

# NOISE ROBUST SPEECH DEREVERBERATION WITH KALMAN SMOOTHER

Masahito Togami and Yohei Kawaguchi

Central Research Laboratory, Hitachi Ltd.  
1-280, Higashi-koigakubo Kokubunji-shi, Tokyo 185-8601, Japan

## ABSTRACT

A speech dereverberation method is proposed that is robust against background noise. In contrast to conventional methods based on the linear prediction of the given microphone input signal, in which the linear prediction coefficients are not fully optimized when there is background noise, the proposed method optimizes the coefficients by linear prediction of the noiseless reverberant speech signal even when there is background noise. The noiseless reverberant speech signal and the parameters are iteratively updated on the basis of the expectation maximization algorithm. In the expectation step, sufficient statistics of latent variables which include noiseless reverberant speech signal are estimated using the Kalman smoother. Unlike the standard Kalman smoother, which uses a time-invariant covariance matrix as a state-transition covariance matrix, the proposed method utilizes a time-varying covariance matrix, enabling it to meet the time-varying speech characteristics. The parameters are updated so that the Q function is increased in the maximization step. Experimental results show that the proposed method is superior to conventional methods under noisy conditions.

**Index Terms**— Noise reduction, dereverberation, kalman smoother

## 1. INTRODUCTION

Room reverberation and background noise are major problems for speech recording and speech communication systems. Their levels must therefore be reduced simultaneously.

Several dereverberation techniques using multichannel inverse filtering have been developed [1][2][3][4][5]. A commonly used dereverberation technique is multi-step linear prediction (MSLP) [3], which reduces late reverberation by linear prediction of the microphone input signal. This technique is easy to implement and reduces reverberation effectively when there is no background noise. If the microphone input signal is contaminated by background noise, the linear prediction coefficients are updated so as to minimize the summation of the residual late reverberation and residual background noise after linear prediction. Therefore, if the background noise is dominant, dereverberation performance is degraded.

We have developed a speech dereverberation technique based on linear prediction of noiseless reverberant speech signal that overcomes this problem. The linear prediction coefficients are optimized so as to minimize only the amount of the residual reverberation after linear prediction even when the background noise is dominant. Since the noiseless reverberant speech signal is not given in advance, the proposed method estimates the noiseless reverberant speech signal and dereverberation parameters in an iterative manner using the expectation maximization (EM) algorithm [6]. In the expectation (E) step, the Kalman smoother [7] is used to obtain sufficient statistics

of latent variables which include the noiseless reverberant speech signal. The state-transition noise is related to the direct path of the speech signal in the state-transition equation. By extracting the state-transition noise, we can obtain the noiseless and dereverberated speech signal. A time-varying covariance matrix is used as a state-transition covariance matrix to meet the time-varying characteristics of the speech sources. The dereverberation parameters are updated so as to increase the Q function in the maximization (M) step. Experimental results show that the proposed method can reduce reverberation and background noise effectively in noisy environments.

## 2. PROBLEM STATEMENT

### 2.1. Microphone input signal model

The focus here is on multichannel processing. Assuming that there is a single speech source and background noise, we can model the multichannel microphone input signal at frequency  $f$  and for frame  $\tau$  as

$$\mathbf{x}_{f,\tau} = \sum_{l=0}^{L_1-1} s_{f,\tau-l} \mathbf{a}_{f,l} + \mathbf{w}_{f,\tau}, \quad (1)$$

where  $\mathbf{x}_{f,\tau} = [x_{1,f,\tau} \ \dots \ x_{N_m,f,\tau}]^T$ ,  $T$  is the transpose operator of a matrix or vector,  $N_m$  is the number of microphones,  $L_1$  is the tap length of the acoustic transfer function (ATF),  $s_{f,\tau}$  is the source signal at each time-frequency point,  $\mathbf{a}_{f,l}$  is the  $l$ th tap of the ATF of the speech source at each time-frequency point, and  $\mathbf{w}_{f,\tau}$  is the background noise term (white Gaussian noise).

