

MULTICHANNEL ACOUSTIC ECHO CANCELLATION IN MULTIPARTY SPATIAL AUDIO CONFERENCING WITH CONSTRAINED KALMAN FILTERING

Zhengyou Zhang, Qin Cai, Jack W. Stokes

Microsoft Research, One Microsoft Way, Redmond, WA 98052, USA

ABSTRACT

This paper proposes a novel scheme for acoustic echo cancellation (AEC) in multi-party, spatialized audio conferencing by exploring the constraints among the adaptive filters for each synthesized remote speaker. The AEC algorithm employs constrained Kalman filtering (CKF-AEC) that takes advantage of the uncorrelated reference signal from each remote channel [1], yet restricts the filter adaptation within the subspace determined by the spatialization functions used to virtualize each remote participant and the number of unique acoustic paths between the loudspeakers and microphone. In addition, the proposed algorithm allows the adaptation for channels without input signals. Experimental results show the proposed algorithm has much higher performance than the scheme proposed in [1] which uses NLMS as the adaptive filter and does not explore the available constraint.

Index Terms— Spatial Audio, Multi-channel AEC, Kalman Filtering, Multi-party conferencing.

INTRODUCTION

Recently, the demand for better quality teleconferencing with remote multi-participants has risen rapidly for reducing travel cost and increasing productivity. It has been shown that spatial mapping of the remote participants' voice signals in multiple virtual positions enhances the collaboration experience during multi-party conferencing [2]. It is well known that stereo AEC suffers from the mis-convergence problem [3]. Researchers have proposed various ways to alleviate the problem including adding nonlinearities [4] or additional noise [3,5] to de-correlate the speaker signals. In the paper, we propose an acoustic echo cancellation algorithm based on the constrained Kalman filter (CKF-AEC) for full-duplex, multi-party, spatialized audio conferencing using Voice over Internet Protocol (VoIP) which does not distort the speaker signal. Kalman filtering has not been used widely for AEC except in [6] which directly considers the varying acoustic paths as the hidden state variable. The multi-party conferencing problem with spatialized audio was previously studied in [1], and we adopt the same synthetic stereophonic structure depicted in Fig. 1. However in [1], each canceller is adapted

independently without realizing that the only unknown parameters are the room impulse responses (RIRs) between the loudspeakers and microphones; the method also ignores the fact that signals are continuously sent to each loudspeaker. If one remote participant has been silent for a long time, noticeable echo will be heard when the participant resumes speaking. Another disadvantage of the method in [1] is that it has very poor performance when multiple remote participants speak simultaneously; the authors in [1] analyzed the convergence for adapting four remote channels talking at the same time using white Gaussian noise. However in our experiments with speech signals, the joint NLMS-based multi-channel AEC failed.

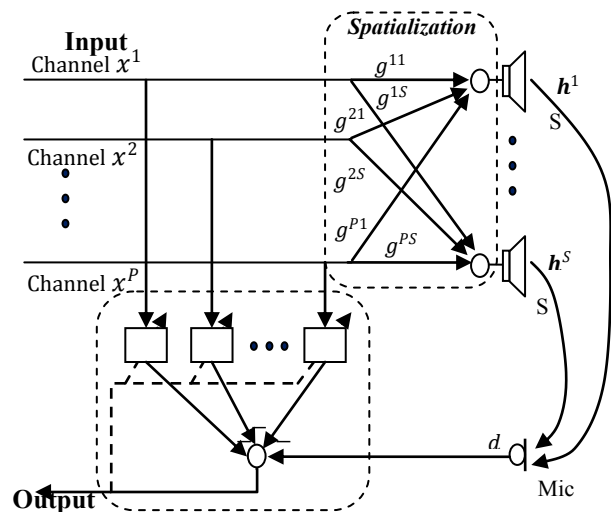


Fig. 1 Synthetic stereo structure with a single canceller per remote-channel

The CKF-AEC algorithm proposed in this paper employs a constrained Kalman filtering mechanism using the signal from each remote channel as reference. It takes advantage that the signals from each remote channel are uncorrelated and the spatialization parameters are known. Not only does the algorithm adapt the per channel impulse response (CIR), it also resolves the RIR between each loudspeaker and microphone. As long as one channel is active, the RIRs and CIRs for the other channels are adapted simultaneously. In other words, if a remote participant is silent for a long time, the adaptation from the other channels will benefit the

