

PROBLEM STATEMENT

- microphone signals corrupted by **reverberation and diffuse noise**
- multichannel Wiener filter (MWF) requires estimates of **relative early transfer functions (RETFs)** and **diffuse power spectral density (PSD)**
- many such estimators are **decoupled**
- goal: develop a **joint RETF and diffuse PSD estimator** under assumption of diffuse noise

SIGNAL MODEL

- microphone signal model in STFT-domain:

$$\mathbf{y}(l) = \mathbf{x}(l) + \mathbf{d}(l) \text{ with } \mathbf{y}(l) = \begin{bmatrix} Y_1(l) \\ \vdots \\ Y_M(l) \end{bmatrix}$$

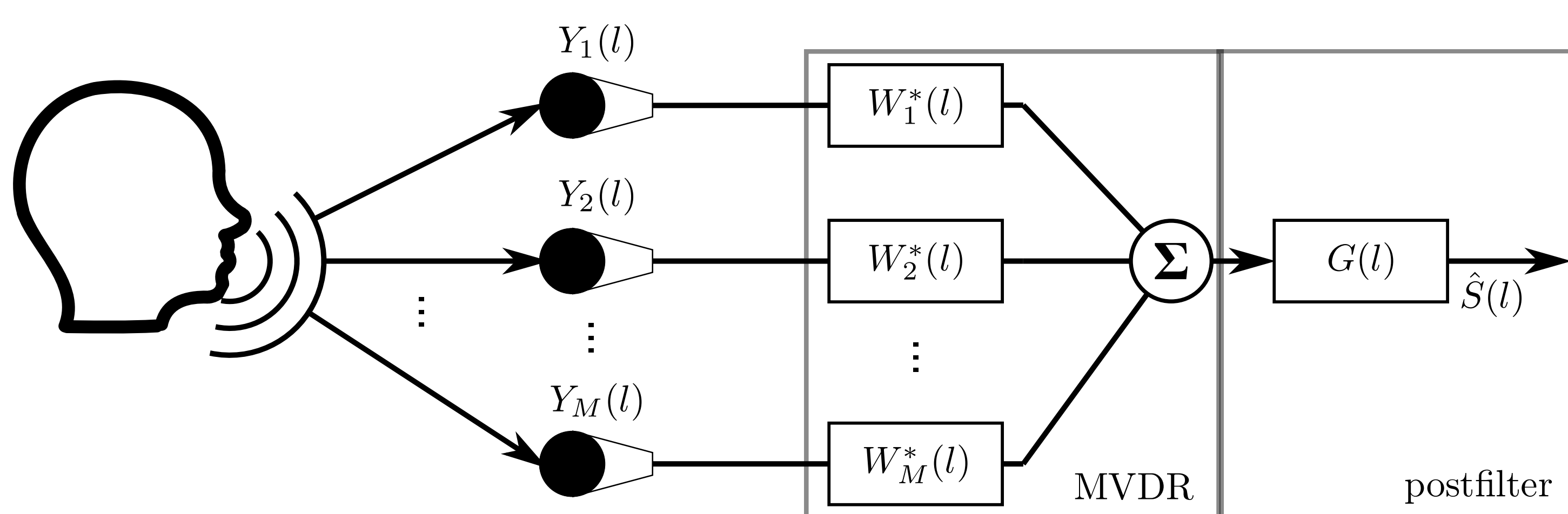
- M microphones, time frame l
- independent processing in each sub-band
- $\mathbf{x}(l)$: direct and early speech component (target)
- $\mathbf{d}(l)$: diffuse noise and reverberation
- assumptions: $\mathbf{x}(l)$ and $\mathbf{d}(l)$ uncorrelated, stationary diffuse noise

⇒ microphone PSD matrix:

$$\Phi_{\mathbf{y}}(l) = \mathbb{E}\{\mathbf{y}(l)\mathbf{y}^H(l)\} = \phi_s(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\Gamma$$

- $\mathbf{a}(l)$: RETF vector, Γ : diffuse coherence matrix
- $\phi_s(l), \phi_d(l)$: target and diffuse PSD

