

# Stereo Acoustic Echo Cancellation Based on Maximum Likelihood Estimation with Inter-channel-Correlated Echo Compensation

Byung Joon Cho and Hyung-Min Park (지능정보처리연구실)

SOGANG UNIVERSITY

## Abstract

This paper presents batch and online algorithms of a stereo acoustic echo cancellation (SAEC) method. In SAEC, the non-uniqueness problem causes performance degradation, especially for highly coherent far-end signals. In our method, this problem can be avoided without an additional decorrelation preprocessor or multi-microphones by overestimating far-end echoes and compensating for the overestimated inter-channel-correlated echo to obtain a desired echo-canceled signal. In addition, our method is based on the maximum likelihood estimation (MLE) criterion of the echo-canceled signal under the assumption that the signal in the time-frequency domain follows a zero-mean complex Gaussian distribution with a time-varying variance. Furthermore, a variable forgetting factor based on the cross spectral density (CSD) between the echo-canceled signal and a far-end echo is presented in the online algorithm to improve the convergence of adaptive filters with a high cancellation performance when converged. Experimental results under various conditions demonstrate that the proposed method can successfully perform SAEC even in the presence of inter-channel correlation, double-talk, and abrupt echo path changes.

## Proposed SAEC algorithm

We assume that the echo-canceled signal in each frequency bin follows a zero-mean complex Gaussian distribution with a time-varying variance to cope with the non-stationarity of the signal.

- The input signal for the  $v$ -th channel

$$X^{(v)}[m, k] = S[m, k] + D_u^{(v)}[m, k] + (\omega_k^{(1)} + \omega_k^{(2)}) D_c[m, k]$$

- The desired echo-canceled signal using  $X^{(v)}[m, k]$

$$Y[m, k] = X^{(v)}[m, k] - \mathbf{g}_k^{(v)H} \mathbf{u}_{m,k}^{(v)}$$

- The SAEC problem to find the unknown parameter sets

$$\hat{\Theta}_{m,k}^{(v)} = \arg \max_{\Theta_{m,k}^{(v)}} \mathcal{L}(\Theta_{m,k}^{(v)})$$

$$\mathcal{L}(\Theta_{m,k}^{(v)}) = \sum_m \log p(Y[m, k]; \Theta_{m,k}^{(v)}) = - \sum_m \left( \log(\pi \lambda_Y[m, k]) + \frac{|X^{(v)}[m, k] - \mathbf{g}_k^{(v)H} \mathbf{u}_{m,k}^{(v)}|^2}{\lambda_Y[m, k]} \right)$$

- The iterate update rule for  $\mathbf{g}_k^{(v)}$

$$\mathbf{g}_k^{(v)} = \mathbf{R}_k^{(v)-1} \mathbf{r}_k^{(v)} \quad \mathbf{R}_k^{(v)} = \sum_m \frac{\mathbf{u}_{m,k}^{(v)} \mathbf{u}_{m,k}^{(v)H}}{\lambda_Y[m, k]} \quad \mathbf{r}_k^{(v)} = \sum_m \frac{\mathbf{u}_{m,k}^{(v)} X^{(v)*}[m, k]}{\lambda_Y[m, k]} \quad \lambda_Y[m, k] = |Y[m, k]|^2$$

- The overestimation ratio

$$\zeta_k = \min \left( \frac{2 |\Phi^{(1)(2)}[k]|}{\Phi^{(1)}[k] + \Phi^{(2)}[k]}, 1 \right) \quad \Phi^{(v)}[k] = \frac{1}{N_m} \sum_m |\mathbf{g}_k^{(v)H} \mathbf{u}_{m,k}^{(v)}|^2$$

$$\Phi^{(1)(2)}[k] = \frac{1}{N_m} \sum_m (\mathbf{g}_k^{(1)H} \mathbf{u}_{m,k}^{(1)})^* (\mathbf{g}_k^{(2)H} \mathbf{u}_{m,k}^{(2)})$$