### 2.2. Conventional autoregressive model for microphone input signal

In the conventional methods [3], (1) is transformed into an autoregressive model of the microphone input signal:

$$\mathbf{x}_{f,\tau} = \sum_{l=0}^{D-1} s_{f,\tau-l} \mathbf{a}_{f,l} + \sum_{l=D}^{L_2-1} \mathbf{W}_{f,l} \mathbf{x}_{f,\tau-l} + \mathbf{b}_{f,\tau}, \quad (2)$$

where  $L_2 = L_1 + L_i - 1$ ,  $L_i$  is the length of the inverse filter of the ATF,  $D$  is the tap-length of the early reflection, and  $\mathbf{b}_{f,\tau}$  is the convolutive background noise term with the ATF of the speech source and the inverse filter.

In the conventional time-varying source-model-based dereverberation techniques, such as that of Nakatani et al. [9], the linear prediction coefficients  $\mathbf{W}_{f,l}$  are estimated using

$$\hat{\mathbf{W}}_{f,l} = \underset{\mathbf{W}_{f,l}}{\operatorname{argmin}} \sum_{\tau=1}^{L_T} \mathcal{G}(\mathbf{W}_{f,l}, \mathbf{x}_{f,\tau}), \quad (3)$$

$$\mathcal{G}(\mathbf{W}_{f,l}, \mathbf{x}_{f,\tau}) = (\mathbf{x}_{f,\tau} - \sum_{l=D}^{L_2-1} \mathbf{W}_{f,l} \mathbf{x}_{f,\tau-l})^H \mathbf{R}_{x,f,\tau}^{-1} (\mathbf{x}_{f,\tau} - \sum_{l=D}^{L_2-1} \mathbf{W}_{f,l} \mathbf{x}_{f,\tau-l}), \quad (4)$$

where  $L_T$  is the number of the frames used for parameter optimization,  $H$  is the Hermite transpose operator of a matrix or vector, and  $\mathbf{R}_{x,f,\tau} = E[\mathbf{x}_{f,\tau} \mathbf{x}_{f,\tau}^H]$ . Let  $\mathbf{x}_{f,\tau}$  be divided into speech term  $\mathbf{c}_{f,\tau} = \sum_{l=0}^{L_1-1} s_{f,\tau-l} \mathbf{a}_{f,\tau,l}$  and background noise term  $\mathbf{w}_{f,\tau}$ :

$$\mathbf{x}_{f,\tau} = \mathbf{c}_{f,\tau} + \mathbf{w}_{f,\tau}. \quad (5)$$

Assuming that  $\mathbf{c}_{f,\tau}$  and  $\mathbf{w}_{f,\tau}$  are uncorrelated, we can approximate  $\sum_{\tau=1}^{L_T} \mathcal{G}(\mathbf{W}_{f,l}, \mathbf{x}_{f,\tau})$ :

$$\sum_{\tau=1}^{L_T} \mathcal{G}(\mathbf{W}_{f,l}, \mathbf{x}_{f,\tau}) \approx \sum_{\tau=1}^{L_T} \mathcal{G}(\mathbf{W}_{f,l}, \mathbf{c}_{f,\tau}) + \sum_{\tau=1}^{L_T} \mathcal{G}(\mathbf{W}_{f,l}, \mathbf{w}_{f,\tau}). \quad (6)$$

The first term is the amount of residual reverberation after linear prediction by  $\mathbf{W}_{f,l}$ , and the second term is the amount of residual background noise. The conventional method thus works by updating the linear prediction coefficients so as to minimize the summation of the residual reverberation and the residual background noise after linear prediction. If the background noise is small, the first term is dominant. Therefore, linear prediction coefficients that can reduce reverberation accurately should be obtained. However, if the background noise is dominant, dereverberation performance is degraded because the linear prediction coefficients are updated so as to minimize the residual background noise.

