

Multichannel Acoustic Echo Cancellation in Multiparty Spatial Audio Conferencing With Constrained Kalman Filtering

Zhengyou Zhang, Qin Cai, Jay Stokes

Communications and Collaboration Systems
Microsoft Research, Redmond, WA, USA



zhang@microsoft.com

Outline

- Motivations
- Audio spatialization
- Two approaches to multichannel AEC
- Constrained Kalman filtering
- Experimental results

Motivation

- Current audio conferencing systems
 - Monaural → adequate for 1-to-1
 - Poor when #people > 2
- Why poor?
 - All the voice streams are intermixed into a single one
 - Huge cognitive load: *Do 2 things simultaneously*
 - *Associate voice signals to the speaker*
 - *Comprehend what is being discussed*

Solutions

- Video conferencing
- Spatial audio conferencing
- Spatial audio + Video
- Immersive conferencing

Solutions

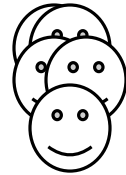
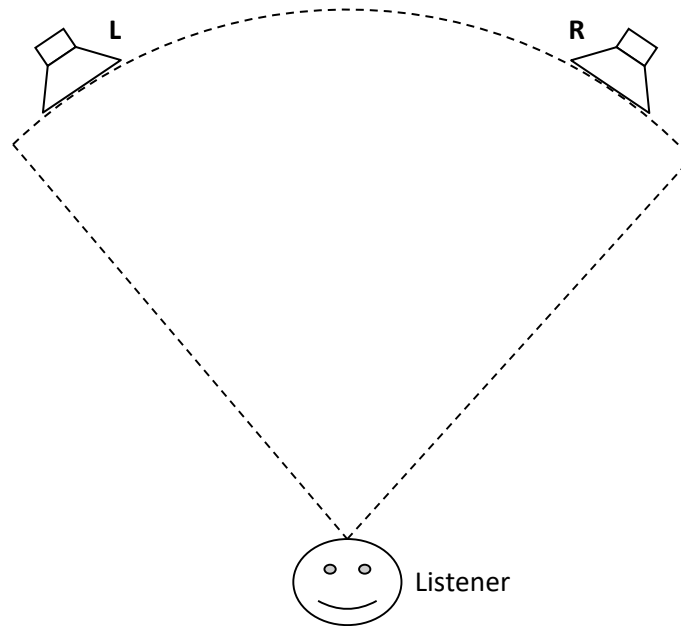
- Video conferencing
- **Spatial audio conferencing**
- Spatial audio + Video
- Immersive conferencing

Benefits of Spatial Audio

- Human's cocktail party effect
 - ▣ Selective attention
 - ▣ Only spend effort on comprehension
- Brain rejects incoherent signals at two ears
 - ▣ Reverberation & noise are disregarded (not for mono!)
- Benefits:
 - ▣ Memory, Focal Assurance, Perceived Comprehension, Listener's Preference
 - ▣ <http://msrweb/users/zhang/ThinkWeekPapers/Spatial%20audio%20conferencing.doc>

Multiparty Spatial Audio Conferencing

□ Virtual seating



Audio Spatialization

□ Delay and Gain Modulation

▣ Delay

$$\Delta_R = D - D \cos(\lambda(\Phi - \Theta))$$

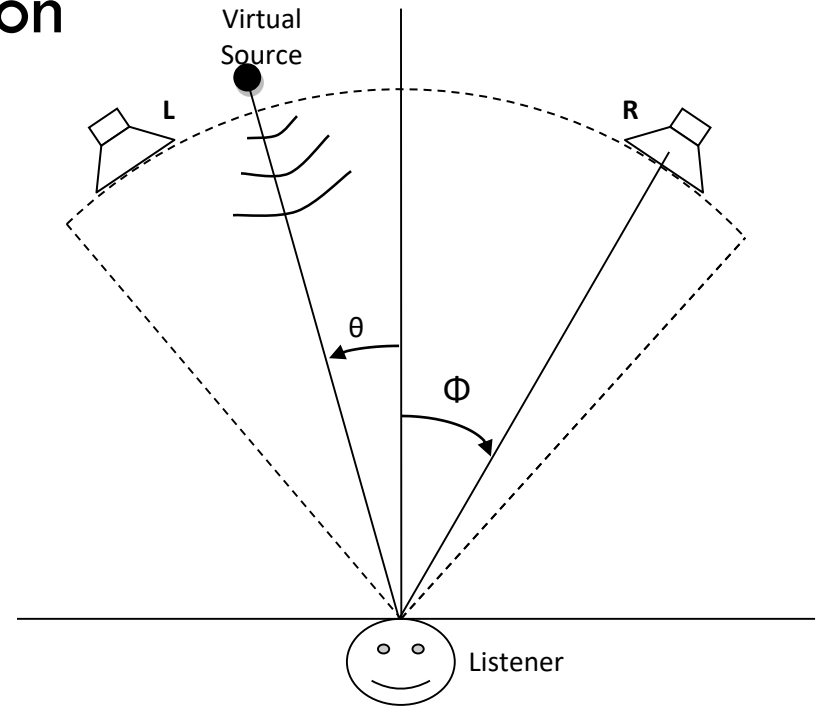
$$\Delta_L = D - D \cos(\lambda(\Phi + \Theta))$$