adaptation for this participant. The CKF-AEC algorithm does not require channel switching as in multiple mono NLMS-based AEC; it even adapts with multiple remote participants speaking simultaneously. Finally to significantly reduce the CPU consumption, we propose an implementation to avoid the matrix inversion that is $O(n^3)$ for the standard KF.

CONSTRAINED KALMAN FILTERING

We describe the proposed algorithm in the frequency domain due to ease of the implementation since the spatial mapping only results in multiplication of a complex number in the frequency domain. The following derivation is conducted in each frequency sub-band.

2.1. Problem Statement

The superscript $*$ denotes the complex conjugate. Column vectors are assumed. The superscript T is the matrix/vector transpose. The superscript H is the Hermitian transpose (complex conjugate of the transpose). The subscript t will be used to indicate time instant (or frame number). Assume there are P remote participants/channels with audio inputs $\{X^i | i = 1, \dots, P\}$ and S loudspeakers. Each remote audio channel X^i is spatialized across the loudspeakers with gain and delay modulation G^{is} for speaker s . The audio played on speaker s is thus $Y^s = \sum_{i=1}^P G^{is} X^i$. A microphone is used to capture the audio in the local room, and is denoted by D . The captured audio includes the audio of the local participant and that played on the speakers. The audio played on speaker s is transmitted through the RIR modeled by an FIR filter with L taps, denoted by $H_t^s = [H_t^s, H_{t-1}^s, \dots, H_{t-L+1}^s]^T$.

Assuming that the local or near end participant is not talking, the audio captured by the microphone at time t , i.e., the echo, should be

$$\begin{aligned} D_t &= \sum_{i=1}^P \sum_{s=1}^S G^{is} (H_t^s X_t^i + H_{t-1}^s X_{t-1}^i + \dots + H_{t-L+1}^s X_{t-L+1}^i) \\ &= \sum_{i=1}^P \sum_{s=1}^S G^{is} \mathbf{H}_t^{sT} \mathbf{X}_t^i \end{aligned}$$

where $\mathbf{X}_t^i = [X_t^i, X_{t-1}^i, \dots, X_{t-L+1}^i]^T$ is speech from remote participant i .

The problem is to design an echo cancellation filter with L taps for each remote participant i ,

$$\mathbf{W}_t^i = [W_t^i, W_{t-1}^i, \dots, W_{t-L+1}^i]^T$$

such that the echo is cancelled. That is, determine \mathbf{W}_t^i 's so

$$D_t - \sum_{i=1}^P \mathbf{W}_t^{iT} \mathbf{X}_t^i = 0. \quad (1)$$

It is clear that we have

$$\mathbf{W}_t^i = \sum_{s=1}^S G^{is} \mathbf{H}_t^s. \quad (2)$$

Therefore, the echo cancellation filters are not mutually independent. We will leverage this constraint to update each echo cancellation filter simultaneously even though the corresponding remote participant is not talking. If the local participant is talking, then the system outputs the echo-cancelled signal: $D_t - \sum_{i=1}^P \mathbf{W}_t^{iT} \mathbf{X}_t^i$.

