

PROBLEM STATEMENT

- microphone signals corrupted by **reverberation and noise**
- multi-channel Wiener filter (MWF) requires PSD estimates
- recently proposed **diffuse PSD estimator based on eigenvalue decomposition (EVD)**, which does not require relative early transfer function (RETF) vector of desired speech source [1]
- goal:** reduce **computational cost** of EVD-based PSD estimator

SIGNAL MODEL

- microphone signal model in STFT-domain, independent processing in each subband

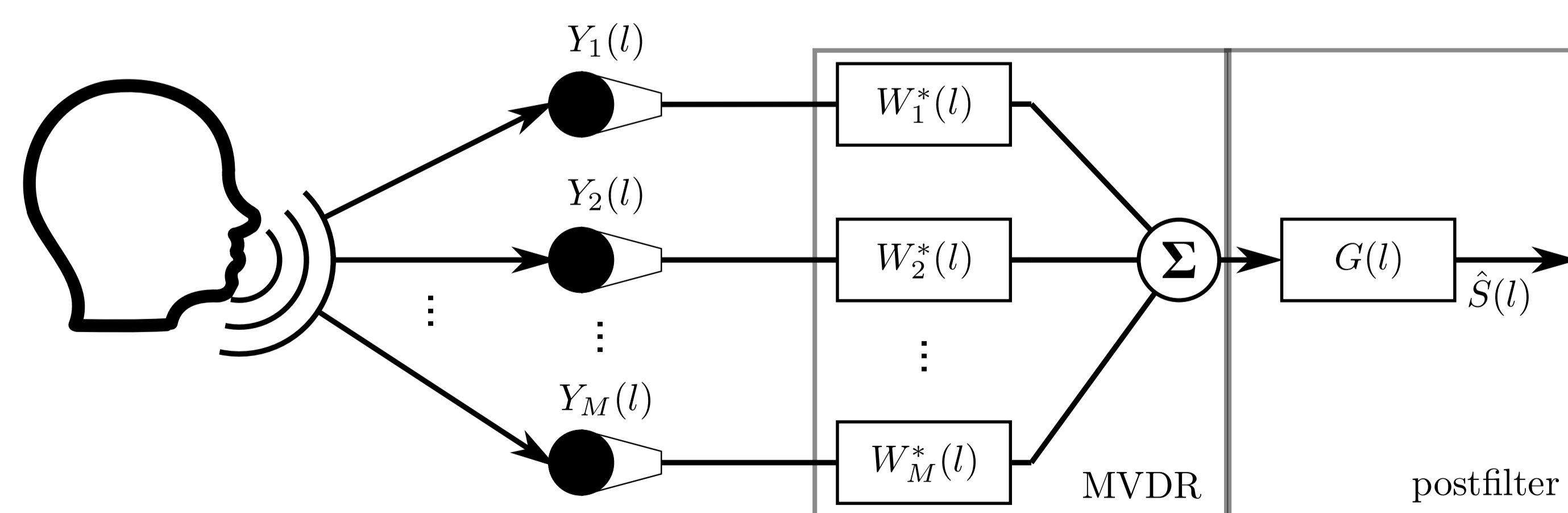
$$\mathbf{y}(l) = \mathbf{x}(l) + \mathbf{d}(l) \text{ with } \mathbf{y}(l) = \begin{bmatrix} Y_1(l) \\ \vdots \\ Y_M(l) \end{bmatrix}$$

- M microphones, time frame l
- $\mathbf{x}(l)$: direct and early speech component
- $\mathbf{d}(l)$: diffuse noise and reverberation
- microphone PSD matrix

$$\Phi_{\mathbf{y}}(l) = \mathcal{E}\{\mathbf{y}(l)\mathbf{y}^H(l)\} = \Phi_s(l)\mathbf{a}(l)\mathbf{a}^H(l) + \Phi_d(l)\Gamma$$

- $\mathbf{a}(l)$: RETF vector, Γ : diffuse coherence matrix
- Φ_s, Φ_d : speech and diffuse PSD

MULTICHANNEL WIENER FILTER



- stage I: MVDR beamformer $\mathbf{w}(l) = \frac{\Gamma^{-1}\mathbf{a}(l)}{\mathbf{a}^H(l)\Gamma^{-1}\mathbf{a}(l)}$
- stage II: spectro-temporal postfilter $G(l) = \frac{\hat{\Phi}_s(l)}{\hat{\Phi}_s(l) + \hat{\Phi}_d(l)/(\mathbf{a}^H(l)\Gamma^{-1}\mathbf{a}(l))}$
requires estimate of speech PSD $\Phi_s(l)$ and diffuse PSD $\Phi_d(l)$