$$D=0.45\text{ms} \quad 1 \leq \lambda \leq \pi/(2\Phi)$$

▣ Gain

$$G_R = \cos(\lambda(\Phi - \Theta)/2)$$

$$G_L = \cos(\lambda(\Phi + \Theta)/2)$$



□ Example:

4 remote participants

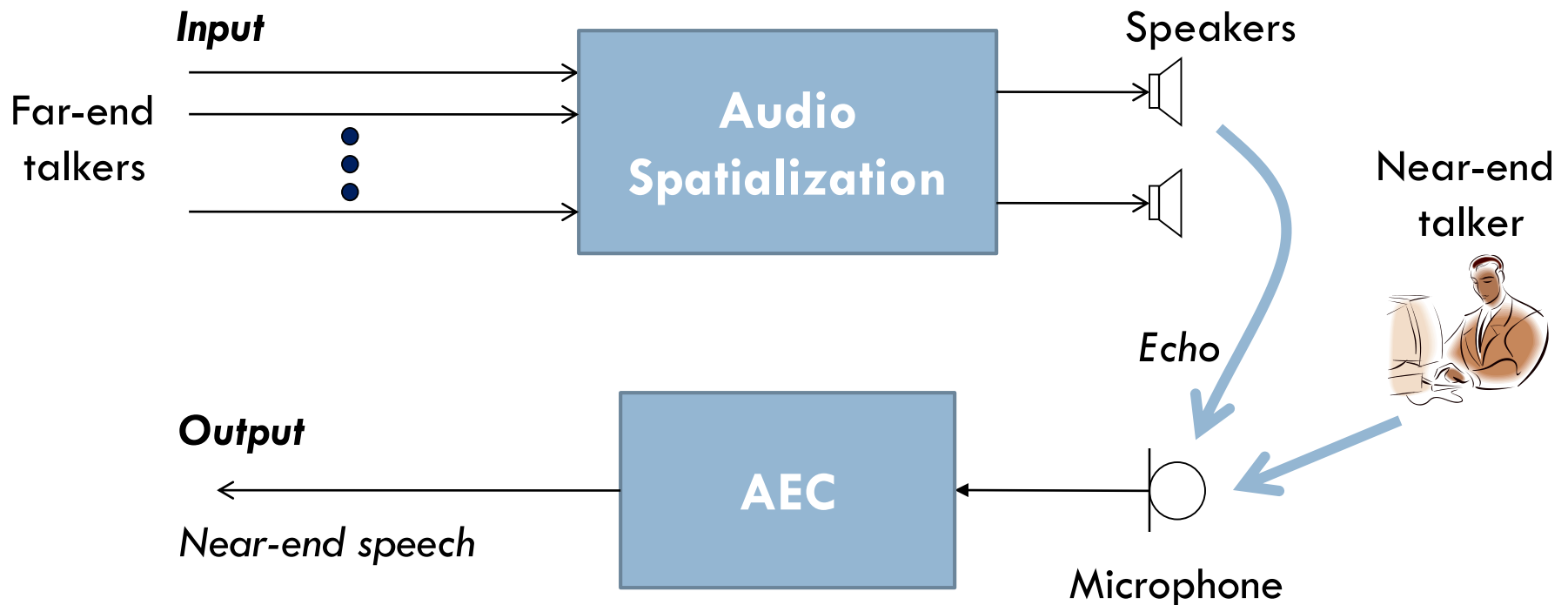
Traditional

(short)

