# A COMPARATIVE STUDY ON TIME DELAY ESTIMATION IN REVERBERANT AND NOISY ENVIRONMENTS

*Jingdong Chen, Yiteng (Arden) Huang*

Bell Laboratories, Lucent Technologies
600 Mountain Avenue
Murray Hill, New Jersey 07974, USA
{jingdong, arden}@research.bell-labs.com

*Jacob Benesty*

Université du Québec, INRS-EMT
800 de la Gauchetière Ouest, Suite 6900
Montréal, Québec, H5A 1K6, Canada
benesty@inrs-emt.uquebec.ca

## ABSTRACT

Time delay estimation (TDE) has been a research topic of significant practical importance in many fields. It is the first stage that feeds into subsequent processing blocks of a multiple-channel system for identifying, localizing, and tracking radiating sources. Various TDE algorithms were developed in the past few decades, and their performance was assessed independently of each other when each algorithm was developed. This paper is to provide a comparative study to illustrate the performance differences among several representative TDE algorithms in room acoustic environments where reverberation, noise, and interference are commonly encountered.

## 1. PROBLEM FORMULATION

The goal of time delay estimation is to measure the relative time difference of arrival (TDOA) among signals received by spatially separated sensors. Two signal models have been widely adopted for describing the TDE problem, *i.e.*, the single-path propagation model and the reverberant model. The first one assumes that the signal acquired by each receiver is a delayed and attenuated version of the original source signal. Suppose that we have an array of $N$ receivers, the received signals are expressed as:

$$x_n[k] = \alpha_n s[k - t - f_n(\tau)] + w_n[k], \qquad (1)$$

where $\alpha_n$, $n = 0, 1, 2, \cdots, N - 1$, are the attenuation factors due to propagation effects, $t$ is the propagation time from the unknown source $s[k]$ to Sensor 0, $w_n[k]$ is an additive noise signal at the $n$th microphone, $\tau$ is the relative delay between Microphones 0 and 1, and $f_n(\tau)$ is the relative delay between Microphones 0 and $n$, with $f_0(\tau) = 0$ and $f_1(\tau) = \tau$. For $n = 2, \ldots, N - 1$, the function $f_n$ depends generally not only on $\tau$ but also on the microphone array geometry. For example, in the far-field case (plane wave propagation), for a linear and equispaced array, we have $f_n(\tau) = n\tau$, and for a linear but non-equispaced array, we have $f_n(\tau) = \left(\sum_{i=0}^{n-1} d_i/d_0\right)\tau$, where $d_i$ is the distance between Microphones $i$ and $i + 1$, $i = 0, 1, 2, \cdots, N - 2$. In the near-field case, $f_n$ depends also on the position of the source. Also note that $f_n(\tau)$ can be a *nonlinear* function of $\tau$ for a nonlinear array geometry, even in the far-field case (e.g., 3 equilateral sensors). In general $\tau$ is not known, but the geometry of the array is known such that the mathematical formulation of $f_n(\tau)$ is well defined or given. It is further assumed that $w_n[k]$ is a zero-mean Gaussian random process that is uncorrelated with $s[k]$ and the noise signals at other sensors. For this model, the TDE problem is formulated to determine an estimate $\hat{\tau}$ of the true time delay $\tau$ using a finite set of observation samples.

In many application scenarios such room acoustic environments, however, each sensor receives, in addition to the direct-path signal, multiple delayed and attenuated replicas of the source signal due to reflections of the wavefront from boundaries and objects. Taking into account this so-called multipath effect, a rever-

berant model was developed recently [1], where an FIR filter is used to model the channel between the source and each receiver. The received signals are expressed as

$$x_n[k] = \mathbf{h}_n^T \mathbf{s}[k] + w_n[k], \quad n = 0, 1, \ldots, N - 1, \qquad (2)$$

where

$$\mathbf{h}_n = \begin{bmatrix} h_{n,0} & h_{n,1} & \ldots & h_{n,L-1} \end{bmatrix}^T,$$
$$\mathbf{s}[k] = \begin{bmatrix} s[k] & s[k-1] & \ldots & s[k-L+1] \end{bmatrix}^T,$$

and $L$ is the length of the longest channel impulse responses among $N$ channels.