2.2. Kalman Filtering Formulation

Let the state vector be a $(P + S)L$ -dimensional vector:

$$\begin{aligned} \mathbf{S}_t &= [\mathbf{W}_t^{1T}, \dots, \mathbf{W}_t^{PT}, \mathbf{H}_t^{1T}, \dots, \mathbf{H}_t^{ST}]^T \\ &= [W_t^1, W_{t-1}^1, \dots, W_{t-L+1}^1, \dots, W_t^P, W_{t-1}^P, \dots, W_{t-L+1}^P, \\ &\quad H_t^1, H_{t-1}^1, \dots, H_{t-L+1}^1, \dots, H_t^S, H_{t-1}^S, \dots, H_{t-L+1}^S]^T. \end{aligned}$$

The state is not expected to remain constant, and the state equation is

$$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{n}_t \quad (3)$$

where \mathbf{n}_t is the system noise. The parts of \mathbf{n}_t corresponding to the H elements model the variation due to changes in the acoustical environment such as the movement of the local participant. The observation equation is:

$$D_t = \mathbf{A}_t^T \mathbf{S}_t + v_t, \quad (4)$$

where v_t is the observation noise (microphone noise and ambient noise) and can also model the inaccuracy of the observation system. The vector \mathbf{A}_t is a $(P + S)L$ -dimensional vector given by

$$\mathbf{A}_t = [\mathbf{X}_t^{1T}, \dots, \mathbf{X}_t^{PT}, \mathbf{0}^{1T}, \dots, \mathbf{0}^{ST}]^T$$

with $\mathbf{0}$ being the L -dimensional zero vector. The constraint on the state parameters (2) is a set of PL linear equations, and can be written as

$$\mathbf{C} \mathbf{S}_t = \mathbf{0} \quad (5)$$

where \mathbf{C} is a $PL \times (P + S)L$ matrix, given by

$$\mathbf{C} = \begin{bmatrix} -\mathbf{1} & & G^{11} \mathbf{1} & \dots & G^{1S} \mathbf{1} \\ & \ddots & \vdots & \ddots & \vdots \\ & & -\mathbf{1} & G^{P1} \mathbf{1} & \dots & G^{PS} \mathbf{1} \end{bmatrix}$$

with $\mathbf{1}$ being the $L \times L$ identity matrix, i.e., $\mathbf{1} = \text{diag}(1, \dots, 1)$. The state constraint can be considered as perfect observation. Thus, we combine the original observation equation and the state constraint into a new observation equation:

$$\mathbf{Y}_t = \mathbf{B}_t \mathbf{S}_t + \mathbf{v}_t \quad (6)$$

where

$$\mathbf{Y}_t = [D_t, 0, \dots, 0]^T, \quad \mathbf{B}_t = \begin{bmatrix} \mathbf{A}_t^T \\ \mathbf{C} \end{bmatrix}, \quad \text{and} \quad \mathbf{v}_t = \begin{bmatrix} v_t \\ \mathbf{u}_t \end{bmatrix}$$

with \mathbf{u}_t being the noise term for the state constraint. If the constraint is exact, then $\mathbf{u}_t = \mathbf{0}$.

We assume that the following conditions are satisfied:

$$E[\mathbf{n}_t] = \mathbf{0}, \quad E[v_t] = 0, \quad E[\mathbf{n}_t \mathbf{n}_t^T] = \mathbf{Q}_t \quad \text{and}$$

$$E[v_t v_t^T] = \mathbf{R}_t = \begin{bmatrix} \sigma_t^2 & \mathbf{0}^T \\ \mathbf{0} & \Lambda_t \end{bmatrix}$$

where Λ_t is the covariance matrix of the noise term for the state constraint. If we indeed want to impose the constraint fully, then $\Lambda_t = \mathbf{0}$. In practice, we sometimes prefer to impose a soft constraint to have a more stable system on account of nonlinearity in loudspeakers and clock drift between loudspeaker and sound capture. We can also start Λ_t with a larger value and gradually decrease it over time.