Spatial

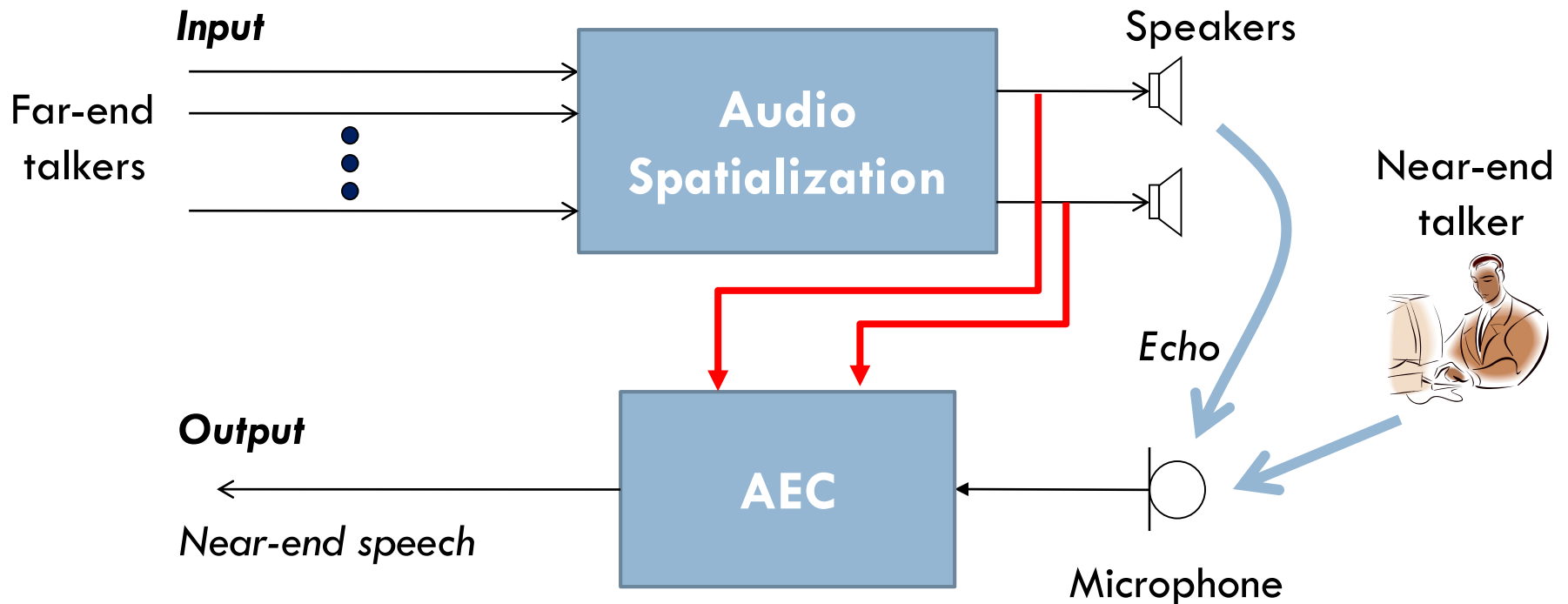
(short)

Multichannel AEC



□ **Question:** Which reference signals to use?