As seen, no time delay is explicitly expressed in (2), hence there is no plain solution to the TDE problem for the reverberant model, unless the channel impulse responses can be accurately (and blindly) identified, which is a very challenging problem.

## 2. TDE ALGORITHMS

Various TDE algorithms were developed in the literature. In this section, we brief some critical techniques.

### 2.1. Generalized Cross-Correlation Method

The generalized cross-correlation (GCC) method, which is developed by Knapp and Carter [2], is perhaps the most popular TDE algorithm thus far [2]. It does not only unify various correlation based algorithms into a general framework, but also provide a mechanism to incorporate knowledge to improve performance of TDE. In this framework, the delay estimate is obtained as

$$\hat{\tau}_{\text{GCC}} = \arg \max_m \hat{\Psi}_{\text{GCC}}[m], \qquad (3)$$

where

$$\hat{\Psi}_{\text{GCC}}[m] = \sum_{\omega=0}^{\Omega-1} \Phi[\omega] S_{x_0 x_1}[\omega] e^{\frac{j2\pi m\omega}{\Omega}}$$

is the generalized cross-correlation function (GCCF), $S_{x_0 x_1}[\omega] = E\{X_0[\omega]X_1^*[\omega]\}$ is the cross spectrum, $E\{\cdot\}$ and $(\cdot)^*$ stand respectively for the mathematical expectation and the complex conjugate operator, $X_n[\omega]$ is the discrete Fourier transform (DFT) of $x_n[k]$, $\Phi[\omega]$ is a weighting function (sometimes called a *prefilter*), and $\Omega$ is the length of DFT.

There are a number of member algorithms in the GCC family depending on how the weighting function $\Phi[\omega]$ is selected. Commonly used weighting functions include the constant weighting [in this case, the GCC becomes a frequency-domain implementation of the traditional cross-correlation (CC) method], the phase transform (PHAT), the maximum likelihood (ML) processor [2], etc. Different weighting functions possesses different properties, as for example the PHAT algorithm where $\Phi_{\text{PHAT}}[\omega] = 1/|S_{x_0 x_1}[\omega]|$. Substituting $\Phi_{\text{PHAT}}[\omega]$ into (3) and neglecting noise effects, one can readily deduce that the weighted cross spectrum is free from the source signal and depends only on the channel responses. Consequently the PHAT algorithm performs more consistently than

many other GCC members when the characteristics of the source signal change over time. It is also observed that the PHAT algorithm is more immune to reverberation than many other cross correlation based methods. Another example is the ML processor with which the delay estimate obtained in the single-path propagation situation is optimal from a statistical point of view since the estimation variance can achieve the Cramèr-Rao lower bound (CRLB). It should be pointed out that in order for the ML processor to achieve the optimal performance, spectra of noise signals have to be known *a priori*. In real applications, this information is not accessible, and can only be estimated. The ML algorithm then becomes suboptimal, like other GCC members.

### 2.2. LMS-Type Adaptive TDE Algorithm

This method, also based on the ideal propagation model with two sensors, was proposed by Reed *et al* in 1981 [3]. Different from the cross-correlation based approaches, this algorithm achieves time delay by minimizing the mean-square error between $x_0[k]$ and a filtered (FIR filter) version of $x_1[k]$, and the delay estimate is obtained as the lag time associated with the largest component of the FIR filter. If we define a signal vector of $x_1[k]$ at time instant $k$ as

$$\mathbf{x_1}[k] = [x_1[k - \tau_{\max}] \ x_1[k - \tau_{\max} + 1] \ \dots \ x_1[k]$$
$$x_1[k+1] \ \dots \ x_1[k + \tau_{\max}]]^T \quad (4)$$

and an FIR filter of length $2\tau_{\max} + 1$ as

$$\mathbf{h}[k] = [h_0 \ h_1 \ \dots \ h_l \ h_{l+1} \ \dots \ h_{2\tau_{\max}}]^T, \quad (5)$$