We can now solve the multichannel AEC problem for spatialized audio using the Kalman filter. We use superscript “ $-$ ” for the prediction, and \mathbf{P} for the covariance matrix of the error in the estimated state vector. The Kalman filter equations are given by:

$$\begin{aligned} \mathbf{S}_t^- &= \mathbf{S}_{t-1} \\ \mathbf{P}_t^- &= \mathbf{P}_{t-1} + \mathbf{Q}_t \\ \mathbf{K}_t &= \mathbf{P}_t^- \mathbf{B}_t^H (\mathbf{B}_t \mathbf{P}_t^- \mathbf{B}_t^H + \mathbf{R}_t)^{-1} \\ \mathbf{S}_t &= \mathbf{S}_t^- + \mathbf{K}_t (\mathbf{Y}_t - \mathbf{B}_t \mathbf{S}_t^-) \\ \mathbf{P}_t &= (\mathbf{I} - \mathbf{K}_t \mathbf{B}_t) \mathbf{P}_t^- \end{aligned}$$

The state vector can be initialized to 0, i.e., $\mathbf{S}_0^- = \mathbf{0}$. The covariance matrix \mathbf{P}_0^- should be set to a large value to reflect that we do not have knowledge of the state vector.

The advantages of the CKF algorithm are: 1) The constraint is taken care of automatically, and can be imposed with varying degrees, 2) All channels are taken into account simultaneously. Thus, overlapping far-end talking is not an issue, 3) The AEC for each channel is updated continuously because of the constraint, even if that channel is not active. Therefore, AEC is always up to date. 4) Ambient noise σ_t^2 can be time varying. A separate noise tracker can be used to provide that information.

However, we are also aware of the major drawback is the need to invert a $(PL + 1) \times (PL + 1)$ matrix, rather than a scalar inversion. Matrix inversion has $O((PL + 1)^3)$ complexity. Thus there is a significant increase in computational cost.

2.3. An Improved Implementation

In our previous discussion, we added the PL constraints to the set of measurement equations, and the constraints are imposed softly by a Gaussian noise vector \mathbf{u}_t with mean $\mathbf{0}$ and covariance matrix Λ_t . In practice, the covariance matrix Λ_t is usually set to be a diagonal matrix by assuming no correlation between the constraints. Thus the covariance matrix of the measurement noise vector \mathbf{R}_t is diagonal.

Let $v_{t,j}$ be the j -th element of the measurement noise vector \mathbf{v}_t , $\mathbf{R}_t = \text{diag}(r_{t,1}, \dots, r_{t,j}, \dots, r_{t,PL})$ be the diagonal covariance matrix, $\mathbf{b}_{t,j}^T$ be the j -th row of the measurement matrix \mathbf{B}_t , and $Y_{t,j}$ be the j -th element of the measurement vector \mathbf{Y}_t , then the vector measurement equation (6) is equivalent to the following $PL + 1$ scalar measurement equations:

$$Y_{t,j} = \mathbf{b}_{t,j}^T \mathbf{S}_t + v_{t,j} \quad \text{for } j = 1, \dots, PL + 1$$

with $E[v_{t,j}] = 0$ and $E[v_{t,j}^2] = r_{t,j}$. Now, we can apply the

Kalman filter sequentially for each of the above scalar measurements, and inversion of a $(PL + 1) \times (PL + 1)$ matrix is avoided. The complexity is reduced from $O((PL + 1)^3)$ to $O(PL + 1)$. This substantially reduces the computational cost.

In summary, the new algorithm can be described as follows:

Prediction: $\mathbf{S}_t^- = \mathbf{S}_{t-1}$
 $\mathbf{P}_t^- = \mathbf{P}_{t-1} + \mathbf{Q}_t$

Update: Let $\mathbf{S}_{t,0}^- \triangleq \mathbf{S}_t^-$, and $\mathbf{P}_{t,0}^- \triangleq \mathbf{P}_t^-$.

For $j = 1, \dots, PL + 1$, do

$$\begin{aligned} \mathbf{K}_{t,j} &= \mathbf{P}_{t,j-1}^- \mathbf{b}_{t,j}^* (\mathbf{b}_{t,j}^T \mathbf{P}_{t,j-1}^- \mathbf{b}_{t,j}^* + r_{t,j})^{-1} \\ \mathbf{S}_{t,j}^- &= \mathbf{S}_t^- + \mathbf{K}_{t,j} (Y_{t,j} - \mathbf{b}_{t,j}^T \mathbf{S}_{t,j-1}^-) \\ \mathbf{P}_{t,j}^- &= (\mathbf{I} - \mathbf{K}_{t,j} \mathbf{b}_{t,j}^T) \mathbf{P}_{t,j-1}^- \end{aligned}$$

Finally, we set $\mathbf{S}_t \triangleq \mathbf{S}_{t,PL+1}^-$ and $\mathbf{P}_t \triangleq \mathbf{P}_{t,PL+1}^-$.