MULTICHANNEL WIENER FILTER



- two stages: minimum variance distortionless response (MVDR) beamformer and **single-channel postfilter**

$$\mathbf{w}_{\text{MWF}}(l) = \underbrace{\left(\frac{\Gamma^{-1}\mathbf{a}(l)}{\mathbf{a}^H(l)\Gamma^{-1}\mathbf{a}(l)} \right)^H}_{\mathbf{w}_{\text{MVDR}}(l)} \underbrace{\frac{\phi_s(l)}{\phi_s(l) + \phi_d(l) / (\mathbf{a}^H(l)\Gamma^{-1}\mathbf{a}(l))}}_{G(l)}$$

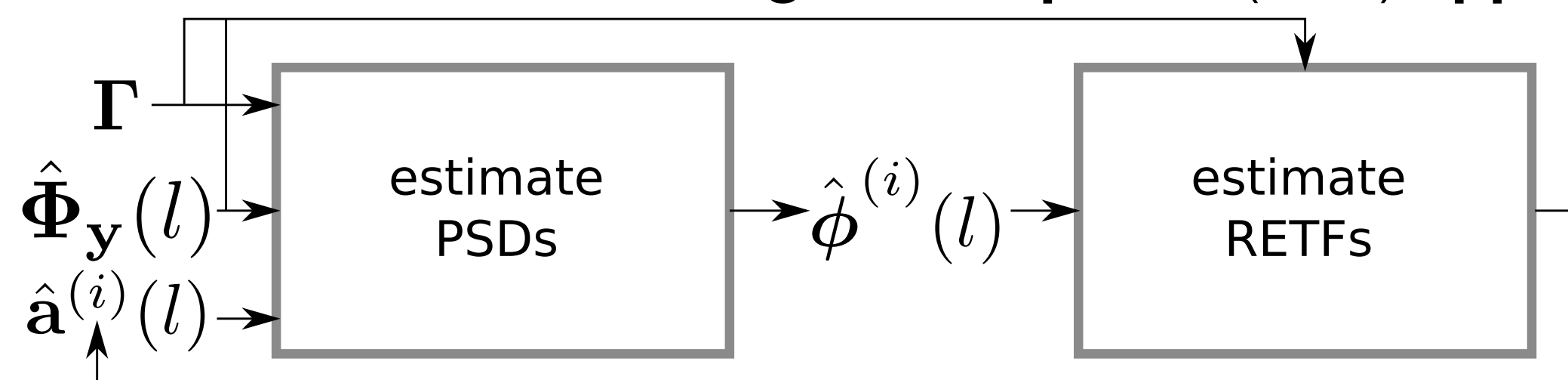
- requires estimates of PSDs $[\phi_s(l), \phi_d(l)]^T =: \phi(l)$ and RETFs $\mathbf{a}(l)$

PROPOSED METHOD

- proposition: based on the **same cost function, jointly estimate PSDs and RETFs**:

$$(\hat{\mathbf{a}}_{\text{LS}}, \hat{\phi}_{\text{LS}}) = \underset{\mathbf{a}, \phi}{\operatorname{argmin}} \left\| \hat{\Phi}_{\mathbf{y}} - (\phi_s \mathbf{a} \mathbf{a}^H + \phi_d \Gamma) \right\|_F^2 =: \underset{\mathbf{a}, \phi}{\operatorname{argmin}} \|\mathbf{E}\|_F^2$$

- no closed-form solution ⇒ **alternating least-squares (ALS) approach**



- initialize RETFs $\hat{\mathbf{a}}_{\text{LS}}^{(0)}$ (e.g., using DOA); set $i = 1$
- estimate PSDs $\hat{\phi}_{\text{LS}}^{(i)}$ using (1), assuming $\mathbf{a} = \hat{\mathbf{a}}_{\text{LS}}^{(i-1)}$
- estimate RETFs $\hat{\mathbf{a}}_{\text{LS}}^{(i)}$ using (2), assuming $\phi = \hat{\phi}_{\text{LS}}^{(i)}$, set $i = i + 1$
- repeat steps (ii) — (iii) until convergence ('ALS iteration')