where again $\tau_{\max}$ is the maximum possible time delay, then an error signal can be formulated as

$$e[k] = x_0[k] - \mathbf{h}^T[k]\mathbf{x_1}[k]. \quad (6)$$

An estimate of $\mathbf{h}[k]$ can be achieved by minimizing $E\{e^2[k]\}$ using either a batch or an adaptive algorithm. For example, with the least-mean-square (LMS) adaptive algorithm, $\mathbf{h}[k]$ can be estimated through

$$\mathbf{h}[k+1] = \mathbf{h}[k] + \mu e[n]\mathbf{x_1}[k], \quad (7)$$

where $\mu$ is a small positive adaptation step size. Given this estimate of $\mathbf{h}[k]$, the delay estimate can be determined as

$$\tau_{\text{LMS}} = \arg\max_l |h_l| - \tau_{\max}. \quad (8)$$

Other adaptive algorithms [4] can also be used, which may lead to a better performance.

### 2.3. Fusion Algorithm Based on Multiple Sensor Pairs

The GCC framework, which can yield reasonable TDE performance in nonreverberant and moderate noisy environments, suffers significant performance degradation in the presence of reverberation. Much attention has been paid to improving the tolerance of TDE against noise and reverberation. Besides using some *a priori* knowledge about the distortion sources, another way of combating noise and reverberation is through exploiting the redundant information provided by multiple sensors. To illustrate the redundancy, let us consider a three-sensor linear array, which can be partitioned into three sensor pairs. Three delay measurements can then be acquired with the observation data, i.e., $\tau_{01}$ (TDOA between Sensor 0 and Sensor 1), $\tau_{12}$ (TDOA between Sensor 1 and Sensor 2), and $\tau_{02}$ (TDOA between Sensor 0 and Sensor 2). Apparently, these three delays are not independent. As a matter of fact, if the source is located in the far field, it is easily seen that $\tau_{02} = \tau_{01} + \tau_{12}$. Such a relation was exploited in [5] to formulate a two-stage TDE algorithm. In the preprocessing stage, three delay measurements were measured independently using the GCC method. A state equation was then formed and the Kalman filter

is used in the post-processing stage to enhance the delay estimate of $\tau_{01}$ and $\tau_{12}$. It was shown that in the far-filed case, the estimation variance of $\tau_{01}$ can be reduced by a factor of 6 in low SNR (SNR $\to$ 0), and of 4 in high SNR ( SNR $\to \infty$) conditions. More recently, several approaches based on multiple sensor pairs was developed to deal with the TDE in room acoustic environments [6], [7], [8]. Different from the Kalman filter method, these approaches fuse the estimation cost functions from multiple sensors pairs before searching the time delay. We shall call such a scheme as information fusion based algorithm. In general, the problem of TDE with the fusion algorithm can be formulated as

$$\hat{\tau}_{\text{FUSION}} = \arg\max_m \sum_{p=1}^{P} \mathcal{F}\{\hat{\Psi}_p[m]\}, \quad (9)$$

where $P$ is the total number of sensor pairs, $\hat{\Psi}_p[m]$ represents some delay cost function measured from the $p$th sensor pair (it can be CCF, GCCF, etc), and $\mathcal{F}\{\cdot\}$ denotes some mathematical transformation, which ensures that the cost functions ($\hat{\Psi}_p[m]$) for all the $P$ sensor pairs, after transformation, have their peaks due to the same source in the same location. Various methods can be formulated by selecting a different $\mathcal{F}\{\cdot\}$ or $\hat{\Psi}$. For example, if all sensor pairs are centered around a same position, by choosing $\mathcal{F}\{x\} = x$, $\hat{\Psi}[m]$ as the GCCF from the PHAT algorithm, one can readily derive the so-called synchronous adding method in [6]. We can also easily derive the consistency method in [7] and the SRP (steered response power)-PHAT algorithm in [8].