### 3. PROPOSED METHOD

#### 3.1. Linear prediction with latent multichannel source signal

Our proposed speech dereverberation method avoids performance degradation due to background noise by minimizing only the residual reverberation after linear prediction instead of minimizing the summation of the residual reverberation and residual background noise. Since the noiseless speech signal cannot be obtained in advance, sufficient statistics of latent variables which include the noiseless speech signal are estimated by using the Kalman smoother [7]. The linear prediction coefficients are obtained by using the estimated statistics as follows:

$$\hat{\mathbf{W}}_{f,l} = \underset{\mathbf{W}_{f,l}}{\operatorname{argmin}} \sum_{\tau=1}^{L_T} E[\mathcal{G}(\mathbf{W}_{f,l}, \mathbf{c}_{f,\tau})], \quad (7)$$

where  $E$  is the operator for mathematical expectation. A block diagram of the proposed method is shown in Fig. 1. The Kalman smoother corresponds to the E step in the EM algorithm [6]. In the M step, the dereverberation parameters are updated using the statistics obtained in the E step so as to increase the Q function. The noiseless speech signal and dereverberation parameters are thereby updated in an iterative manner.

#### 3.2. Sufficient statistics estimation

Sufficient statistics of latent variables which include the noiseless speech signal are the minimum mean square error (MMSE) estimate and the mean square error (MSE) of the latent variables. These statistics are effectively calculated by the Kalman smoother. At first, the microphone input signal is converted into a state-transition model and an observation model.

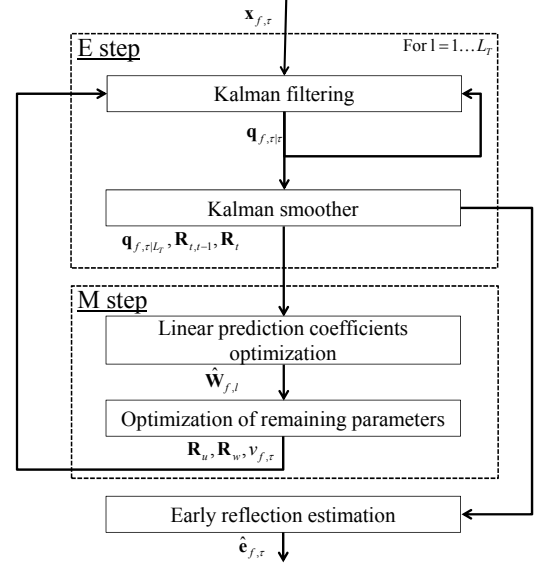


Fig. 1. Block diagram of proposed method

#### 3.2.1. State-transition model

The noiseless reverberant speech signal in the microphone input signal can be expressed by the following state-transition model as

$$\mathbf{q}_{f,\tau} = \mathbf{A}_f \mathbf{q}_{f,\tau-1} + \mathbf{u}_{f,\tau}, \quad (8)$$

where  $\mathbf{A}_f$  is the state-transition matrix,  $\mathbf{q}_{f,\tau}$  is a state vector,  $\mathbf{q}_{f,\tau} = [\mathbf{c}_{f,\tau}^H, \mathbf{c}_{f,\tau-1}^H, \dots, \mathbf{c}_{f,\tau-L_2+2}^H]^H$ ,  $\mathbf{A}_f$  is defined as

$$\mathbf{A}_f = \begin{bmatrix} \mathbf{0}_{N_m \times N_m(D-1)} & \mathbf{W}_{f,D} & \dots & \mathbf{W}_{f,L_2-1} \\ \mathbf{I}_{N_m \times N_m} & \mathbf{0} & \dots & \mathbf{0} \\ \mathbf{0}_{N_m \times N_m} & \mathbf{I}_{N_m \times N_m} & \dots & \mathbf{0} \\ \mathbf{0}_{N_m \times N_m} & \mathbf{0}_{N_m \times N_m} & \mathbf{I}_{N_m \times N_m} & \mathbf{0} \\ \vdots & \vdots & \vdots & \vdots \end{bmatrix}, \quad (9)$$

