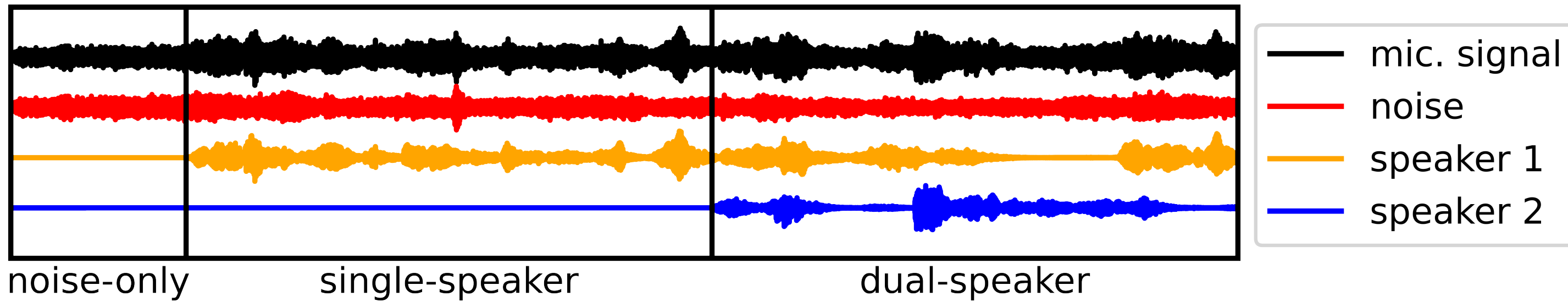


Introduction

- Common multi-microphone noise reduction methods, e.g., linearly constrained minimum variance (LCMV) beamforming, rely on estimates of the relative transfer function (RTF) vectors of all speakers
- In this work: acoustic scenario with two successively activating speakers



- Objective:** Estimate RTF vector of second speaker during overlapping speech segments

MAIN IDEAS

Covariance Blocking and Whitening (CBW) method for estimating the RTF vector of the second speaker

- Block the initial speaker** to isolate information about the second speaker
- Whiten the noise** to minimize its influence

Problem Statement

Signal Model in STFT-domain

- Two speakers and noise recorded with M microphones
- Noisy covariance matrix

$$\mathbf{R}_y = \underbrace{\mathbf{h}\phi_x\mathbf{h}^H}_{\text{speaker 2}} + \underbrace{\mathbf{g}\phi_u\mathbf{g}^H}_{\text{speaker 1}} + \mathbf{R}_n \in \mathbb{C}^{M \times M}$$

- Noise covariance matrix \mathbf{R}_n and RTF vector \mathbf{g} of **speaker 1** can be estimated in noise-only and single-speaker segments
- Power spectral densities ϕ_u and ϕ_x are unknown and time-varying

GOAL

Estimate RTF vector \mathbf{h} of **speaker 2** in dual-speaker segment using noisy covariance matrix \mathbf{R}_y and estimates of \mathbf{R}_n and \mathbf{g}

Conventional Methods

1.) Covariance Whitening (CWu) [1]

- Jointly whiten **speaker 1** and **noise** with undesired covariance matrix \mathbf{R}_v

$$\tilde{\mathbf{h}}^{(CW)} = \tilde{\mathbf{h}}/e_r^T \tilde{\mathbf{h}} \quad \text{with} \quad \tilde{\mathbf{h}} = \mathbf{R}_v^{H/2} \mathcal{P}\{\mathbf{R}_v^{-H/2} \mathbf{R}_y \mathbf{R}_v^{-1/2}\}$$

with $\mathcal{P}\{\cdot\}$ denoting principal eigenvector

2.) Blind Oblique Projection (BOP) [2]

- Noise is neglected** by assuming a sufficiently high SNR
- Block speaker 2** using parameterized oblique projection matrix $\mathbf{P}_{g\theta}^\perp$ while **keeping speaker 1 distortionless** and minimizing the power

$$\tilde{\mathbf{h}}^{(BOP)} = \tilde{\mathbf{h}}/e_r^T \tilde{\mathbf{h}} \quad \text{with} \quad \tilde{\mathbf{h}} = \arg\min(\text{Tr}\{\mathbf{P}_{g\theta}^\perp \mathbf{R}_y \mathbf{P}_{g\theta}^{\perp H}\}) \quad \text{with} \quad \mathbf{P}_{g\theta}^\perp = \mathbf{g}(\mathbf{P}_\theta^\perp \mathbf{g})^+$$

