filter performance

The CMWF $\mathbf{w}_\alpha^T(t)$ must satisfy the Wiener-Hopf equation [25],

$$\mathbf{r}_{xd}^T(t - \alpha) = \int_0^\infty \mathbf{w}_\alpha^T(u)\mathbf{r}_x(t - u) du, \quad 0 < t < \infty. \quad (4)$$

The MSE between $y_\alpha(t)$ and $d_\alpha(t)$ is

$$\mathcal{E}(\alpha) = r_d(0) - \int_{-\infty}^\infty \mathbf{w}_\alpha^T(t)\mathbf{r}_{xd}(t - \alpha) dt. \quad (5)$$

The *noncausal* ($\alpha \rightarrow \infty$) solution to (4) and its error power are readily expressed in the frequency domain:

$$\mathbf{W}_{\text{nc}}^T(\omega) = \mathbf{R}_{xd}^H(\omega)\mathbf{R}_x^{-1}(\omega) \quad (6)$$

$$\mathcal{E}_{\text{nc}} = \int_{-\infty}^\infty \left[R_d(\omega) - \mathbf{R}_{xd}^H(\omega)\mathbf{R}_x^{-1}(\omega)\mathbf{R}_{xd}(\omega) \right] \frac{d\omega}{2\pi}. \quad (7)$$

For finite α , we can solve (4) by first decomposing $\mathbf{R}_x(\omega)$ into its *spectral factors* [28],

$$\mathbf{R}_x(\omega) = \mathbf{G}(\omega)\mathbf{G}^H(\omega), \quad (8)$$

where $\mathbf{G}(\omega)$ and its inverse are both causal. We proceed by decorrelating $\mathbf{x}(t)$ using $\mathbf{G}^{-1}(\omega)$ and then solving (4) for the decorrelated signals [29] to find the causal filter

$$\mathbf{W}_\alpha^T(\omega) = \left[e^{-j\omega\alpha} \mathbf{R}_{xd}^H(\omega)(\mathbf{G}^H(\omega))^{-1} \right]_+ \mathbf{G}^{-1}(\omega), \quad (9)$$

where $[\cdot]_+$ denotes the causal part of the argument, that is, time-domain truncation from $t = 0$. Let $\tilde{\mathbf{R}}^T(\omega) = \mathbf{R}_{xd}^H(\omega)(\mathbf{G}^H(\omega))^{-1}$. For the listening enhancement application, this vector can be written

$$\tilde{\mathbf{R}}^T(\omega) = A_{1,1}(\omega)R_{s_1}(\omega)\mathbf{A}_1^H(\omega)(\mathbf{G}^H(\omega))^{-1}. \quad (10)$$

Let $\tilde{\mathbf{r}}^T(t)$ be the inverse Fourier transform of $\tilde{\mathbf{R}}^T(\omega)$. Substituting \mathbf{w}_α^T from (9) into (5), using the spectral factorization (8) and Parseval’s identity, and rearranging terms [27], we can show that

$$\mathcal{E}(\alpha) = \mathcal{E}_{\text{nc}} + \int_{-\infty}^{-\alpha} \tilde{\mathbf{r}}^T(t)\tilde{\mathbf{r}}(t) dt. \quad (11)$$

Thus, *the error penalty due to causality is the energy in $\tilde{\mathbf{r}}(t)$ for $t < -\alpha$* . Our goal is to understand how $\mathcal{E}(\alpha)$ depends on the spatial and spectral characteristics of the source signals. While multivariate spectral factorizations are often difficult to compute in practice [30], we can find exact expressions for certain special cases that provide insight about the delay-constrained array processing problem.