Online SAEC based on frame-by-frame processing is required to cope with various environments in which the echo path may change. Similar to the batch parameter update rules, online parameter update rules can be derived based on RLS. The log-likelihood function to derive the online algorithm is defined as follows:

$$\mathcal{L}(\Theta_{m,k}^{(v)}) = \sum_{\mu=1}^m \gamma^{m-\mu} \log p(Y[\mu, k]; \Theta_{\mu,k}^{(v)}) = - \sum_{\mu=1}^m \gamma^{m-\mu} \left( \log(\pi \lambda_Y[\mu, k]) + \frac{|X^{(v)}[\mu, k] - \mathbf{g}_{\mu,k}^{(v)H} \mathbf{u}_{\mu,k}^{(v)}|^2}{\lambda_Y[\mu, k]} \right)$$

- Update of the adaptive filter for the first channel

$$1 : \hat{\Phi}^{(v)}[m, k] = \alpha \Phi^{(v)}[m-1, k] + (1-\alpha) |\mathbf{g}_{m-1,k}^{(v)H} \mathbf{u}_{m,k}^{(v)}|^2, v \in \{1, 2\}.$$

$$2 : \hat{\Phi}^{(1)(2)}[m, k] = \alpha \Phi^{(1)(2)}[m-1, k] + (1-\alpha) (\mathbf{g}_{m-1,k}^{(1)H} \mathbf{u}_{m,k}^{(1)})^* (\mathbf{g}_{m-1,k}^{(2)H} \mathbf{u}_{m,k}^{(2)})$$

$$3 : \hat{\zeta}_{m,k} = \min \left( \frac{2 |\hat{\Phi}^{(1)(2)}[m, k]|}{\hat{\Phi}^{(1)}[m, k] + \hat{\Phi}^{(2)}[m, k]}, 1 \right)$$

$$4 : \hat{Y}[m, k] = X[m, k] - \sum_{v \in \nu} \mathbf{g}_{m-1,k}^{(v)H} \mathbf{u}_{m,k}^{(v)} + \frac{\hat{\zeta}_{m,k}}{2} \sum_{v \in \nu} (\mathbf{g}_{m-1,k}^{(v)H} \mathbf{u}_{m,k}^{(v)})^* (\mathbf{g}_{m-1,k}^{(v)H} \mathbf{u}_{m,k}^{(v)})$$

$$5 : \hat{\lambda}_Y[m, k] = |\hat{Y}[m, k]|^2$$

$$6 : \mathbf{k}_{m,k}^{(1)} = \frac{\Psi_{m-1,k}^{(1)-1} \mathbf{u}_{m,k}^{(1)}}{\gamma \hat{\lambda}_Y[m, k] + \mathbf{u}_{m,k}^{(1)H} \Psi_{m-1,k}^{(1)-1} \mathbf{u}_{m,k}^{(1)}}$$

$$7 : \Psi_{m,k}^{(1)-1} = \gamma^{-1} (\Psi_{m-1,k}^{(1)-1} - \mathbf{k}_{m,k}^{(1)} \mathbf{u}_{m,k}^{(1)H} \Psi_{m-1,k}^{(1)-1})$$

$$8 : \mathbf{g}_{m,k}^{(1)} = \mathbf{g}_{m-1,k}^{(1)} + \mathbf{k}_{m,k}^{(1)} \hat{Y}^*[m, k]$$

- Update of the adaptive filter for the second channel

$$9 : \Phi^{(1)}[m, k] = \alpha \Phi^{(1)}[m-1, k] + (1-\alpha) |\mathbf{g}_{m,k}^{(1)H} \mathbf{u}_{m,k}^{(1)}|^2$$

$$10 : \hat{\Phi}^{(1)(2)}[m, k] = \alpha \Phi^{(1)(2)}[m-1, k] + (1-\alpha) (\mathbf{g}_{m,k}^{(1)H} \mathbf{u}_{m,k}^{(1)})^* (\mathbf{g}_{m,k}^{(2)H} \mathbf{u}_{m,k}^{(2)})$$

$$11 : \hat{\zeta}_{m,k} = \min \left( \frac{2 |\hat{\Phi}^{(1)(2)}[m, k]|}{\hat{\Phi}^{(1)}[m, k] + \hat{\Phi}^{(2)}[m, k]}, 1 \right)$$