Proposed Method

Covariance Blocking and Whitening (CBW)

- Block speaker 1** using orthogonal projection matrix $\mathbf{P}_g^\perp = \mathbf{I}_M - \frac{\mathbf{g}\mathbf{g}^H}{(\mathbf{g}^H\mathbf{g})}$
 - Noise whitening requires full column rank \rightarrow remove one column

$$\mathbf{R}_y \mathbf{P}_{g,r}^\perp = \mathbf{h}\phi_x\mathbf{h}^H \mathbf{P}_{g,r}^\perp + \mathbf{R}_n \mathbf{P}_{g,r}^\perp \in \mathbb{C}^{M \times M-1}$$

- Whiten the noise** using pseudo-inverse of blocked noise covariance matrix

$$(\mathbf{R}_n \mathbf{P}_{g,r}^\perp)^+ \mathbf{R}_y \mathbf{P}_{g,r}^\perp - \mathbf{I}_{M-1} = \underbrace{(\mathbf{R}_n \mathbf{P}_{g,r}^\perp)^+}_{\propto \mathbf{q}_L} \mathbf{h}\phi_x \underbrace{\mathbf{h}^H \mathbf{P}_{g,r}^\perp}_{\propto \mathbf{q}_r^H} \in \mathbb{C}^{M-1 \times M-1}$$

Covariance Blocking and Whitening (CBW)

- Set up non-linear equation system using left & right principal singular vectors \mathbf{q}_L and \mathbf{q}_R and unknown scaling factor α

$$\begin{bmatrix} \mathbf{q}_L \\ \mathbf{q}_R \alpha \end{bmatrix} \stackrel{\perp}{=} \mathbf{B} \tilde{\mathbf{h}} \quad \text{with} \quad \mathbf{B} = \begin{bmatrix} (\mathbf{R}_n \mathbf{P}_{g,r}^\perp)^+ \\ (\mathbf{P}_{g,r}^\perp)^H \end{bmatrix}$$