2.2. Uniform linear array

First, consider a plane wave incident upon a uniform linear array of M sensors with the reference at one end. Let τ be the time difference of arrival (TDOA) between adjacent microphones, let $R_s(\omega) = 1$ and let $\mathbf{R}_z(\omega) = \sigma^2 \mathbf{I}$, so that

$$\mathbf{R}_{xd}^H = [1 \quad e^{+j\omega\tau} \quad \dots \quad e^{+j\omega(M-1)\tau}] \quad \text{and} \quad (12)$$

$$\mathbf{R}_x(\omega) = \begin{bmatrix} \sigma^2 + 1 & e^{+j\omega\tau} & \dots & e^{+j\omega(M-1)\tau} \\ e^{-j\omega\tau} & \sigma^2 + 1 & \dots & e^{+j\omega(M-2)\tau} \\ \vdots & \vdots & \ddots & \vdots \\ e^{-j\omega(M-1)\tau} & e^{-j\omega(M-2)\tau} & \dots & \sigma^2 + 1 \end{bmatrix}. \quad (13)$$

A convenient spectral factor is the lower triangular matrix

$$\mathbf{G}(\omega) = \begin{bmatrix} b_1(\sigma^2 + 1) & 0 & \dots & 0 \\ b_1 e^{-j\omega\tau} & b_2(\sigma^2 + 2) & & 0 \\ \vdots & \vdots & \ddots & \vdots \\ b_1 e^{-j\omega(M-1)\tau} & b_2 e^{-j\omega(M-2)\tau} & \dots & b_M(\sigma^2 + M) \end{bmatrix} \quad (14)$$

where $b_m = \sqrt{\sigma^2 / ((\sigma^2 + m)(\sigma^2 + m - 1))}$. Applying (10) and taking the inverse Fourier transform, we have

$$\tilde{\mathbf{r}}^T(t) = [b_1 \quad b_2 \delta(t + \tau) \quad \dots \quad b_M \delta(t + (M - 1)\tau)]. \quad (15)$$

Finally, from (11), the MSE is

$$\begin{aligned} \mathcal{E}(\alpha) &= \frac{\sigma^2}{\sigma^2 + M} + \sum_{m=0}^{M-1} b_{m+1}^2 u(m\tau - \alpha) \\ &= \frac{\sigma^2}{\sigma^2 + \sum_{m=0}^{M-1} \bar{u}(\alpha - m\tau)}, \end{aligned} \quad (16)$$

where $u(t) = 1$ if $t > 0$ and $\bar{u}(t) = 1$ if $t \geq 0$. Thus, the error is reduced for each microphone that the source reaches within time α of reaching the reference. The delay-error curve is a piecewise constant function with steps of width $|\tau|$ and decreasing heights that depend on σ^2 .

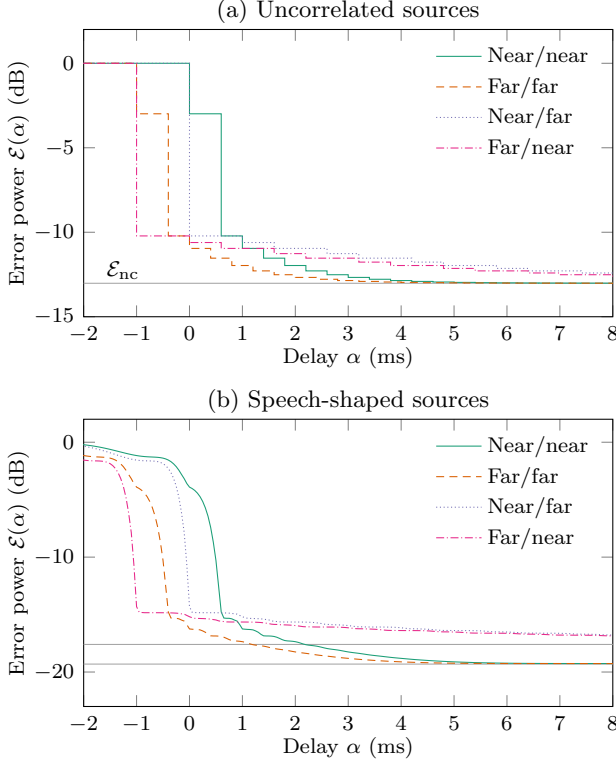


Fig. 2. Delay-error curves for two plane wave sources and two sensors with $|\tau_1| = 1$ ms, $|\tau_2| = 0.6$ ms, and $\sigma^2 = -20$ dB. The legend indicates the placement of the target/interference sources with respect to the reference.