If a constraint is already satisfied, then $Y_{t,j} - \mathbf{b}_{t,j}^T \mathbf{S}_{t,j-1}^- = 0$, leaving the state vector unaffected. If not, a correction is added by projecting to the direction $\mathbf{K}_{t,j} \propto \mathbf{P}_{t,j-1}^- \mathbf{b}_{t,j}^*$.

EXPERIMENTAL RESULTS

We consider four remote participants with spatialization on two loudspeakers with virtual positions at $[-30^\circ, 30^\circ, 0^\circ, -45^\circ]$, which is the angle between the virtual speaker and the center line originating from the microphone, i.e., 0° gives an illusion the participant is positioned in the middle front. We used 16kHz speech. Each remote participant speaks for 4 seconds. A known RIR is used to generate the microphone input and -20dB white Gaussian noise is added. The number of sub-bands is 512 with 256 samples per frame. The best result is achieved with $\sigma_t = -20\text{dB}$ at the lowest bands, $\Lambda_t = 1e-2$ and $\mathbf{Q}_t = \mathbf{0}$. Fig. 2 shows the comparison between our CKF-AEC and the four mono NLMS filters proposed in [1] (but with our frequency domain implementation). Whenever a new talker starts, the Error Return Loss Enhancement (ERLE) drops significantly in NLMS. However after the second talker, thanks to the constraint, the CKF-AEC handles new participants gracefully. In Fig. 3, we conducted another experiment with two participants talking simultaneously in consecutive positions and compared the CKF-AEC algorithm with the joint NLMS which adapts parameters for four channels simultaneously, described in [1]. It is shown that joint NLMS performs poorly while the CKF-AEC maintains good performance. The authors in [1] only showed convergence with white Gaussian noise input, but not with speech signals. Finally, we add -30dB change on RIRs for every 50 frames and add \mathbf{Q}_t accordingly. Fig 4 shows the case where the first channel is re-activated at the end of the third channel or at 12 seconds. This is to check whether CKF-AEC can update a canceller even when the

corresponding channel is silent. Indeed, CKF-AEC handles the transition more smoothly than the four mono NLMS after the second participant. We repeat the same experiment with Kalman filtering without the constraints, glitches similar to NLMS happened for each transition, as shown in Fig 5.