and  $\mathbf{u}_{f,\tau}$  is defined as

$$\mathbf{u}_{f,\tau} = [\mathbf{e}_{f,\tau}^T, \mathbf{0}, \mathbf{0}]^T, \quad (10)$$

where  $\mathbf{e}_{f,\tau}$  is an  $M \times M$  matrix. This matrix is the summation of the direct path and the early reflection term:

$$\mathbf{e}_{f,\tau} = \sum_{l=0}^{D-1} s_{f,\tau-l} \mathbf{a}_{f,l}. \quad (11)$$

The speech sources can be modeled as non-stationary source signals that are located at the same location while talking. Under these conditions, and similar to the modeling used in the conventional method [5], the probabilistic distribution of  $\mathbf{e}_{f,\tau}$  is modeled as a multichannel time-varying Gaussian distribution [10] with a 0-mean vector:

$$p(\mathbf{e}_{f,\tau}) = \mathcal{N}(\mathbf{e}_{f,\tau}; \mathbf{0}, v_{f,\tau} \mathbf{R}_{u,f}), \quad (12)$$

where  $v_{f,\tau}$  is a time-varying scalar coefficient which reflects the time-varying characteristics of the speech source and  $\mathbf{R}_{u,f}$  is a time-invariant matrix which reflects the time-invariant characteristics of the speech source location.

### 3.2.2. Observation model

The microphone input signal which consists of the noiseless reverberant speech and the background noise can be derived as the following observation model:

$$\mathbf{x}_{f,\tau} = \mathbf{F}\mathbf{q}_{f,\tau} + \mathbf{w}_{f,\tau}, \quad (13)$$

where

$$\mathbf{F} = \begin{bmatrix} \mathbf{I}_{N_m \times N_m} & \mathbf{0} \\ \mathbf{0} & \mathbf{0} \end{bmatrix}. \quad (14)$$

The probabilistic distribution of observation noise  $\mathbf{w}_{f,\tau}$  is approximately modeled as a stationary Gaussian distribution with a  $\mathbf{0}$ -mean vector as follows:

$$p(\mathbf{w}_{f,\tau}) = \mathcal{N}(\mathbf{w}_{f,\tau}; \mathbf{0}, \mathbf{R}_{w,f}), \quad (15)$$

where  $\mathbf{R}_{w,f}$  is the time-invariant covariance matrix of  $\mathbf{w}_{f,\tau}$ .

### 3.2.3. Kalman smoother for sufficient statistics estimation

The MMSE estimate and the MSE of the state vector  $\mathbf{q}_{f,\tau}$  are observed under the condition that the microphone input signal  $\mathbf{x}_{f,\tau}$  ( $\tau = 1 \dots L_T$ ) is given.

**Kalman filter:**

$$\mathbf{q}_{f,1|0} = \boldsymbol{\pi}_f, \quad (16)$$

$$\mathbf{R}_{f,1|0} = \mathbf{z}_f, \quad (17)$$

$$\mathbf{q}_{f,\tau|\tau-1} = \mathbf{A}_f \mathbf{q}_{f,\tau-1|\tau-1}, \quad (18)$$

$$\mathbf{q}_{f,\tau|\tau} = \mathbf{q}_{f,\tau|\tau-1} + \mathbf{K}_{f,\tau}(\mathbf{x}_{f,\tau} - \mathbf{F}\mathbf{q}_{f,\tau|\tau-1}), \quad (19)$$

$$\mathbf{K}_{f,\tau} = \mathbf{R}_{f,\tau|\tau-1} \mathbf{F}^H \mathbf{R}_{x,f,\tau}^{-1}, \quad (20)$$