2.3. Two-source, two-microphone separation

We can follow a similar procedure with multiple sources. Consider a scenario with two plane wave sources and two microphones. Let τ_1 and $\tau_2 \neq \tau_1$ be the TDOAs of the sources, let $\mathbf{R}_s(\omega) = \mathbf{I}$ and let $\mathbf{R}_z(\omega) = \sigma^2 \mathbf{I}$ with $\sigma^2 > 0$, so that

$$\mathbf{R}_{xd}^H(\omega) = [1 \ e^{+j\omega\tau_1}], \text{ and} \quad (18)$$

$$\mathbf{R}_x(\omega) = \begin{bmatrix} 2 + \sigma^2 & e^{+j\omega\tau_1} + e^{+j\omega\tau_2} \\ e^{-j\omega\tau_1} + e^{-j\omega\tau_2} & 2 + \sigma^2 \end{bmatrix}. \quad (19)$$

The determinant of $\mathbf{R}_x(\omega)$ can be written

$$\det \mathbf{R}_x(\omega) = \gamma^{-1} \left| 1 - \gamma e^{-j\omega(\tau_1 - \tau_2)} \right|^2, \quad (20)$$

where γ is a scalar that depends only on σ^2 . The spectral factorization of $\mathbf{R}_x(\omega)$ takes different forms depending on the signs of τ_1 and τ_2 , but $\hat{\mathbf{R}}^T(\omega)$ always includes a term of the form $(1 - \gamma e^{+j\omega|\tau_1 - \tau_2|})^{-1}$, which results in an infinite-duration $\tilde{\mathbf{r}}^T(t)$. Applying (11), we find that

$$\mathcal{E}(\alpha) = \begin{cases} \mathcal{E}_{nc} + \frac{u(t_0 - \alpha) + c_1^2 \gamma u(t_1 - |\tau_1 - \tau_2| - \alpha) + c_2^2 f(t_1)}{\sigma^2 + 2}, & \text{if } \tau_1 \tau_2 > 0 \\ \mathcal{E}_{nc} + \sqrt{\gamma} u(t_0 - \alpha) + f(t_0 - |\tau_1|) + \gamma f(t_1), & \text{if } \tau_1 \tau_2 \leq 0 \end{cases} \quad (21)$$

where $t_0 = \min(0, \tau_1)$, $t_1 = \max(0, \tau_1, \tau_2, \tau_1 - \tau_2)$,

$$f(t) = \gamma^{1+2 \max(0, \lfloor (\alpha - t) / |\tau_2 - \tau_1| \rfloor + 1)} / (1 - \gamma^2), \text{ and} \quad (22)$$

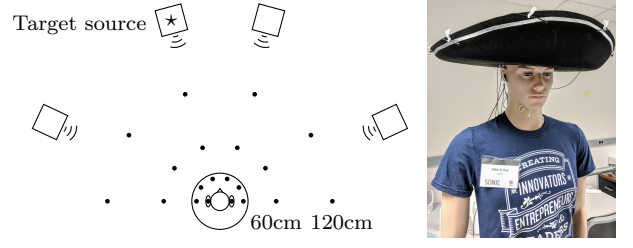


Fig. 3. Left: Recording setup. Circles are microphones and squares are loudspeakers. Right: Hat-mounted microphone array.