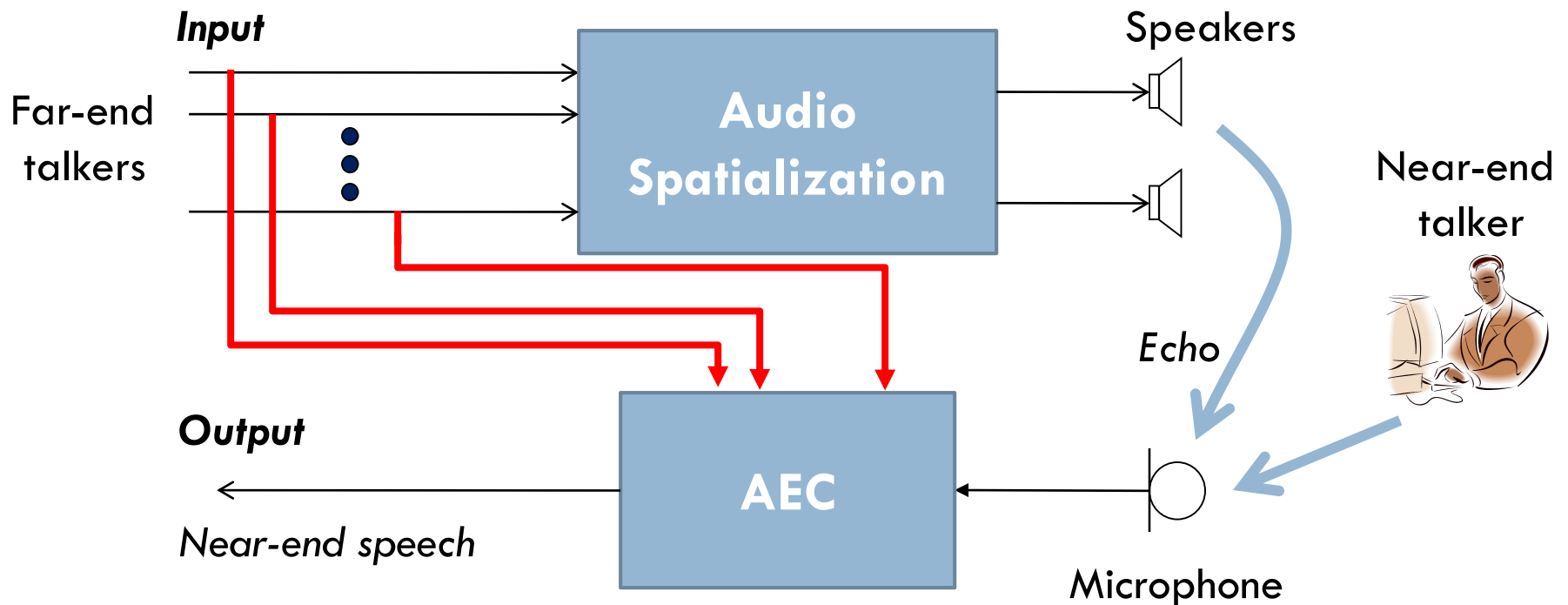
Approach 1: Use Speaker Signals



□ Possible problem:

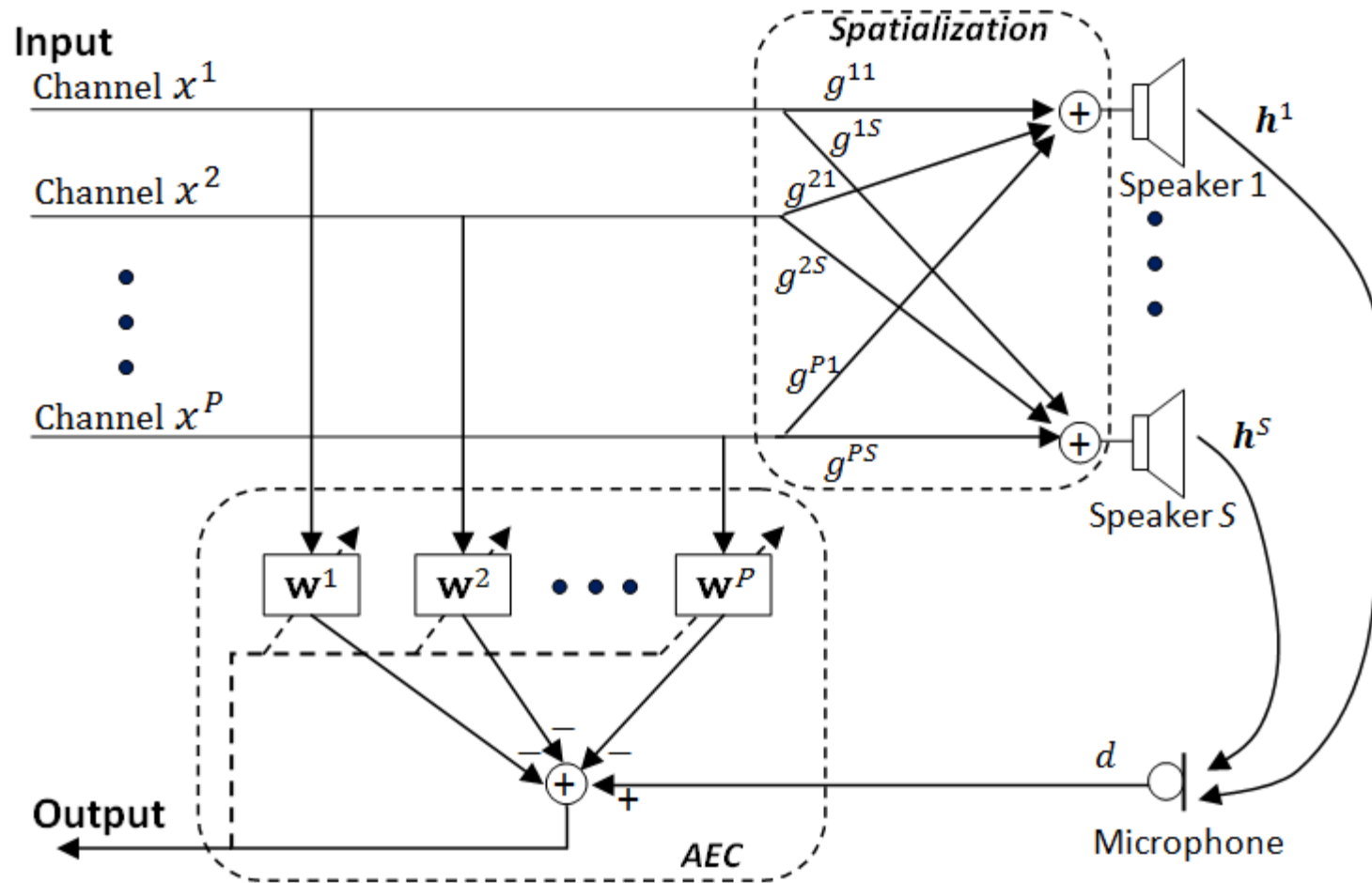
▣ Correlation between speaker signals

Approach 2: Use Far-End Channels



- Cancel each individual far-end speech
- **Our solution:** Constrained Kalman filtering

Multichannel AEC: Diagram



Problem Statement

□ Remote channels: $\{X^i \mid i = 1, \dots, P\}$

□ Spatialization on S speakers: $Y^s = \sum_{i=1}^P G^{is} X^i$

□ Speaker's room response: L -tap filter

$$\mathbf{H}_t^s = [H_t^s, H_{t-1}^s, \dots, H_{t-L-1}^s]^T$$

□ Microphone input: Echo

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

Problem Statement (cont'd)

- Determine the echo cancellers:
 - one per remote channel i : L -tap filter

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

such that echo is cancelled, i.e.,

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

- Constraint: \mathbf{W}_t^i 's are not mutually independent

$$D_t = \sum_{i=1}^P \sum_{s=1}^S G^{is} \mathbf{H}_t^{sT} \mathbf{X}_t^i$$



$$\mathbf{W}_t^i = \sum_{s=1}^S G^{is} \mathbf{H}_t^s$$

Constrained Kalman Filtering

- State Vector: Echo cancellers + Speaker RIR filters

$$\mathbf{s}_t = [\mathbf{W}_t^{1T}, \dots, \mathbf{W}_t^{PT}, \mathbf{H}_t^{1T}, \dots, \mathbf{H}_t^{ST}]^T$$

- System equation: $\mathbf{s}_t = \mathbf{s}_{t-1} + \mathbf{n}_t$

- Observation equation: $D_t = \mathbf{A}_t^T \mathbf{s}_t + v_t$ with $\mathbf{A}_t = [\mathbf{X}_t^{1T}, \dots, \mathbf{X}_t^{PT}, \mathbf{0}^{1T}, \dots, \mathbf{0}^{ST}]^T$

- Constraint: $\mathbf{C} \mathbf{s}_t = \mathbf{0}$ with $\mathbf{C} = \begin{bmatrix} -1 & & G^{11} \mathbf{1} & \dots & G^{1S} \mathbf{1} \\ & \ddots & \vdots & \ddots & \vdots \\ & & -1 & G^{P1} \mathbf{1} & \dots & G^{PS} \mathbf{1} \end{bmatrix}$

- New observation equation: observation + constraint

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

Constrained Kalman Filtering (cont'd)

□ Assumptions

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

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

→ Tuning parameter
to control how hard
the constraint be satisfied

□ Equations

$$\mathbf{S}_t^- = \mathbf{S}_{t-1}$$

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

$$\mathbf{K}_t = \mathbf{P}_t^- \mathbf{B}_t^H (\mathbf{B}_t \mathbf{P}_t^- \mathbf{B}_t^H + \mathbf{R}_t)^{-1}$$

$$\mathbf{S}_t = \mathbf{S}_t^- + \mathbf{K}_t (\mathbf{Y}_t - \mathbf{B}_t \mathbf{S}_t^-)$$