ESTIMATING PSDs AND RETFs

- PSDs

- assume $\mathbf{a} = \hat{\mathbf{a}}_{\text{LS}}^{(i-1)} \Rightarrow$ estimate PSDs by minimizing Frobenius norm of error matrix $\mathbf{E}^{(i)} = \hat{\Phi}_{\mathbf{y}} - \left(\hat{\phi}_{\text{s,LS}}^{(i)} \hat{\mathbf{a}}_{\text{LS}}^{(i-1)} \hat{\mathbf{a}}_{\text{LS}}^{(i-1)H} + \hat{\phi}_{\text{d,LS}}^{(i)} \Gamma \right)$ [1]:

$$\hat{\phi}_{\text{LS}}^{(i)} = \underset{\phi}{\operatorname{argmin}} \|\mathbf{E}^{(i)}\|_F^2 = \mathbf{A}^{-1} \mathbf{b}, \text{ where } \mathbf{A} = \begin{bmatrix} (\hat{\mathbf{a}}^H \hat{\mathbf{a}})^2 & \hat{\mathbf{a}}^H \Gamma \hat{\mathbf{a}} \\ \hat{\mathbf{a}}^H \Gamma \hat{\mathbf{a}} & \operatorname{trace}\{\Gamma^H \Gamma\} \end{bmatrix}, \mathbf{b} = \begin{bmatrix} \operatorname{Re}\{\hat{\mathbf{a}}^H \hat{\Phi}_{\mathbf{y}} \hat{\mathbf{a}}\} \\ \operatorname{Re}\{\operatorname{trace}\{\hat{\Phi}_{\mathbf{y}} \Gamma^H\}\} \end{bmatrix} \quad (1)$$

- constrain PSDs to valid range: $0 \leq \{\hat{\phi}_{\text{s,LS}}^{(i)}, \hat{\phi}_{\text{d,LS}}^{(i)}\} \leq \frac{1}{M} \mathbf{y}^H \mathbf{y}$

- RETFs

- assume $\phi = \hat{\phi}_{\text{LS}}^{(i)}$

- identify $\hat{\Phi}_{\mathbf{y}} - \hat{\phi}_{\text{d,LS}}^{(i)} \Gamma =: \hat{\Phi}_{\mathbf{x}}^{(i)}$

$$\Rightarrow \hat{\mathbf{a}}_{\text{LS}}^{(i)} = \underset{\mathbf{a}}{\operatorname{argmin}} \left\| \hat{\Phi}_{\mathbf{x}}^{(i)} - \hat{\phi}_{\text{s,LS}}^{(i)} \mathbf{a} \mathbf{a}^H \right\|_F^2 \triangleq \text{scaled rank-1 approx. of } \hat{\Phi}_{\mathbf{x}}^{(i)} \quad [2]$$

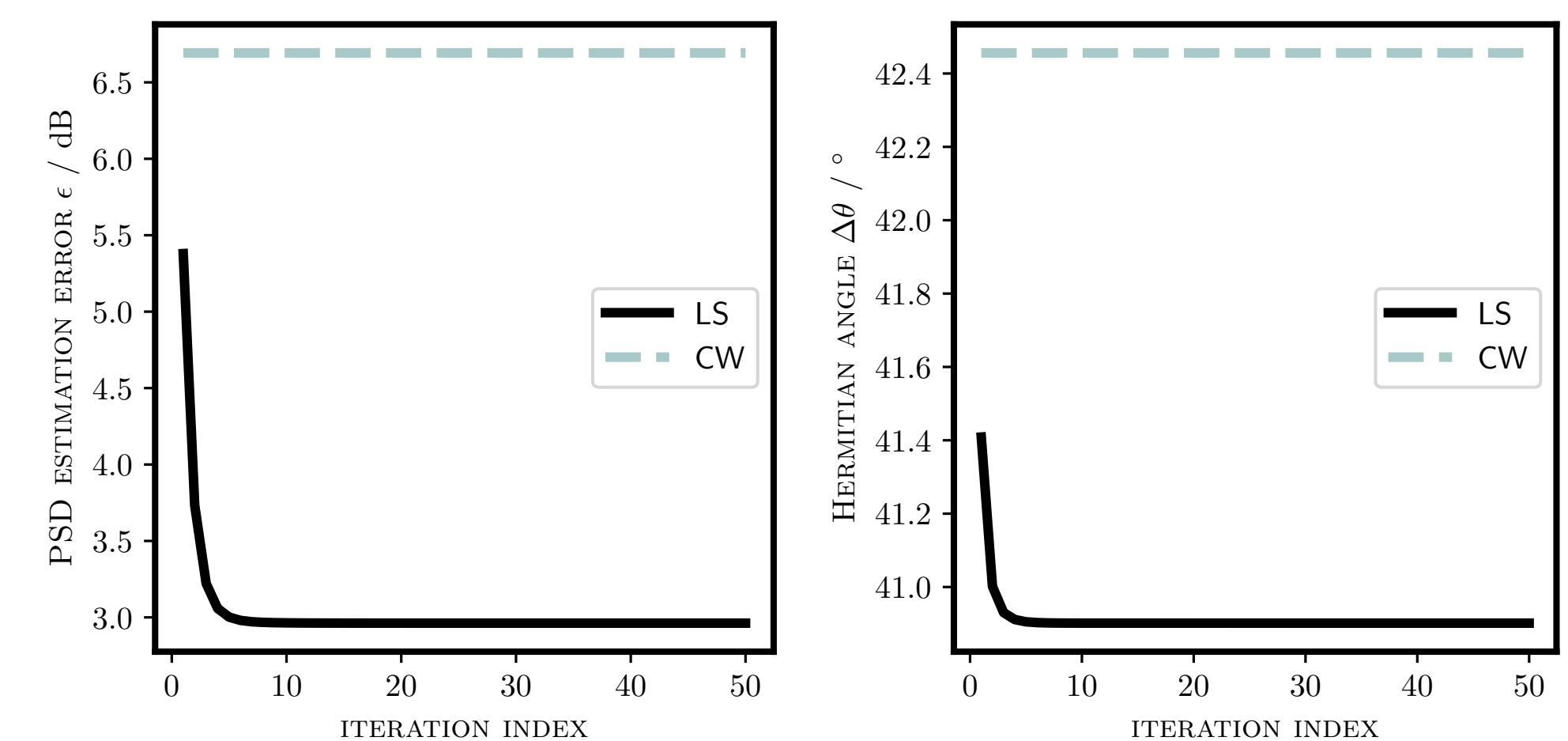
$$= \sqrt{\frac{\lambda_1}{\hat{\phi}_{\text{s,LS}}^{(i)}}} \mathbf{u}_1 \quad (2)$$