$$(c_1, c_2) = \begin{cases} (0, 0), & \text{if } |\tau_1| = |\tau_2| \\ (\sigma^2 + 1, \gamma + \gamma\sigma^2 - 1), & \text{if } |\tau_1| < |\tau_2| \\ (1, \sigma^2 + 1 - \gamma), & \text{if } |\tau_1| > |\tau_2|. \end{cases} \quad (23)$$

This delay-error curve is also piecewise constant, but has a geometric “tail” that decays with a rate of roughly $\gamma^{2/|\tau_2 - \tau_1|}$. The height of the steps is determined by σ^2 and the width is determined by $|\tau_2 - \tau_1|$, which depends on the distance between the sources. For large positive α , $\mathcal{E}(\alpha)$ approaches \mathcal{E}_{nc} .

Figure 2(a) shows $\mathcal{E}(\alpha)$ for four combinations of source placement. The causality penalty takes a different form depending on the relative placement of sources and microphones. For example, if both the target and interference source are closer to microphone 1 than microphone 2 (near/near), then the second microphone does not contribute any information at $\alpha = 0$. If the sources are on opposite sides, then the difference in TDOAs, $|\tau_1 - \tau_2|$ is larger, and therefore $\mathcal{E}(\alpha)$ decays more slowly.

2.4. Temporally correlated signals

The expressions above were derived for uncorrelated source and noise processes. In many applications, however, the signals of interest are correlated and can therefore be separated spectrally as well as spatially. It is difficult in general to predict the effects of signal correlation on the delay-error curve. However, if the entries of $\mathbf{R}_x(\omega)$ share a common spectral factor—for example, if the sources are identically distributed and are recorded by identical microphones—then we can write $\mathbf{R}_x(\omega) = H(\omega) \hat{\mathbf{G}}^H(\omega) \hat{\mathbf{G}}^H(\omega) H^*(\omega)$ and $\mathbf{R}_{xd}^H(\omega) = H(\omega) \hat{\mathbf{R}}_{xd}^H(\omega) H^*(\omega)$, where $H(\omega) H^*(\omega)$ is the scalar spectral factorization of the common factor. Then we have

$$\mathbf{R}_{xd}^H(\omega) (\mathbf{G}^H(\omega))^{-1} = H(\omega) \hat{\mathbf{R}}_{xd}^H(\omega) (\hat{\mathbf{G}}^H(\omega))^{-1} \quad (24)$$

$$\tilde{\mathbf{r}}^T(t) = (h * \hat{\mathbf{r}}^T)(t). \quad (25)$$

Since $h(t)$ is causal, it spreads the energy of $\tilde{\mathbf{r}}(t)$ forward in time. Figure 2(b) shows the same scenario as in the previous section, but with identically distributed speech-shaped sources. The error is lower overall, the steps are smoother, and the filter can begin to separate the signals even before they reach either microphone.

3. EXPERIMENTS

To evaluate delay-performance tradeoffs in realistic conditions, we recorded audio mixtures using a wearable microphone array in a cocktail party scenario at the Augmented Listening Laboratory at the University of Illinois at Urbana-Champaign, which has a reverberation time of around $T_{60} = 300$ ms. The recording setup, shown in Figure 3, consisted of twenty omnidirectional lavalier microphones:

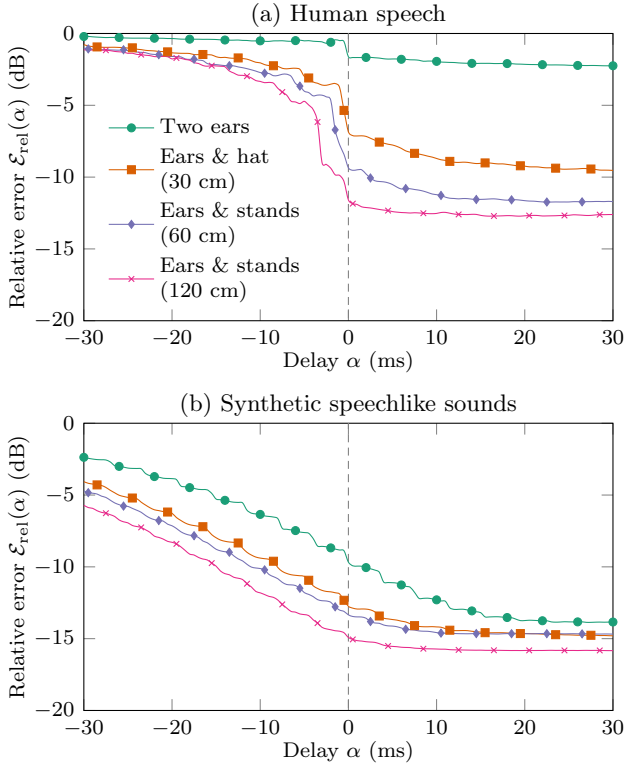


Fig. 4. Experimental delay-error results for isolating a single target source from a mixture of four sources.

two at the left and right ears of a mannequin “listener,” six along the perimeter of a hat with radius 30 cm, and twelve mounted on stands at 60 cm and 120 cm distances from the listener. The reference microphone is that in the left ear. Source signals were produced by loudspeakers two meters away from the listener. The acoustic impulse responses between the loudspeakers and microphones were measured using linear sweeps. All data was sampled at 16 kHz.

The signals were separated using the discrete-time, finite-length version of the CMWF. Let $\bar{\mathbf{x}}[k] = [\mathbf{x}^T[k], \dots, \mathbf{x}^T[k-L+1]]^T$ and $\bar{\mathbf{w}}_\alpha^T = [\mathbf{w}_\alpha^T[0], \dots, \mathbf{w}_\alpha^T[L-1]]$ be stacked vectors of the sampled multichannel signals and the finite impulse response filter coefficients, respectively. Let $y_\alpha[k] = \bar{\mathbf{w}}_\alpha^T \bar{\mathbf{x}}[k]$ be the filter output sequence and let $d_\alpha[k]$ be the desired output sequence. Let $\bar{\mathbf{r}}_x = \mathbb{E}[\bar{\mathbf{x}}[n]\bar{\mathbf{x}}^T[n]]$ and $\bar{\mathbf{r}}_{xd}(\alpha) = \mathbb{E}[\bar{\mathbf{x}}[n]d_\alpha[n]]$, where $\mathbb{E}[\cdot]$ is expectation. The linear minimum MSE filter coefficients are [10]

$$\bar{\mathbf{w}}_\alpha^T = \bar{\mathbf{r}}_{xd}^T(\alpha)\bar{\mathbf{r}}_x^{-1}. \quad (26)$$

In our experiments, $\bar{\mathbf{r}}_x$ was computed using truncated impulse response measurements. We applied diagonal loading comparable to the source power to account for modeling errors and ambient noise. We used discrete-time filters with length $L = 2048$ samples (128 ms). For each experiment we report the sample MSE relative to the source power, computed as $\mathcal{E}_{\text{rel}}(\alpha) = 10 \log_{10} \sum_k (y_\alpha[k] - d_\alpha[k])^2 / \sum_k d_\alpha^2[k]$.

Figure 4(a) shows delay-error curves for four simultaneous talkers using arrays of up to eight microphones at varying distances. The speech signals were twenty-second clips taken from the VCTK dataset [31] and the filters were designed using a single approximate long-term average speech autocorrelation. Because we model

the signals as identically distributed, the filters must rely on spatial rather than spectral diversity to separate the sources. As the radius of the array increases, the curves move downward and to the left, indicating that the larger-aperture arrays can achieve similar performance with lower delay compared to the smaller-aperture arrays. In fact, since the source signals reach the listener several milliseconds before they reach the microphone stands, the system could operate with negative delay.