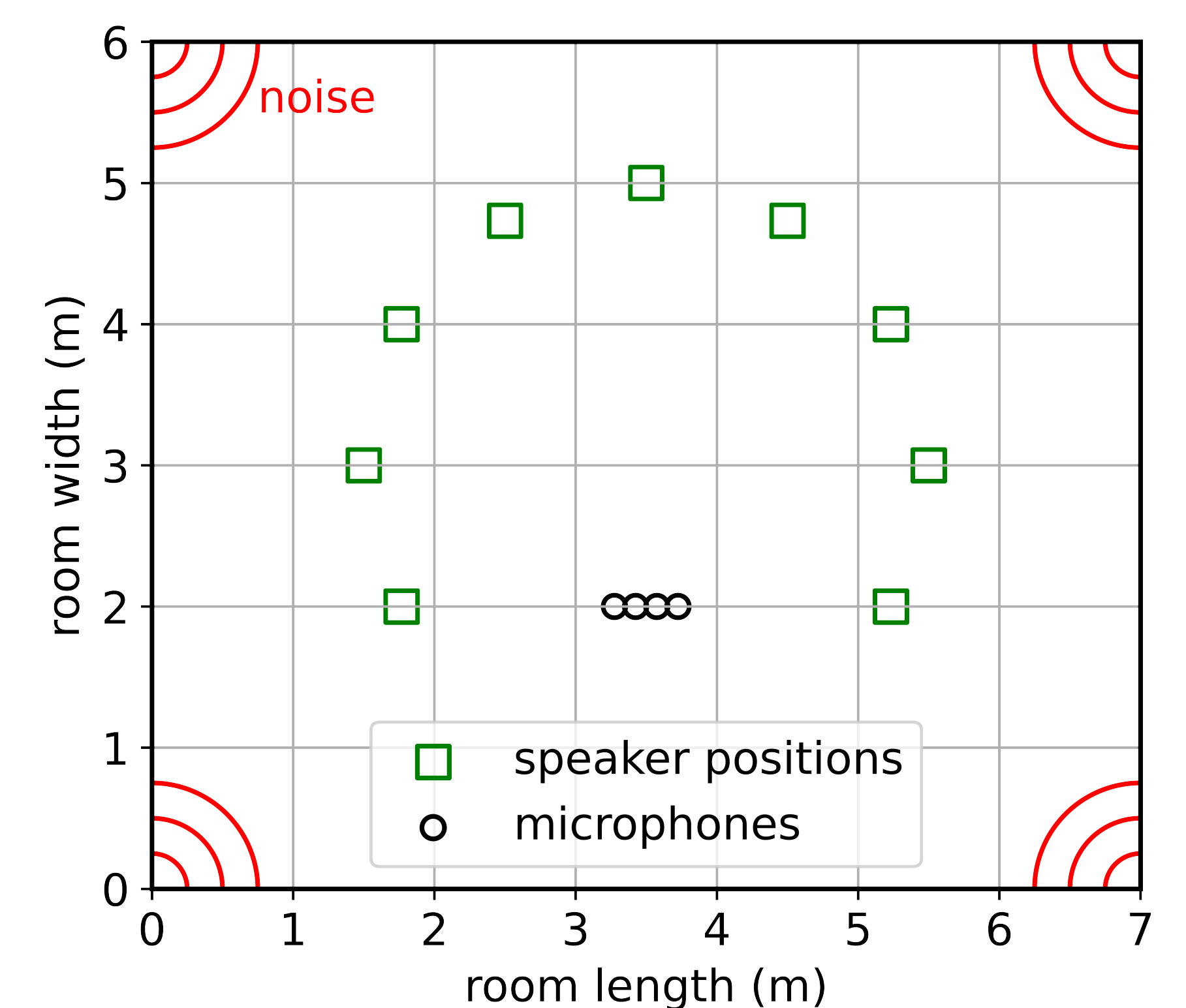
$$\tilde{\mathbf{h}}^{(CBW)} = \tilde{\mathbf{h}}/e_r^T \tilde{\mathbf{h}} \quad \text{with} \quad \tilde{\mathbf{h}} = \mathbf{B}^+ \begin{bmatrix} \mathbf{q}_L \\ -\mathbf{q}_R (\mathbf{P}_B^{\perp,R} \mathbf{q}_R)^+ \mathbf{P}_B^{\perp,L} \mathbf{q}_L \end{bmatrix}$$

Method Overview

	CWu [1]	BOP [2]	CBW (prop.)
1.) Blocking	X	speaker 2	speaker 1
2.) Whitening	speaker 1 & noise	X	noise
required estimates	\mathbf{R}_v	\mathbf{g}	\mathbf{R}_n & \mathbf{g}