$$\mathbf{R}_{f,\tau|\tau-1} = \mathbf{R}_{f,\tau-1|\tau-1} + v_{f,\tau} \mathbf{R}_{u,f}, \quad (21)$$

$$\mathbf{R}_{x,f,\tau} = \mathbf{R}_{w,f} + \mathbf{F} \mathbf{R}_{f,\tau|\tau-1} \mathbf{F}^H, \quad (22)$$

$$\mathbf{R}_{f,\tau|\tau} = (\mathbf{I} - \mathbf{K}_{f,\tau} \mathbf{F})(\mathbf{R}_{f,\tau|\tau-1}), \quad (23)$$

where  $\boldsymbol{\pi}_f$  is the initial MMSE estimate of the latent variable,  $\mathbf{z}_f$  is the initial MSE,  $\mathbf{q}_{f,\tau_1|\tau_2}$  is the MMSE estimate of  $\mathbf{q}_{f,\tau_1}$  under the condition that  $\mathbf{x}_{f,\tau}$  ( $\tau = 1 \dots \tau_2$ ) are observed, and  $\mathbf{R}_{f,\tau}$  is the MSE matrix.

**Kalman smoother:**

$$\begin{aligned} \mathbf{R}_{f,\tau|L_T} &= \mathbf{R}_{f,\tau|\tau} \\ &- \mathbf{B}_{f,\tau}(\mathbf{R}_{f,\tau+1|\tau} - \mathbf{R}_{f,\tau+1|L_T})\mathbf{B}_{f,\tau}^H, \end{aligned} \quad (24)$$

$$\mathbf{q}_{f,\tau|L_T} = \mathbf{q}_{f,\tau|\tau} + \mathbf{B}_{f,\tau}(\mathbf{q}_{f,\tau+1|L_T} - \mathbf{q}_{f,\tau+1|\tau}), \quad (25)$$

$$\mathbf{P}_{f,\tau} = \mathbf{R}_{f,\tau|L_T} + \mathbf{q}_{f,\tau|L_T} \mathbf{q}_{f,\tau|L_T}^H, \quad (26)$$

$$\mathbf{P}_{f,\tau,\tau-1} = \mathbf{R}_{f,\tau,\tau-1|L_T} + \mathbf{q}_{f,\tau|L_T} \mathbf{q}_{f,\tau-1|L_T}^H, \quad (27)$$

$$\begin{aligned} \mathbf{R}_{f,\tau,\tau-1|L_T} &= \mathbf{R}_{f,\tau|\tau} \mathbf{B}_{f,\tau-1}^H \\ &- \mathbf{B}_{f,\tau}(\mathbf{A}_f \mathbf{R}_{f,\tau|\tau} - \mathbf{R}_{f,\tau+1,\tau|L_T})\mathbf{B}_{f,\tau-1}^H, \end{aligned} \quad (28)$$

$$\mathbf{R}_{f,L_T,L_T-1|L_T} = (\mathbf{I} - \mathbf{K}_{f,L_T} \mathbf{F}) \mathbf{A}_f \mathbf{R}_{f,L_T-1|L_T-1}, \quad (29)$$

$$\mathbf{B}_{f,\tau} = \mathbf{R}_{f,\tau|L_T} \mathbf{A}_f^H \mathbf{R}_{f,\tau+1|L_T}^{-1}, \quad (30)$$

The dereverberated noiseless speech signal is estimated as the following state-transition noise:

$$\hat{\mathbf{u}}_{f,\tau} = \mathbf{q}_{f,\tau|L_T} - \mathbf{A}_f \mathbf{q}_{f,\tau-1|L_T}, \quad (31)$$

where the first  $N_m$ th row of  $\hat{\mathbf{u}}_{f,\tau}$  is the dereverberated speech source estimate.