The two-channel filter performs poorly in this experiment because it does not take advantage of the time-frequency sparsity of speech signals, which many speech enhancement algorithms exploit. To account for the benefits of sparsity within the stationary estimation framework of this paper, we repeated the experiment with four stationary speechlike sounds generated using the Vocaloid music synthesis software. Each ten-second source signal represents a different vowel sung in a different key. Although the signals are deterministic and periodic, the filters were designed based on 50 ms von Hann-windowed autocorrelation sequences. Figure 1 shows the delay-error curves for single-channel mixtures of these sources and Figure 4(b) compares the separation performance of multichannel filters with different array sizes. Because the sources are approximately disjoint in the frequency domain, a one- or two-channel filter can separate them effectively, but requires a delay to do so. The larger microphone arrays also benefit from longer delay, but perform better for small α . For example, the performance of the hat-mounted array with zero delay matches that of the binaural microphones with about 10 ms delay, which would be perceptible to many listeners.

4. CONCLUSIONS

The theoretical and experimental results presented here suggest that larger arrays can separate sound sources with lower delay and that the delay-performance tradeoff depends on both the spatial and temporal correlation structure of the observed signals. When microphones are located between the listener and sound sources, those sensors receive the signals before the listening device, shifting the delay-performance curve to the left. Arrays also provide spatial gain, which improves overall performance regardless of delay. When signals are spectrally distinct, a single-channel filter could separate them effectively given a long enough delay, but an array can achieve the same performance with little or no delay.

Much remains to be understood about delay-constrained array processing. For example, equations (10) and (11) tell us little in general about the effects of reverberation and signal spectra on delay. Furthermore, because many signals of interest are nonstationary, we must also consider time-varying causal array processing. Finally, to realize the benefits of spatial diversity in delay-constrained listening enhancement, listening devices must use larger microphone arrays than they do today. Large wearable and distributed arrays could allow us to apply stronger noise reduction while meeting the strict delay constraints of real-time listening applications.

5. REFERENCES

- [1] J. M. Kates, *Digital Hearing Aids*. Plural Publishing, 2008.
- [2] V. Valimaki, A. Franck, J. Ramo, H. Gamper, and L. Savioja, “Assisted listening using a headset: Enhancing audio perception in real, augmented, and virtual environments,” *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 92–99, 2015.
- [3] J. Agnew and J. M. Thornton, “Just noticeable and objection-