### 2.4. Multichannel Cross-Correlation Algorithm

Recently, a squared multichannel cross-correlation coefficient (MCCC) was derived from the theory of spatial interpolation [9]. Consider the signal model given in (1) with a total of $N$ sensors. At time instant $k$, the MCCC is defined as :

$$\varrho_N^2(k,m) = 1 - \det\left[\widetilde{\mathbf{R}}(k,m)\right], \quad (10)$$

where "det" stands for *determinant* of a matrix,

$$\widetilde{\mathbf{R}}(k,m) = \begin{bmatrix} 1 & \rho_{01}(k,m) & \cdots & \rho_{0N-1}(k,m) \\ \rho_{10}(k,m) & 1 & \cdots & \rho_{1N-1}(k,m) \\ \vdots & \ddots & \ddots & \vdots \\ \rho_{N-10}(k,m) & \rho_{N-11}(k,m) & \cdots & 1 \end{bmatrix}, \quad (11)$$

is the normalized signal covariance matrix,

$$\rho_{ij}(k,m) = \frac{r_{ij}(k,m)}{\sqrt{r_{ii}(k,m)r_{jj}(k,m)}} \quad (12)$$

is the cross-correlation coefficient between $x_i$ and $x_j$,

$$r_{ij}(k,m) = \sum_{p=0}^{k} \lambda^{k-p} x_i[p + f_j(m)]x_j[p + f_i(m)]\}, \quad (13)$$

and $i, j = 0, 1, \cdots, N - 1$.

Just like the cross-correlation coefficient between two signals, this definition of cross correlation among multiple channels possesses quite a few good properties, and can be treated as a natural generalization of the traditional cross-correlation coefficient to the multichannel case. The problem of TDE at time instant $k$, based on this new definition, can be formulated as

$$\hat{\tau}_{\text{MCCC}} = \arg\max_m \varrho_N^2(k,m)$$
$$= \arg\min_m \left\{ \det\left[\widetilde{\mathbf{R}}(m,k)\right] \right\}. \quad (14)$$

For two-sensor case, it can be easily checked that this method is same as the cross-correlation method. When we have more than

2005 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics

two sensors, this method can be viewed as a natural generalization of the cross-correlation method to the multichannel case, which can take advantage of the redundancy among multiple sensors to improve the time delay estimate between two sensors. It is worth to mention that a prewhitening process can be applied to the observation signals before delay estimation. In this case, the MCCC algorithm can be treated as a generalized version of the PHAT algorithm.

### 2.5. Adaptive Eigenvalue Decomposition Algorithm

All the algorithms outlined in the previous sections are derived from the single-path propagation model. Recently, an adaptive eigenvalue decomposition (AED) algorithm was proposed to deal with TDE in room reverberant environment [1]. This algorithm first identifies the channel impulse responses from the source to the two sensors ($\mathbf{h}_0$ and $\mathbf{h}_1$). The delay estimate is then determined by finding the direct paths from the two measured impulse responses. In brief, at time instant $k$, the channel impulse response vector $\mathbf{u} = [\ \mathbf{h}_0^T \ -\mathbf{h}_1^T \ ]^T$ is estimated as the eigenvector of the covariance matrix $\mathbf{R}[k]$ associated with the smallest eigenvalue, where $\mathbf{R}[k] = E\{\mathbf{x}[k]\mathbf{x}^T[k]\}$. In an adaptive way, $\mathbf{u}$ can be estimated via

$$\hat{\mathbf{u}}[k+1] = \frac{\hat{\mathbf{u}}[k] - \mu e[k]\mathbf{x}[k]}{\|\hat{\mathbf{u}}[k] - \mu e[k]\mathbf{x}[k]\|}, \tag{15}$$

with the constraint that $\|\hat{\mathbf{u}}[k]\| = 1$, where

$$e[k] = \hat{\mathbf{u}}^T[k]\mathbf{x}[k] \tag{16}$$

is an error signal, $\|\cdot\|$ denotes the $l_2$ norm of a vector or matrix, and $\mu$, the adaptation step, is a positive constant.