$$12 : \hat{Y}[m, k] = X[m, k] - \mathbf{g}_{m,k}^{(1)H} \mathbf{u}_{m,k}^{(1)} - \mathbf{g}_{m,k}^{(2)H} \mathbf{u}_{m,k}^{(2)} + \frac{\hat{\zeta}_{m,k}}{2} (\mathbf{g}_{m,k}^{(1)H} \mathbf{u}_{m,k}^{(1)})^* (\mathbf{g}_{m,k}^{(2)H} \mathbf{u}_{m,k}^{(2)})$$

$$13 : \text{Repeat Steps 5-8 with } \hat{Y}[m, k] \text{ and (1) replaced by } \hat{Y}[m, k] \text{ and (2).}$$

- Calculation of the echo canceled signal

$$14 : \Phi^{(2)}[m, k] = \alpha \Phi^{(2)}[m-1, k] + (1-\alpha) |\mathbf{g}_{m,k}^{(2)H} \mathbf{u}_{m,k}^{(2)}|^2$$

$$15 : \Phi^{(1)(2)}[m, k] = \alpha \Phi^{(1)(2)}[m-1, k] + (1-\alpha) (\mathbf{g}_{m,k}^{(1)H} \mathbf{u}_{m,k}^{(1)})^* (\mathbf{g}_{m,k}^{(2)H} \mathbf{u}_{m,k}^{(2)})$$

$$16 : \zeta_{m,k} = \min \left( \frac{2 |\Phi^{(1)(2)}[m, k]|}{\Phi^{(1)}[m, k] + \Phi^{(2)}[m, k]}, 1 \right) \quad 17 : Y[m, k] = X[m, k] - \sum_{v \in \nu} \mathbf{g}_{m,k}^{(v)H} \mathbf{u}_{m,k}^{(v)} + \frac{\zeta_{m,k}}{2} \sum_{v \in \nu} \mathbf{g}_{m,k}^{(v)H} \mathbf{u}_{m,k}^{(v)}$$

## Problem Formulation

The microphone input signal

$$X[m, k] \approx S[m, k] + \sum_{v \in \nu} \sum_{l=0}^{L-1} H^{(v)*}[l, k] U^{(v)}[m-l, k]$$

The echo-canceled signal

$$Y[m, k] \approx X[m, k] - \sum_{v \in \nu} \sum_{l=0}^{L-1} \hat{H}^{(v)*}[l, k] U^{(v)}[m-l, k]$$

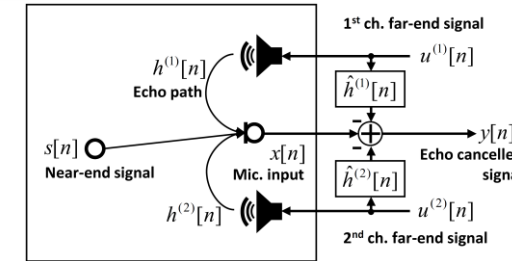


Fig. 1. Block diagram of the SAEC in the time domain

We estimate far-end echoes without considering the inter-channel correlation first, to deal with SAEC without the additional decorrelation. However, if the echo of each channel is estimated, the inter-channel-correlated echo is overestimated due to the overlap between the two estimated echoes. Therefore, we compensate for the overestimated inter-channel-correlated echo based on an overestimation ratio to avoid an oversubtraction of the inter-channel-correlated echo. The echo-canceled signal is calculated as follows:

$$Y[m, k] = X[m, k] - \sum_{v \in \nu} \mathbf{g}_k^{(v)H} \mathbf{u}_{m,k}^{(v)} + \frac{\zeta_k}{2} \sum_{v \in \nu} \mathbf{g}_k^{(v)H} \mathbf{u}_{m,k}^{(v)} \quad - \zeta_k : \text{The overestimation ratio to compensate for the inter-channel-correlated echo}$$