$$\mathbf{P}_t = (\mathbf{I} - \mathbf{K}_t \mathbf{B}_t) \mathbf{P}_t^-$$

Benefits of Constrained KF

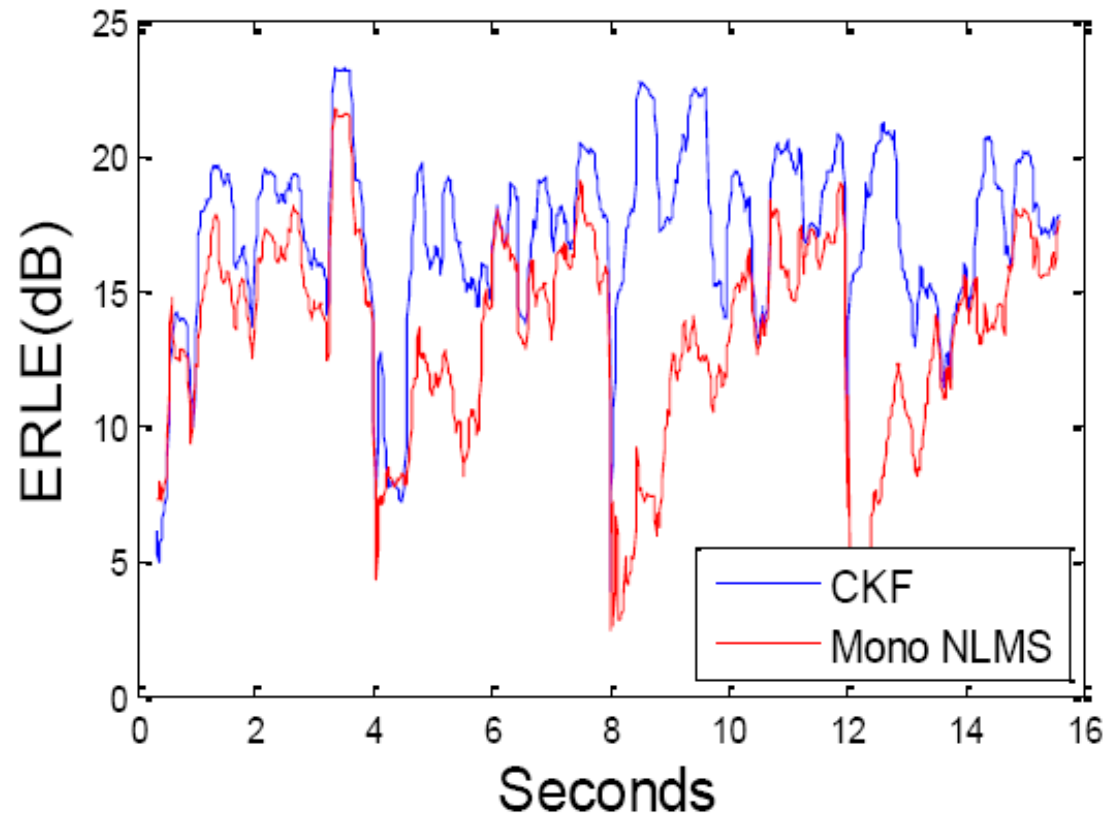
- The constraint is taken care of automatically, and can be imposed *with varying degrees*.
- All channels are taken into account simultaneously.
→ Overlapping far-end talking is not an issue
- The AEC for each channel is updated continuously because of the constraint, even if it is inactive.
→ AEC's are always up to date
- Ambient noise can be time varying.
→ Use a separate noise tracker

Comparison with Prior Art

- T.N. Yensen, R.A. Goubran, and I. Lambadaris, “Synthetic Stereo Acoustic Echo Cancellation Structure for Multiple Participant VoIP Conferences”, IEEE Transactions on Speech and Audio Processing, Vol. 9, No. 2, pp. 168-174, Feb. 2001.
- Same: One canceller per remote channel
- Differences:
 - ▣ Constrained vs. independent cancellers
 - Additional canceller is initiated before being active
 - A canceller is updated even if it is not active
 - ▣ Frequency vs. time domain
 - ▣ KF (RLS) vs. NLMS