With the identified impulse responses $\hat{\mathbf{h}}_0$ and $\hat{\mathbf{h}}_1$, the time delay estimate is determined as the difference between two direct paths, i.e.,

$$\hat{\tau}_{\text{AED}} = \min\left\{\underset{l}{\text{argmax}}^q |h_{1,l}|\right\} - \min\left\{\underset{l}{\text{argmax}}^q |h_{0,l}|\right\} \tag{17}$$

where $\max^q$ computes the $q^{th}$ largest element.

### 2.6. Adaptive Multichannel Time Delay Estimation

In the AED algorithm, the delay estimate is obtained by blindly identifying two channel impulse responses. It requires that the two channels do not share any common zeros, which is usually true for systems with short impulse responses. In many application scenarios such as room acoustic environments, however, the channel impulse response from the source to the microphone sensor could be very long. As a result, the likelihood for two impulse responses not sharing common zeros tends to be low and the AED algorithm often fails when a zero is shared between two channels or some zeros of the two channels are close. One way to overcome this problem is to employ more channels in the system, since it would be less likely for all channels to share a common zero when the number of sensors is large. This idea leads to an adaptive multichannel (AMC) time delay estimation approach based on a blind channel identification technique [10].

Considering the reverberant model in (2), we can define a cost function among all the $N$ channels, at time instant $k+1$, as

$$J[k+1] = \sum_{i=0}^{N-2} \sum_{j=i+1}^{N-1} e_{ij}^2[k+1], \tag{18}$$

where

$$e_{ij}[k+1] = \frac{\mathbf{x}_i^T[k+1]\hat{\mathbf{h}}_j[k] - \mathbf{x}_j^T[k+1]\hat{\mathbf{h}}_i[k]}{\left\|\hat{\mathbf{h}}[k]\right\|}$$
$$i, j = 0, 1, \cdots, N-1, \tag{19}$$

is an error signal between Sensor $i$ and Sensor $j$ at time $k+1$, $\hat{\mathbf{h}}_n[k]$ is the modeling filter of $\mathbf{h}_n$, and

$$\hat{\mathbf{h}}[k] = \left[\ \hat{\mathbf{h}}_0^T[k] \quad \hat{\mathbf{h}}_1^T[k] \quad \cdots \hat{\mathbf{h}}_{N-1}^T[k] \ \right]^T. \tag{20}$$

It follows immediately that various adaptive algorithms can be used to estimate the channel impulse responses. For example, a multichannel LMS (MCLMS) algorithm and a normalized multichannel frequency-domain LMS (NMCFLMS) algorithm were developed [10] to estimate $\hat{\mathbf{h}}$ by minimizing $J[k+1]$. While the former performs estimation in the time domain, the latter operates in the frequency domain on a block-by-block basis, which enables a faster convergence rate. We will adopt the NMCFLMS algorithm in our experiment. Once $\hat{\mathbf{h}}$ is achieved, time delay between the $i$th and $j$th sensors is determined as

$$\hat{\tau}_{ij} = \min\left\{\underset{l}{\text{argmax}}^q |h_{i,l}|\right\} - \min\left\{\underset{l}{\text{argmax}}^q |h_{j,l}|\right\}. \tag{21}$$

### 3. EXPERIMENTS

In an attempt to simulate real reverberant acoustic environments, the image model technology [11] is used. We consider a rectangular room of size $120 \times 180 \times 150$ inches ($x \times y \times z$). A point omnidirectional source is located at (100, 100, 40). A linear array which consists of four (4) ideal point microphones is placed in parallel with the x-axis. Four microphones are located at (20, 10, 40), (28, 10, 40), (36, 10, 40), and (44, 10, 40), respectively. The directivity pattern of each microphone is assumed to be omnidirectional.

A low-pass sampled version of the impulse response of the acoustic transmission channel between the source and each microphone is generated using the image method. A speech signal from a female speaker, digitized with 16-bit resolution at 16 kHz, is then convolved with the synthetic impulse responses. Finally, mutually independent white Gaussian noise is properly scaled and added to each microphone signal to control the SNR.