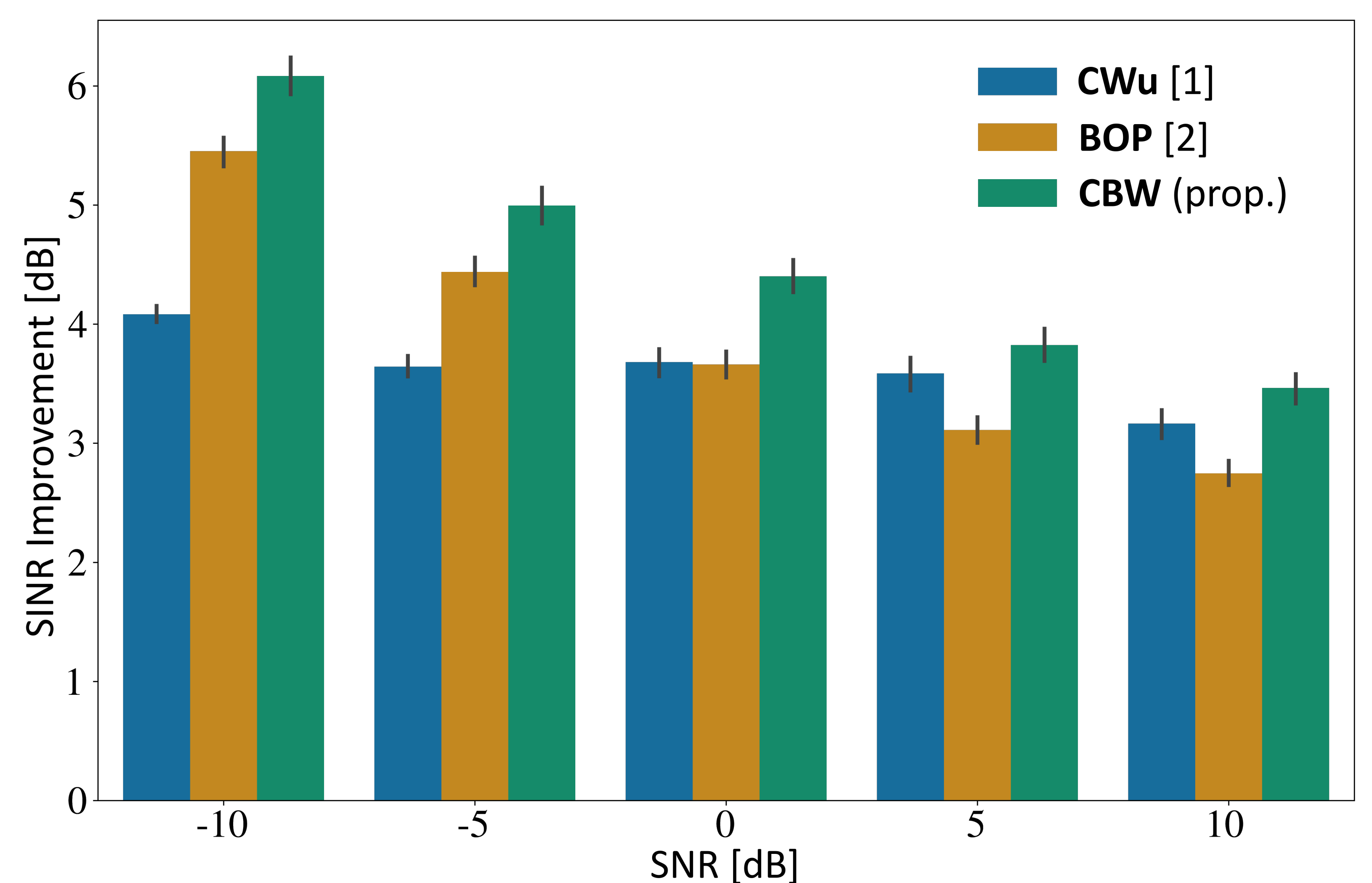
Evaluation

- Linear array with $M = 4$ microphones ($d = 2$ cm)
- Clean signals convolved with measured room impulse responses ($T_{60} \approx 500$ ms)
- 72 combinations of dual-speaker positions
- Quasi-diffuse babble noise with Signal-to-noise ratio (SNR): $-10 : 5 : 10$ dB
- Signal-to-interferer ratio (SIR): $-10 : 5 : 10$ dB
- $f_s = 16$ kHz
- STFT framework: frame length 200 ms, 75% overlap
- segment borders are assumed to be known



Results

- Signal-to-interferer-and-noise ratio (SINR) improvement of LCMV beamformer using estimated noise covariance matrix and RTF vectors of both speakers



Conclusions

- The proposed CBW method combines blocking of the initial speaker and whitening of the noise to estimate the RTF vector of the second speaker
- In terms of SINR improvement, the proposed CBW method outperforms conventional RTF vector estimation methods

References

[1] E. Worsitz and R. Haeb-Umbach, "Blind acoustic beamforming based on generalized eigenvalue decomposition," *IEEE Trans. Audio, Speech, and Language Processing*, vol. 15, no. 5, pp. 1529–1539, 2007.

[2] D. Cherkassky and S. Gannot, "Successive Relative Transfer Function Identification Using Blind Oblique Projection," *IEEE/ACM Trans. Audio, Speech, and Language Processing*, vol. 28, pp. 474–486, 2020.