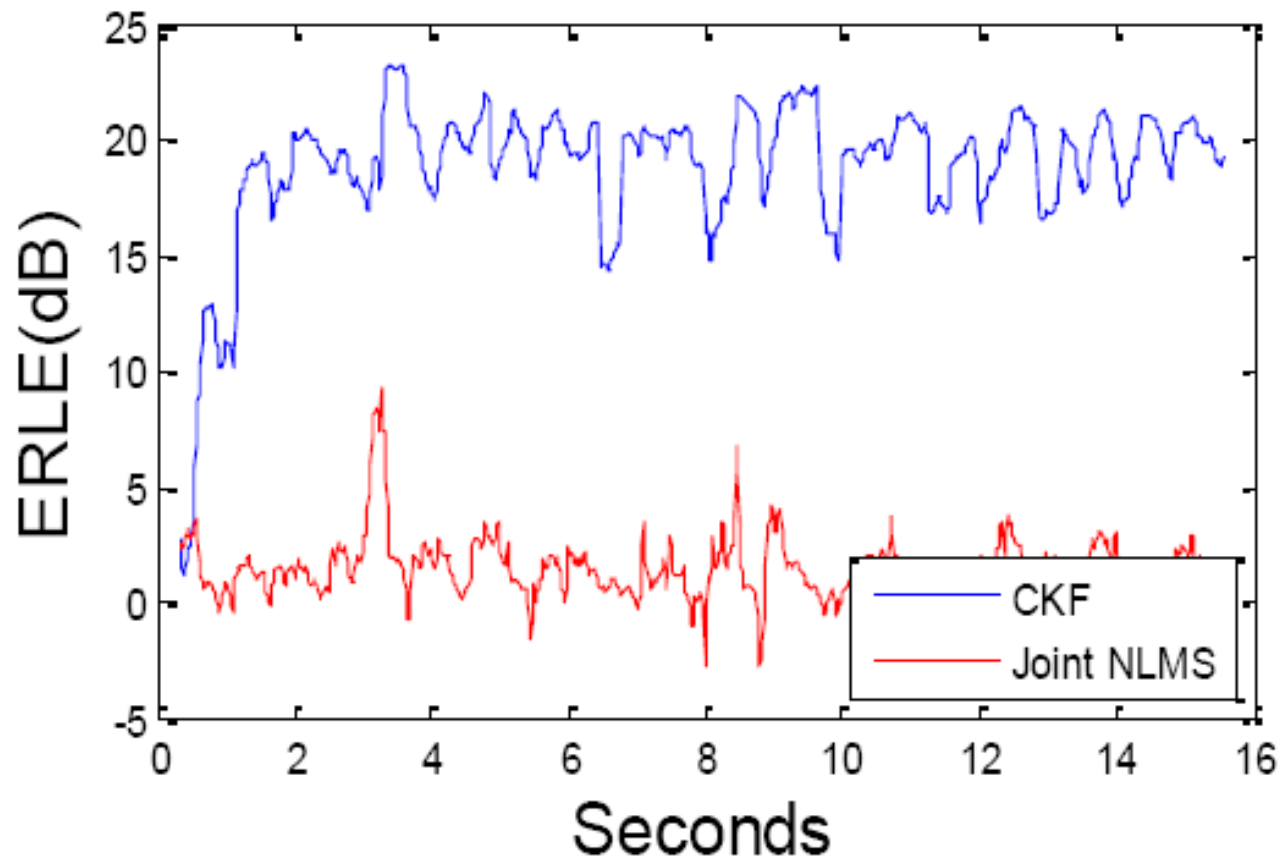
Experimental Results

- Simulation setup:
 - ▣ 4 remote talkers at $[-30^\circ, 30^\circ, 0^\circ, -45^\circ]$
 - ▣ Each talks for 4s
 - ▣ Noise: -20dB
 - ▣ Fixed RIR
- Comparison
 - ▣ Constrained KF
 - ▣ Multiple mono NLMS