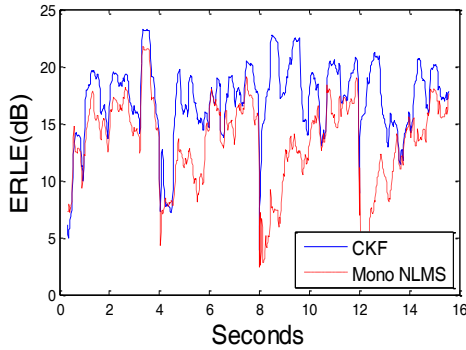


Fig. 2 AEC Performance between CKF and 4 mono NLMS

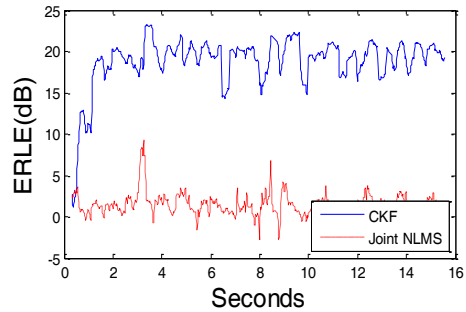


Fig. 3 AEC Performance between CKF and joint NLMS

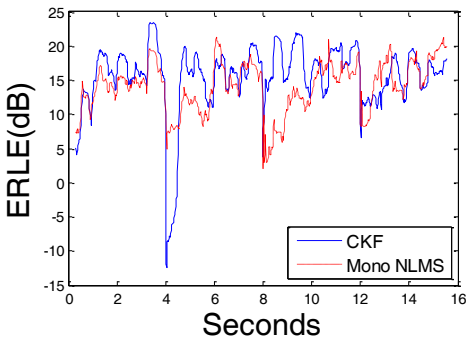


Fig. 4 AEC performance with changing RIRs

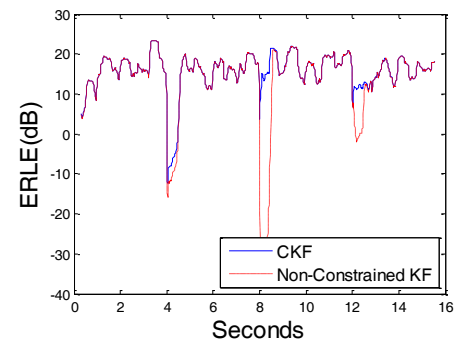


Fig 5. Performance between CKF and non-constrained KF

To verify if the adaptive filters converge to the optimal solution, we calculated the mean values of misalignments using known S_t in the frequency domain. Two data sets were used. Table 1 shows that CKF-AEC converges for both data sets. Mono NLMS also converges for each CIR, but without estimation of the RIRs.

	CKF	Mono NLMS
$\overline{\text{misalign}(w_i)}$ (set 1)	-22	-21
$\overline{\text{misalign}(h_i)}$ (set 1)	-19	N/A
$\overline{\text{misalign}(w_i)}$ (set 2)	-24	-18
$\overline{\text{misalign}(h_i)}$ (set 2)	-23	N/A

Table 1 Misalignment evaluation in dB.

CONCLUSION

We proposed a novel algorithm using Kalman filtering for multi-channel AEC during spatial audio conferencing by imposing constraints among each canceller during adaptation. The imposed constraints guide the adaptation toward the optimal solution. They allow for estimation of the RIRs as a byproduct and for adaptation for channels without input signals. Experimental results show that the algorithm outperforms the array of mono NLMS when a new channel starts. It also adapts when speech co-exists among multiple channels without need of channel switching. As our future work, we plan to compare our algorithm with stereo AEC algorithms such as the one proposed in [4].

ACKNOWLEDGMENTS

The authors thank Philip A. Chou, Dinei Florencio, and Zicheng Liu for many fruitful discussions.

REFERENCES

- [1] T. N. Yensen, R. A. Goubran, and I. Lambadaris, "Synthetic Stereo Acoustic Echo Cancellation Structure for Multiple Participant VoIP Conference" *IEEE Transaction on Speech and Audio Processing*, vol. 9, no. 2, pp. 168-174, Feb. 2001.
- [2] J. Baldis, "Effects of spatial audio on memory, compression, and preference during desktop conferences," *Proceedings of the CHI Conference on Human Factors in Computing Systems*, ACM Press, Seattle, 2001.
- [3] M. Sonfhi, D. Morgan, and J. Hall, "Stereophonic acoustic echo cancellation – An overview of the fundamental problem," *IEEE Speech Processing Letter*, vol. 2, pp. 148-151, 1995.
- [4] J. Benesty, D. Morgan, and M. Sonfhi, "A better understanding and an improved solution to the problem of stereophonic acoustic echo cancellation," *IEEE Transaction on Speech and Audio Processing*, vol. 6, pp. 156-165, 1998.
- [5] A. Gilloire, and V. Turbin "Using auditory properties to improve the behavior of stereophonic acoustic cancellers," *Proceedings of IEEE Intl. Conf. Acoustic, Speech and Signal Processing*, pp. 3681-3684, 1998.
- [6] G. Enzner, and P. Vary, "Frequency-domain adaptive Kalman filter for acoustic echo control in hands-free telephones," *Signal Processing*, vol. 86, pp. 1140-1156, Oct, 2006.