### 3.3. Parameter optimization

In the M step, the parameters are updated by using the sufficient statistics obtained using the Kalman smoother. The optimization algorithm for the state-space model with the time-varying state-transition covariance matrix is slightly modified by the parameter optimization algorithm for the state-space model with the time-invariant state-transition covariance matrix cited by [8]. The linear prediction coefficients are obtained as

$$[\mathbf{W}_{f,D} \dots \mathbf{W}_{f,L_2-1}] = \left( \sum_{\tau=2}^{L_T} \frac{\mathbf{P}_{f,\tau,\tau-1}^{(2)}}{v_{f,\tau}} \right) \left( \sum_{\tau=2}^{L_T} \frac{\mathbf{P}_{f,\tau-1}^{(2)}}{v_{f,\tau}} \right)^{-1}, \quad (32)$$

where  $\mathbf{P}_{f,\tau-1}^{(2)}$  is the submatrix of  $\mathbf{P}_{f,\tau-1}$  from the  $\{N_m(D-1) + 1\}$ th column to the last column and from the  $N_m(D-1) + 1$ th row to the last row, and  $\mathbf{P}_{f,\tau,\tau-1}^{(2)}$  is the submatrix of  $\mathbf{P}_{f,\tau,\tau-1}$  from the  $\{N_m(D-1) + 1\}$ th column to the last column and from the first row to the  $N_m$ th row. The optimized linear prediction coefficients and sufficient statistics are used to obtain the remaining parameters:

$$v_t = \frac{1}{N_m} \text{trace}\{\mathbf{R}_u^{-1} \mathbf{Q}_{f,\tau}^{(2)}\}, \quad (33)$$

$$\mathbf{R}_{u,f} = \frac{1}{L_T - 1} \sum_{t=2}^{L_T} \frac{1}{v_t} \mathbf{Q}_{f,\tau}^{(2)}, \quad (34)$$

$$\mathbf{R}_{w,f} = \frac{1}{L_T} \sum_{\tau=1}^{L_T} \{\mathbf{x}_{f,\tau} \mathbf{x}_{f,\tau}^H - \mathbf{F} \mathbf{q}_{f,\tau|L_T} \mathbf{q}_{f,\tau|L_T}^H\}, \quad (35)$$

where  $\mathbf{Q}_{f,\tau}^{(2)}$  is the submatrix of  $\mathbf{Q}_{f,\tau}$  from the first column to the  $N_m$ th column and from the first row to the  $N_m$ th row.  $\mathbf{Q}_{f,\tau}$  is defined as

$$\mathbf{Q}_{f,\tau} = \mathbf{P}_{f,\tau} - \mathbf{A}_f \mathbf{P}_{f,\tau,\tau-1}^H - \mathbf{P}_{f,\tau,\tau-1} \mathbf{A}_f^H + \mathbf{H} \mathbf{P}_{f,\tau-1} \mathbf{A}_f^H. \quad (36)$$

The initial MMSE estimate and the initial MSE of the state vector are updated as follows:

$$\boldsymbol{\pi}_f = \mathbf{q}_{f,1|L_T}, \quad (37)$$

$$\mathbf{z}_f = \mathbf{R}_{f,1|L_T}. \quad (38)$$

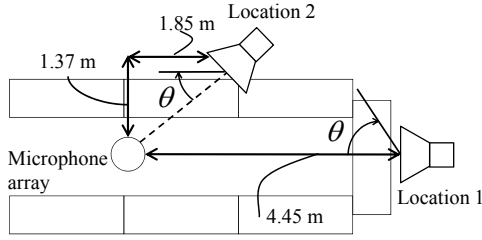
## 4. EVALUATION

### 4.1. Setup

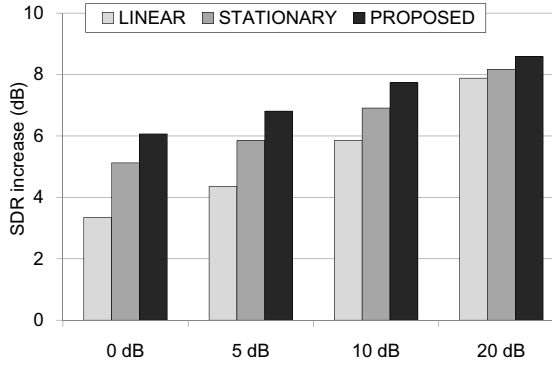
The proposed method was evaluated by using impulse responses measured in a room (Fig. 2,  $RT_{60} = 430$  ms) and a microphone array with three elements. The impulse responses were measured at two locations, with the direction of sound radiation from the loudspeaker  $\theta$  set to 0 degrees. The background noise was also recorded in the same environment using the same microphone array. Original source signals were extracted from the TIMIT database [11] for 34 speakers (one utterance each). The evaluation metric was SDR increase:

$$\begin{aligned} \text{SDR increase} &= 10 \log_{10} \frac{\sum_{n=1}^{L_T} \sum_{m=1}^{N_m} s_m(n)^2}{\sum_{n=1}^{L_T} \sum_{m=1}^{N_m} (s_m(n) - y_m(n))^2} \\ &- 10 \log_{10} \frac{\sum_{n=1}^{L_T} \sum_{m=1}^{N_m} s_m(n)^2}{\sum_{n=1}^{L_T} \sum_{m=1}^{N_m} (s_m(n) - x_m(n))^2}, \end{aligned}$$

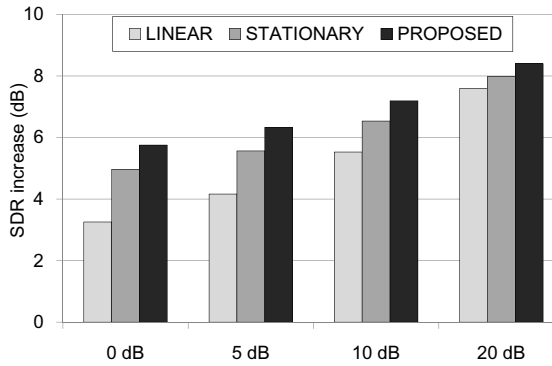
where  $x_m(n)$  is the  $m$ th microphone input signal at the  $n$ th point,  $s_m(n)$  is the desired signal in the  $m$ th microphone at the  $n$ th point



**Fig. 2.** Experimental environment: sound radiation direction  $\theta$  set to 0 degrees.



**Fig. 3.** Increase in SDR at location 1



**Fig. 4.** Increase in SDR at location 2

(which is the summation of the direct sound and early reflection), and  $y_m(n)$  is the dereverberated signal for the  $m$ th channel. The average value of  $\Delta\text{SNR}$  for 34 utterances is used as the metric. The SNR of the microphone input signal was set to 0, 5, 10, or 20 dB. The other parameters were set as are shown in Table 1.

**Table 1.** Evaluation conditions

Sampling rate (Hz)	16,000
Frame size (pt)	1024
Frame shift (pt)	256
Number of microphones $N_m$	3
Number of EM iterations	20

The proposed method (PROPOSED) was compared with two other methods

- **LINEAR:** Dereverberation using linear prediction of microphone input signal. Covariance matrix of microphone input signal is modeled as  $\mathbf{R}_{x,f,\tau} = v_{f,\tau}\mathbf{R}_{u,f}$ , where  $v_{f,\tau}$  and  $\mathbf{R}_{u,f}$  are updated using EM algorithm.
- **STATIONARY:** Dereverberation using linear prediction of microphone input signal and multichannel Wiener-filtering-based stationary background noise reduction. Covariance matrix of microphone input signal is modeled as  $\mathbf{R}_{x,f,\tau} = v_{f,\tau}\mathbf{R}_{u,f} + \mathbf{R}_{w,f}$  where  $v_{f,\tau}$ ,  $\mathbf{R}_{u,f}$ , and  $\mathbf{R}_{w,f}$  are updated using EM algorithm.