## Experimental Results

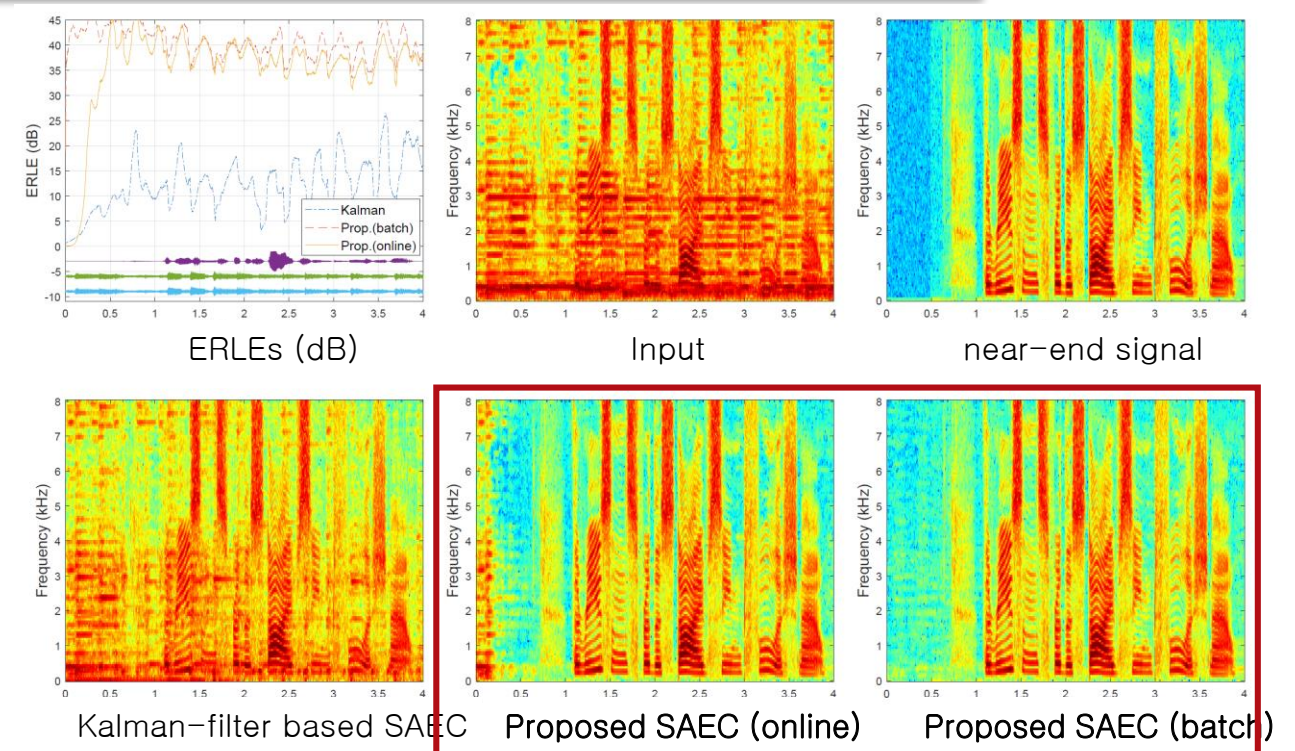


Fig. 2. ERLEs and spectrograms when the far-end signals were identical music signals.

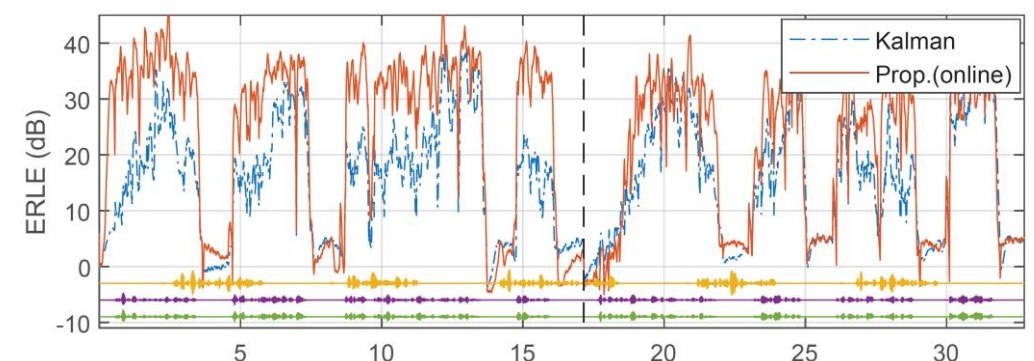


Fig. 2. ERLEs of the Kalman-filter-based and proposed online SAEC methods without the decorrelation preprocessor. Both echo paths changed in approximately 17.15 s, as denoted by the black dashed line.