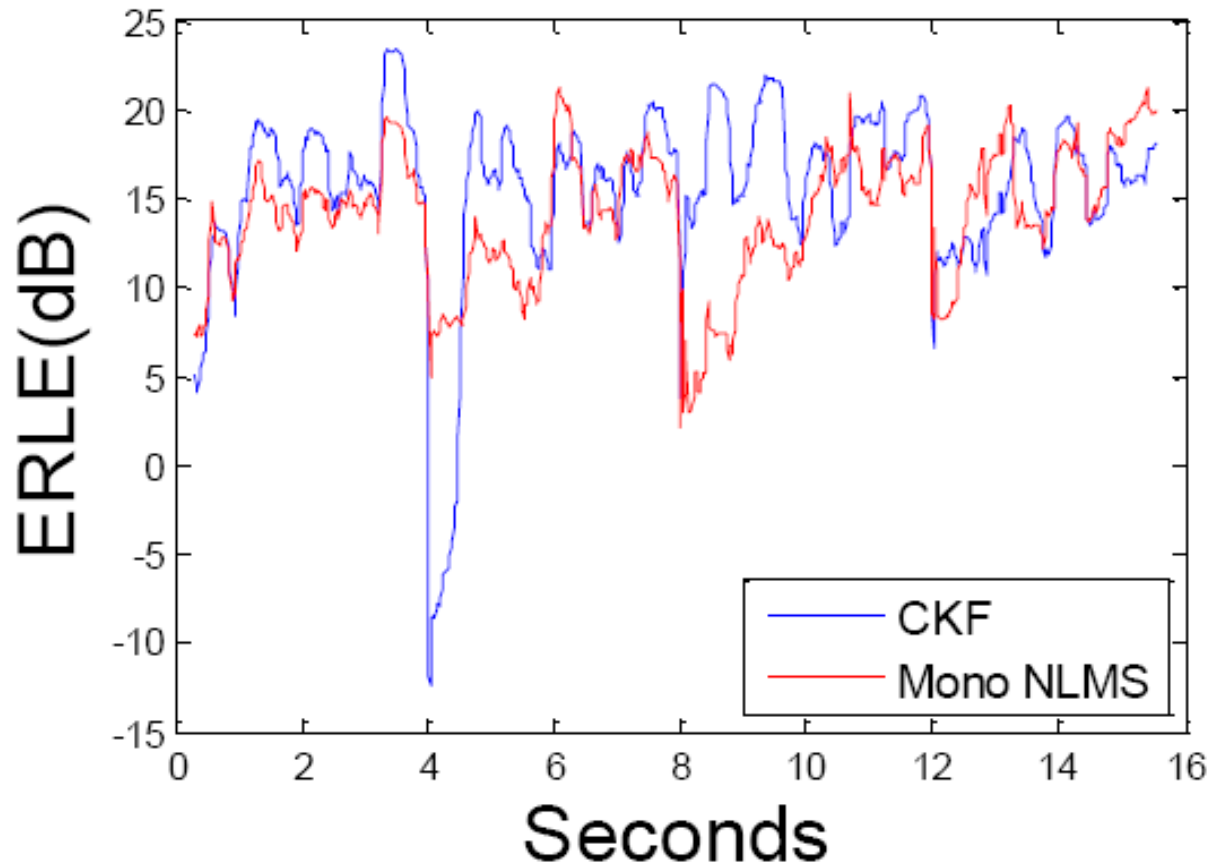
Experiment: Overlapping Talkers

- Two simultaneous remote talkers



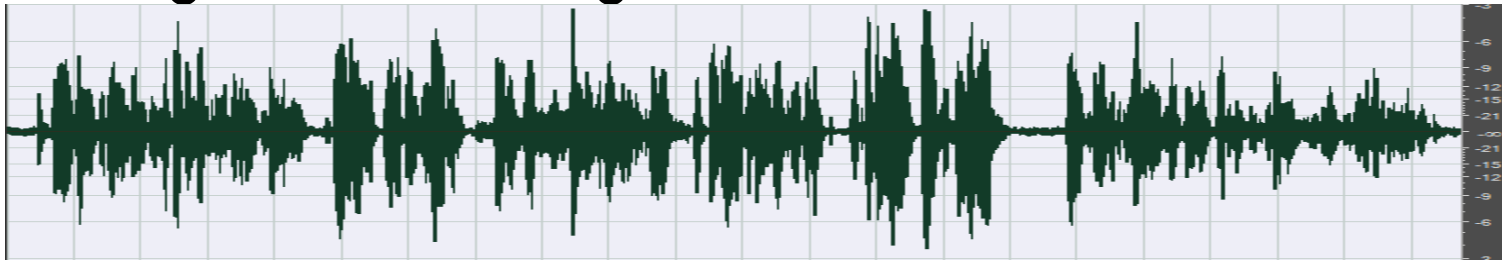
Experiment: Changing RIR

- -30dB change in RIR every 0.5 seconds



Experiment with real data

- Original recording with near-end talker



- AEC with multiple mono NLMSs



- AEC with CKF



Conclusions

- Constrained Kalman Filtering for multichannel AEC
- Outperform over multiple independent mono AECs
 - Additional canceller is initiated before being active
 - A canceller is updated even if it is not active
- Naturally works with multiple simultaneous remote talkers without resort to channel switching

Thank you !

Q & A