Delay estimates were obtained on a frame-by-frame basis. The frame size used in all experiments is 64 ms. To reduce the temporal effect of noise on TDE performance, the cost function of each algorithm is smoothed using a single-pole recursion as follows:

$$\bar{\Psi}_k = \gamma\bar{\Psi}_{k-1} + (1-\gamma)\hat{\Psi}_k, \tag{22}$$

where $\hat{\Psi}_k$ denotes the cost function estimated using the $k$th frame of observation data, $\bar{\Psi}_k$ is a smoothed version of the cost function, based on which the delay estimates were obtained. For the MCCC algorithm, the signal was prewhitened before computing the cost function. Therefore, this method, in the case of two sensors, is equivalent to the PHAT algorithm. For the ML method, we assume that the noise spectrum is know *a priori*. The fusion algorithm implemented here is the consistency method presented in [7]).

It is not always easy to compare fairly different algorithms. In our experiments, we optimized each individual algorithm in a non-reverberant and favorable noisy ($\text{SNR} = 25$ dB) environment to its best performance. We then test and compare all the algorithms in reverberation and different noise conditions. Such a process should, in generally, not favor any specific algorithm.

Several experiments were performed. Due to space limitations, we present one set of results, as shown in Fig. 1, where $\text{SNR} = 15$ dB. It can be seen that, in the first environment, all the algorithms can accurately identify the time delay. When reverberation time is increased to 580 ms, both the CC and the ML methods suffer significant performance degradation, showing that

these two approaches are sensitive to reverberation. The PHAT algorithm, though also belongs to the GCC family like the CC and ML methods, still yields a reasonable performance, implying its robustness with respect to reverberation. Among the five techniques that use two sensors (i.e., CC, PHAT, ML, AED, LMS), the AED algorithm delivers the best performance. This indicates that taking it into account in the signal model is an effective way in dealing with reverberation. Comparing the MCCC, AMC, and fusion algorithms with dual-sensor techniques, one can easily see the advantage of using multiple sensors. Since the AMC algorithm was formulated from the reverberant signal model and using multiple sensors, it is not surprising to see that it achieves the best performance in this strong reverberant environment.

## 4. SUMMARY

This paper presented a comparative study of TDE techniques in adverse environments. Broadly, the studied techniques can be classified into two categories: cross-correlation based methods and system identification based approaches. Both categories can be implemented either based on two sensors, or using multiple sensors. We evaluated eight algorithms, including five dual-channel techniques and three multiple-channel techniques, in both reverberant and noisy environments. Among the five studied dual-channel techniques, the adaptive eigenvalue decomposition algorithm demonstrated the best performance in both noise and reverberation conditions, showing its great potential for real applications. In general, more sensors will lead to a higher robustness because of the redundancy. However, it should be pointed out that attention has to be paid to implementing the multichannel cross-correlation algorithm and the fusion method. Both need to synchronize either the signals observed at different sensors, or the cost functions from different sensor pairs. In case that the true delay is not integral multiple of the sampling rate, we will have to either increase the sampling rate or use interpolation, which may significantly increase the computational complexity. In case that the observation signals or the cost functions are not properly aligned, we may not achieve much improvement.

## REFERENCES

[1] J. Benesty, "Adaptive eigenvalue decomposition algorithm for passive acoustic source localization," *J. Acoust. Soc. Am.*, vol. 107, pp. 384–391, Jan. 2000.

[2] C. H. Knapp and G. C. Carter, "The generalized correlation method for estimation of time delay," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-24, pp. 320–327, Aug. 1976.

[3] F. A. Reed, P. L. Feintuch, and N. J. Bershad, "Time delay estimation using the LMS adaptive filter–static behavior," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. ASSP-29, pp. 561–571, June 1981.

[4] S. Haykin, *Adaptive Filter Theory*. Fourth Edition, Upper Saddle River, NJ: Prentice Hall, 2002.

[5] R. L. Kirlin, D. F. Moore, and R. F. Kubichek, "Improvement of delay measurements from sonar arrays via sequential state estimation," *IEEE Trans. Acoust., Speech, Signal Processing*, vol. 29, pp. 514–519, June 1981.