- able group delays in digital hearing aids,” *Journal of the American Academy of Audiology*, vol. 11, no. 6, pp. 330–336, 2000.
- [4] M. A. Stone and B. C. Moore, “Tolerable hearing aid delays. I. Estimation of limits imposed by the auditory path alone using simulated hearing losses,” *Ear and Hearing*, vol. 20, no. 3, pp. 182–192, 1999.
- [5] M. A. Stone, B. C. Moore, K. Meisenbacher, and R. P. Derleth, “Tolerable hearing aid delays. V. Estimation of limits for open canal fittings,” *Ear and Hearing*, vol. 29, no. 4, pp. 601–617, 2008.
- [6] M. A. Stone and B. C. Moore, “Tolerable hearing aid delays. II. Estimation of limits imposed during speech production,” *Ear and Hearing*, vol. 23, no. 4, pp. 325–338, 2002.
- [7] S. Makino, ed., *Audio Source Separation*. Springer, 2018.
- [8] O. Yilmaz and S. Rickard, “Blind separation of speech mixtures via time-frequency masking,” *IEEE Transactions on Signal Processing*, vol. 52, no. 7, pp. 1830–1847, 2004.
- [9] A. Ozerov and C. Févotte, “Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 3, pp. 550–563, 2010.
- [10] J. Benesty, J. Chen, and Y. Huang, *Microphone Array Signal Processing*, vol. 1. Springer, 2008.
- [11] M. Brandstein and D. Ward, *Microphone Arrays: Signal Processing Techniques and Applications*. Springer, 2013.
- [12] S. Gannot, E. Vincent, S. Markovich-Golan, and A. Ozerov, “A consolidated perspective on multimicrophone speech enhancement and source separation,” *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 4, pp. 692–730, 2017.
- [13] M. S. Pedersen, J. Larsen, U. Kjems, and L. C. Parra, “Convolutional blind source separation methods,” in *Springer Handbook of Speech Processing*, pp. 1065–1094, Springer, 2008.
- [14] S. Doclo, W. Kellermann, S. Makino, and S. E. Nordholm, “Multichannel signal enhancement algorithms for assisted listening devices,” *IEEE Signal Processing Magazine*, vol. 32, no. 2, pp. 18–30, 2015.
- [15] S. Doclo, S. Gannot, M. Moonen, and A. Spriet, “Acoustic beamforming for hearing aid applications,” in *Handbook on Array Processing and Sensor Networks* (S. Haykin and K. R. Liu, eds.), pp. 269–302, Wiley, 2008.
- [16] P. Naylor and N. D. Gaubitch, *Speech Dereverberation*. Springer, 2010.
- [17] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B.-H. Juang, “Blind speech dereverberation with multi-channel linear prediction based on short time Fourier transform representation,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 85–88, 2008.
- [18] B. Schwartz, S. Gannot, and E. Habets, “Online speech dereverberation using Kalman filter and EM algorithm,” *IEEE/ACM Transactions on Audio, Speech and Language Processing*, vol. 23, no. 2, pp. 394–406, 2015.
- [19] J. Benesty, J. Chen, Y. Huang, and J. Dmochowski, “On microphone-array beamforming from a MIMO acoustic signal processing perspective,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 3, pp. 1053–1065, 2007.
- [20] E. Hänsler and G. Schmidt, *Acoustic Echo and Noise Control: A Practical Approach*. Wiley, 2005.
- [21] J. Chen, J. Benesty, and Y. Huang, “A minimum distortion noise reduction algorithm with multiple microphones,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 3, pp. 481–493, 2008.
- [22] J. M. Kates and K. H. Arehart, “Multichannel dynamic-range compression using digital frequency warping,” *EURASIP Journal on Applied Signal Processing*, vol. 2005, pp. 3003–3014, 2005.
- [23] H. W. Löllmann and P. Vary, “Low delay noise reduction and dereverberation for hearing aids,” *EURASIP Journal on Advances in Signal Processing*, vol. 2009, no. 1, p. 437807, 2009.
- [24] T. Barker, T. Virtanen, and N. H. Pontoppidan, “Low-latency sound-source-separation using non-negative matrix factorisation with coupled analysis and synthesis dictionaries,” in *IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 241–245, 2015.
- [25] N. Wiener, *Extrapolation, Interpolation, and Smoothing of Stationary Time Series*. Wiley, 1949.
- [26] H. W. Bode and C. E. Shannon, “A simplified derivation of linear least square smoothing and prediction theory,” *Proceedings of the IRE*, vol. 38, no. 4, pp. 417–425, 1950.
- [27] H. L. Van Trees, *Detection, Estimation, and Modulation Theory, Part I*. Wiley, 2004.
- [28] N. Wiener and P. Masani, “The prediction theory of multivariate stochastic processes, II,” *Acta Mathematica*, vol. 99, no. 1, pp. 93–137, 1958.
- [29] E. Wong and J. Thomas, “On the multidimensional prediction and filtering problem and the factorization of spectral matrices,” *Journal of the Franklin Institute*, vol. 272, no. 2, pp. 87–99, 1961.
- [30] V. Kucera, “Factorization of rational spectral matrices: a survey of methods,” in *International Conference on Control*, pp. 1074–1078, 1991.
- [31] C. Veaux, J. Yamagishi, and K. MacDonald, “CSTR VCTK corpus: English multi-speaker corpus for CSTR voice cloning toolkit,” 2017.