SIMULATIONS ON ARTIFICIAL DATA

- generated according to signal model ⇒ oracle information available
- test robustness to model errors: additional uncorrelated white noise
- evaluation measures: average diffuse PSD mismatch and Hermitian angle between estimate and oracle
- baseline: EVD-based PSD estimator combined with covariance whitening-based RETF estimator (denoted CW) [3]
- mismatches vs. number of iterations

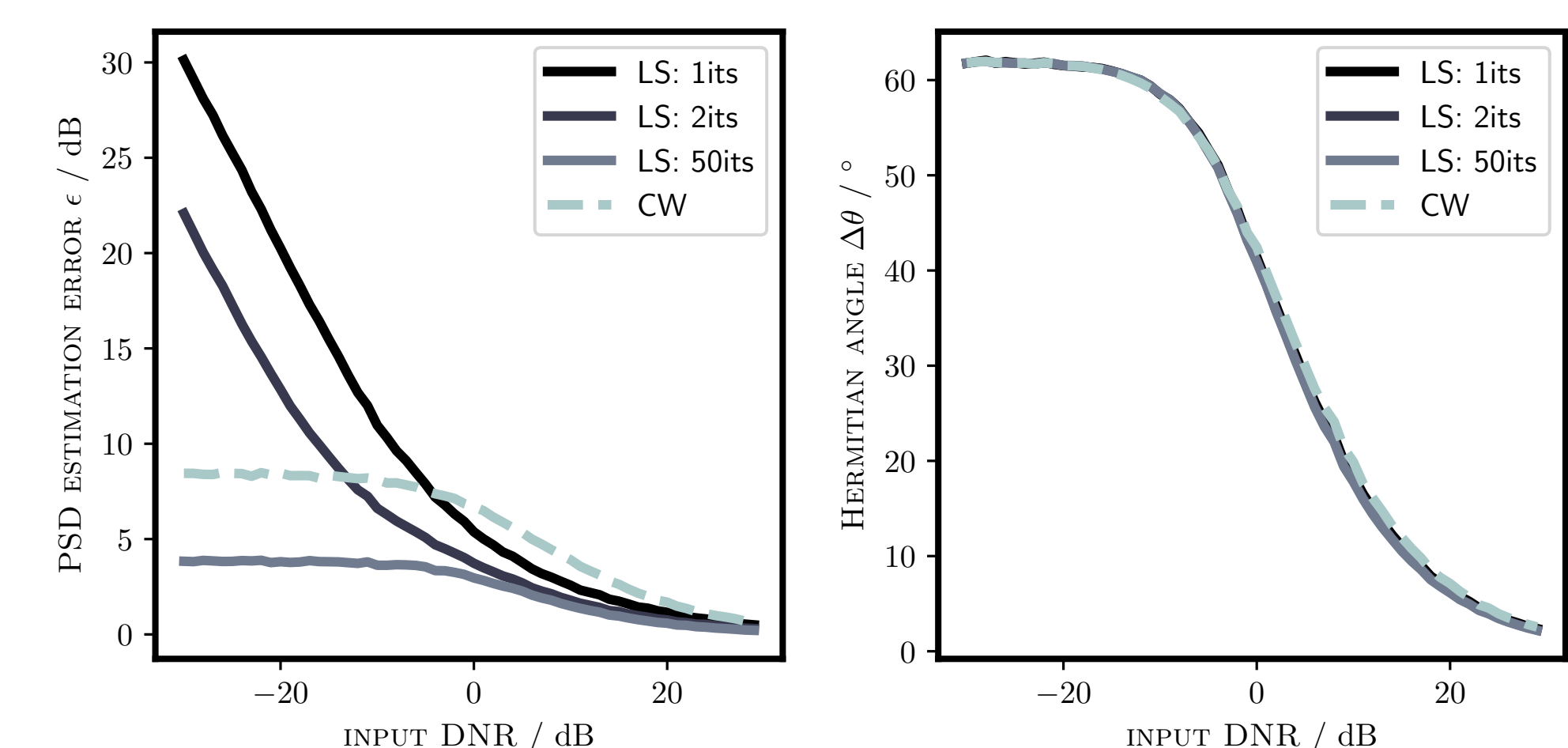
more ALS iterations:

- no improved RETF estimate
- significantly smaller diffuse PSD mismatch



- mismatches vs. input DNR

- similar RETF estimation accuracy
- after convergence: overperformance at low DNR, same at high DNR



SIMULATIONS ON MEASURED DATA

- $M = 4$ microphones, $f_s = 16$ kHz
- diffuse babble noise [4] at 10 dB SNR wrt. reference microphone

	array geometry	mic. distance	θ	T_{60}
AS ₁	linear	$d = 8$ cm	45°	0.61 s
AS ₂	circular	$r = 10$ cm	45°	0.73 s
AS ₃	linear	$d = 6$ cm	-15°	1.25 s

- additional **spatially uncorrelated white noise** at $\{0, 10, 20, 30\}$ dB DNR
- STFT: 64 ms frame length ($N_{\text{FFT}} = 1024$), 16 ms shift
- $\Phi_{\mathbf{y}}(l)$ estimated using recursive averaging, 40 ms smoothing constant
- target PSD $\phi_s(l)$ estimated using **decision-directed approach** [5] (LS estimate not used in postfilter)
- performance measures: **frequency-weighted segmental SNR, PESQ**

uncorrelated noise DNR / dB	ΔPESQ				ΔfwsSNR			
	0	10	20	30	0	10	20	30
CW	0.0547	0.1632	-0.0561	-0.1037	1.7112	2.2500	3.6697	3.4651
LS	0.3823	0.5264	0.1818	0.2331	4.5457	4.9639	4.7897	3.7432

References

- O. Schwartz, S. Gannot, and E. A. P. Habets, Joint estimation of late reverberant and speech power spectral densities in noisy environments using Frobenius norm. In *Proc. European Signal Processing Conference*, pages 1123–1127, Budapest, Hungary, September 2016.
- C. Eckart and G. Young, The approximation of one matrix by another of lower rank. *Psychometrika*, 1(3):211–218, 1936.
- I. Kodrasi and S. Doclo, EVD-based multi-channel dereverberation of a moving speaker using different RETF estimation methods. In *Proc. Joint Workshop on Hands-Free Speech Communication and Microphone Arrays*, pages 116–120, San Francisco, USA, March 2017.

[4] E. A. P. Habets, I. Cohen, and S. Gannot, Generating nonstationary multisensor signals under a spatial coherence constraint. *Journal of the Acoustical Society of America*, 124(5):2911–2917, November 2008.

[5] Y. Ephraim and D. Malah, Speech enhancement using a minimum mean-square error short-time spectral amplitude estimator. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 32(6):1109–1121, December 1984.