As shown in Figs. 3 and 4, the lower the SNR of the microphone input signal, the more effective the proposed method at each location of the speech source. This means the proposed method can reduce reverberation and the background noise effectively even when the microphone input signal is recorded in noisy environments.

## 5. CONCLUSION

Our proposed noise robust speech dereverberation method is based on linear prediction of the noiseless speech signal, which is estimated using the Kalman smoother with the time-varying covariance matrix as the state-transition covariance matrix. Noiseless speech signal estimation and parameter optimization are performed in an iterative manner on the basis of the EM algorithm. Testing showed that the proposed method is better than conventional methods under noisy conditions.

## 6. RELATION TO PRIOR WORK

Previous linear prediction based dereverberation methods [3][4][5] use a blind speech dereverberation technique: linear prediction of the microphone input signal. In contrast, our proposed method uses linear prediction of the noiseless speech signal to reduce degradation of speech dereverberation performance in noisy environments. While a time-varying covariance-matrix-based source separation method has been proposed [10], it is only for source separation without dereverberation. The optimization scheme used in the proposed method is based on optimization using the conventional Kalman smoother [8]. The conventional method uses a time-invariant covariance matrix for the state-transition covariance matrix while our proposed method uses a time-varying covariance matrix so as to meet the time-varying speech characteristics.

## 7. REFERENCES

- [1] M. Miyoshi and Y. Kaneda, "Inverse filtering of room acoustics," *IEEE Trans. ASSP*, vol. 30, no. 2, pp. 145-152, Feb. 1988.
- [2] S. Gannot and M. Moonen, "Subspace methods for multi microphone speech dereverberation," *EURASIP J. Appl. Signal Process.*, vol. 2003, no. 11, pp. 1074-1090, 2003.
- [3] K. Kinoshita, M. Delcroix, T. Nakatani, and M. Miyoshi, "Suppression of Late Reverberation Effect on Speech Signal Using Long-term Multiple-step Linear Prediction," *IEEE Trans. ASLP*, vol. 17, no. 4, pp. 534-545, 2009.
- [4] T. Yoshioka, T. Nakatani, M. Miyoshi, and H.G. Okuno, "Blind separation and dereverberation of speech mixtures by joint optimization," *IEEE Trans. ASLP*, vol. 19, no. 1, pp. 69-84, Jan. 2011.
- [5] M. Togami, Y. Kawaguchi, R. Takeda, Y. Obuchi, and N. Nukaga, "Multichannel speech dereverberation and separation with optimized combination of linear and non-linear filtering," in *Proc. ICASSP2012, 2012*, pp. 4057-4060.
- [6] A.P. Dempster et al., "Maximum likelihood from incomplete data via the EM algorithm," *J. of the Royal Statistic Society, Series B* 39(1), pp. 1-38, 1977.
- [7] A.H. Jazwinski, *Stochastic Processes and Filtering Theory*. Academic Press, 1970.
- [8] Z. Ghahramani and G.E. Hinton, "Parameter estimation for linear dynamical systems," Technical Report CRG-TR-96-2, Department of Computer Science, University of Toronto, 1996.
- [9] T. Nakatani, T. Yoshioka, K. Kinoshita, M. Miyoshi, and B.-H. Juang, "Speech dereverberation based on variance-normalized delayed linear prediction," *IEEE Trans. Speech Audio Process.*, vol. 17, no. 7, pp. 1717-1731, Sep. 2010.
- [10] N.Q.K. Duong, E. Vincent, R. Gribonval, "Under-determined reverberant audio source separation using a full-rank spatial covariance model," *IEEE Trans. Speech Audio Process.*, vol. 18, no. 7, pp. 1830-1840, 2010/9.
- [11] TIMIT corpus, <http://www.ldc.upenn.edu/Catalog/CatalogEntry.jsp?catalogId=LDC93S1>.