[6] T. Nishiura, T. Yamada, S. Nakamura, and K. Shikano, "Localization of multiple sound sources based on a CSP analysis with a microphone array," *Proc. IEEE ICASP*, 2000, pp. 1053–1055.

[7] S. M. Griebel and M. S. Brandstein, "Microphone array source localization using realizable delay vectors," in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2001, pp. 71–74.

[8] J. DiBiase, H. Silverman, and M. Branstein, "Robust localization in reverberant rooms," in *Microphone Arrays: Signal Processing Techniques and Applications*, M. Branstein and D. Ward, Eds., Springer, New York, 2001.

[9] J. Chen, J. Benesty, and Y. Huang, "Robust time delay estimation exploiting redundancy among multiple microphoens," *IEEE Trans. Speech Audio Processing*, vol. 11, pp. 549–557, Nov. 2003.

[10] Y. Huang and J. Benesty, "Adaptive multichannel time delay estimation based on blind system identification for acoustic source localization," in *Adaptive Signal Processing–Applications to Real-World Problems*, J. Benesty and Y. Huang, Eds., Springer, New York, 2003.

[11] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *J. Acoust. Soc. Am.*, vol. 65, no. 4, pp. 943–950, Apr. 1979.
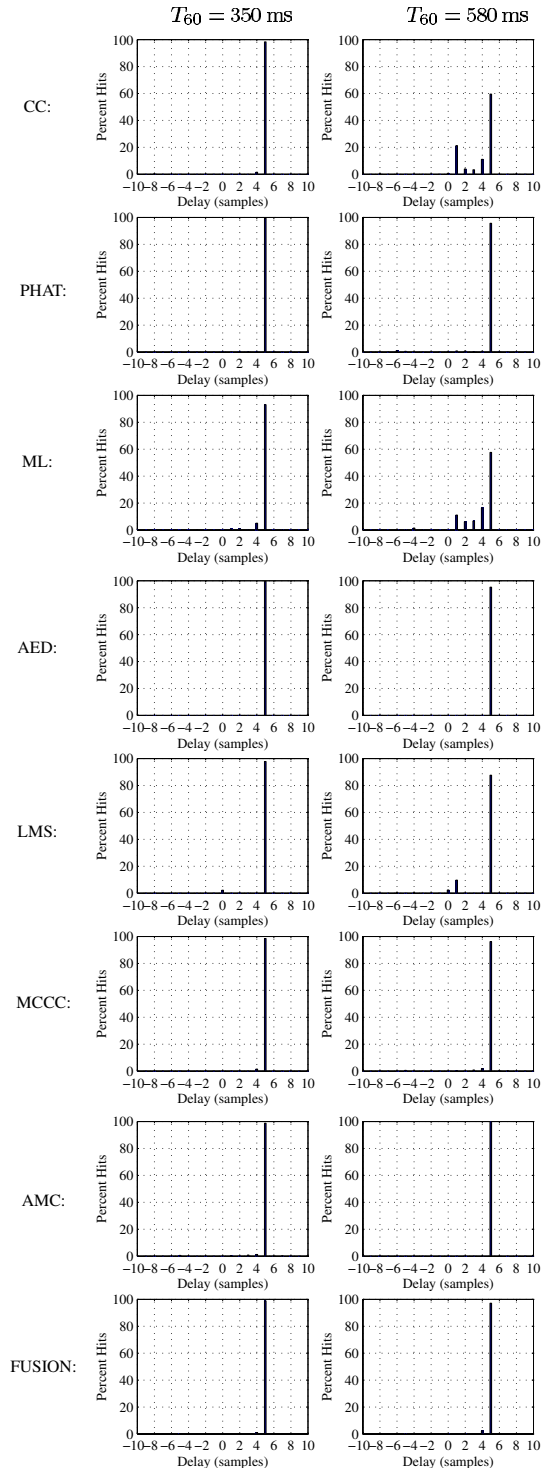
**Fig. 1**. TDE performances in moderate noisy and reverberant environments, where $\mathrm{SNR} = 15\,\mathrm{dB}$, and $T_{60} = 350\,\mathrm{ms}$, and $580\,\mathrm{ms}$ respectively.