Table 1. ERLE (dB) averaged over 370 utterances

Input SER	IR type	Far-end signals	Method		
			Kalman	Prop. (online)	Prop. (batch)
0 dB	Random IR	Speech	8.08	26.08	<b>32.62</b>
		Music	12.03	31.57	<b>39.02</b>
		Speech, Music	7.44	21.42	<b>33.07</b>
	RIR	Speech	16.66	26.12	<b>32.44</b>
		Music	19.50	32.10	<b>39.20</b>
		Speech, Music	17.03	22.18	<b>32.70</b>
5 dB	Random IR	Speech	7.00	24.07	<b>30.56</b>
		Music	10.98	29.91	<b>36.91</b>
		Speech, Music	6.70	20.29	<b>32.24</b>
	RIR	Speech	14.82	23.54	<b>30.30</b>
		Music	17.73	30.87	<b>37.10</b>
		Speech, Music	15.34	21.05	<b>31.89</b>
10 dB	Random IR	Speech	5.88	21.74	<b>27.55</b>
		Music	9.85	27.90	<b>33.67</b>
		Speech, Music	5.91	18.92	<b>29.49</b>
	RIR	Speech	12.88	20.67	<b>27.24</b>
		Music	16.03	29.27	<b>33.89</b>
		Speech, Music	13.71	19.61	<b>29.36</b>

Table 2. PESQ scores averaged over 370 utterances

Input SER	IR type	Far-end signals	Method			
			No proc.	Kalman	Prop. (online)	Prop. (batch)
0 dB	Random IR	Speech	1.87	2.33	3.86	<b>4.06</b>
		Music	1.91	2.52	4.05	<b>4.23</b>
		Speech, Music	1.80	2.08	3.31	<b>3.98</b>
	RIR	Speech	1.83	2.95	3.77	<b>4.06</b>
		Music	1.66	2.94	3.85	<b>4.16</b>
		Speech, Music	1.63	2.74	3.20	<b>3.90</b>
5 dB	Random IR	Speech	2.20	2.58	4.01	<b>4.19</b>
		Music	2.27	2.80	4.17	<b>4.30</b>
		Speech, Music	2.16	2.40	3.56	<b>4.15</b>
	RIR	Speech	2.15	3.16	3.91	<b>4.19</b>
		Music	2.00	3.15	4.04	<b>4.27</b>
		Speech, Music	1.98	2.97	3.47	<b>4.09</b>
10 dB	Random IR	Speech	2.53	2.84	4.11	<b>4.26</b>
		Music	2.61	3.07	4.25	<b>4.34</b>
		Speech, Music	2.51	2.69	3.78	<b>4.22</b>
	RIR	Speech	2.49	3.37	4.02	<b>4.26</b>
		Music	2.35	3.35	4.18	<b>4.32</b>
		Speech, Music	2.33	3.19	3.70	<b>4.20</b>

## Conclusion

In this paper, we proposed and derived batch and online algorithms of an SAEC method based on the MLE of an echo-canceled signal by assuming that the signal in the time-frequency domain follows a zero-mean complex Gaussian distribution with a time-varying variance. To avoid the non-uniqueness problem without an additional decorrelation preprocessor or multi-microphones even for highly coherent far-end signals, the proposed method obtained an echo-canceled signal by overestimating far-end echoes and compensating for an overestimated inter-channel-correlated echo based on the overestimation ratio. In addition, a variable forgetting factor based on the CSD between the echo-canceled signal and a far-end echo in the online algorithm was presented to improve the convergence of the adaptive filters significantly with a high ERLE when converged. Experimental results under various conditions demonstrated that the proposed method achieved a superior performance compared with the Kalman-filter-based method even in the presence of inter-channel correlation, double-talk, and abrupt echo path changes.