EVD-BASED DIFFUSE PSD ESTIMATOR

- exploit PSD matrix structure
 - prewhitening based on Cholesky decomposition: $\Gamma = \mathbf{L}\mathbf{L}^H$
$$\Phi_{\mathbf{y}}^w(l) = \mathbf{L}^{-1}\Phi_{\mathbf{y}}(l)\mathbf{L}^{-H} = \underbrace{\Phi_s(l)\mathbf{L}^{-1}\mathbf{a}(l)}_{\mathbf{b}(l)}\underbrace{\mathbf{a}^H(l)\mathbf{L}^{-H}}_{\mathbf{b}^H(l)} + \Phi_d(l)\mathbf{I}_M$$
 - rank-1 term $\Phi_s(l)\mathbf{b}(l)\mathbf{b}^H(l)$ has one non-zero eigenvalue $\sigma(l)$
 - term $\Phi_d(l)\mathbf{I}_M$ adds offset $\Phi_d(l)$ to eigenvalues
$$\lambda_1\{\Phi_{\mathbf{y}}^w(l)\} = \sigma(l) + \Phi_d(l), \quad \lambda_i\{\Phi_{\mathbf{y}}^w(l)\} = \Phi_d(l), \quad i \in \{2, \dots, M\}$$
 - in practice, **model not perfect** $\rightarrow \lambda_i \neq \lambda_j, \quad i, j \in \{2, \dots, M\}, i \neq j$
- estimate $\Phi_d(l)$ using either second eigenvalue or mean of smallest $M - 1$ eigenvalues

$$\hat{\Phi}_{d,\text{EIG2}}(l) = \lambda_2\{\Phi_{\mathbf{y}}^w(l)\}$$

$$\hat{\Phi}_{d,\text{EIG1}}(l) = \frac{1}{M-1} (\text{trace}\{\Phi_{\mathbf{y}}^w(l)\} - \lambda_1\{\Phi_{\mathbf{y}}^w(l)\})$$

- eigenvalue decomposition of $M \times M$ matrix required for each STFT bin

high performance high complexity

POWER METHOD

- since only first or second eigenvalue of $\Phi_{\mathbf{y}}^w(l)$ are required, computational cost can be reduced using **power method**
- power method iteratively estimates **dominant eigenvalue** provided that

$$|\lambda_1| > |\lambda_2| \geq \dots \geq |\lambda_M| \quad \checkmark$$

and can also estimate **additional eigenvalues** using rank reduction

- initialization:** random or estimate from previous frame
- convergence speed:** depends on $|\lambda_1|/|\lambda_2|$
- complexity (flops):**
 - λ_1 : $N(16M^2 + 2M - 2)$
 - λ_2 : $2N(16M^2 + 2M - 2) + 5M^2$
 - full EVD using Hessenberg QR algorithm: $4/3M^3 + \mathcal{O}(M^2)$

In: $\Phi_{\mathbf{y}}^w(l) \in \mathbb{C}^{M \times M}$, number of iterations N
Out: estimates $\hat{\lambda}_1\{\Phi_{\mathbf{y}}^w(l)\}, \hat{\lambda}_2\{\Phi_{\mathbf{y}}^w(l)\}$
for $m = 1$ to 2 do
 initialize $\mathbf{u}_m^{(0)} \in \mathbb{C}^M$;
 for $n = 1$ to N do
 $\mathbf{t} = \Phi_{\mathbf{y}}^w(l)\mathbf{u}_m^{(n-1)}$; // power iteration
 $\mathbf{u}_m^{(n)} = \mathbf{t}/\|\mathbf{t}\|_2$;
 $\lambda_m^{(n)} = \mathbf{u}_m^{(n)H}\Phi_{\mathbf{y}}^w(l)\mathbf{u}_m^{(n)}$; // RAYLEIGH quotient
 end
 $\hat{\lambda}_m\{\Phi_{\mathbf{y}}^w(l)\} = \lambda_m^{(N)}$;
 // matrix rank reduction
 $\Phi_{\mathbf{y}}^w(l) = \Phi_{\mathbf{y}}^w(l) - \hat{\lambda}_m\{\Phi_{\mathbf{y}}^w(l)\}\mathbf{u}_m^{(N)}\mathbf{u}_m^{(N)H}$;
end

SIMULATION SETUP

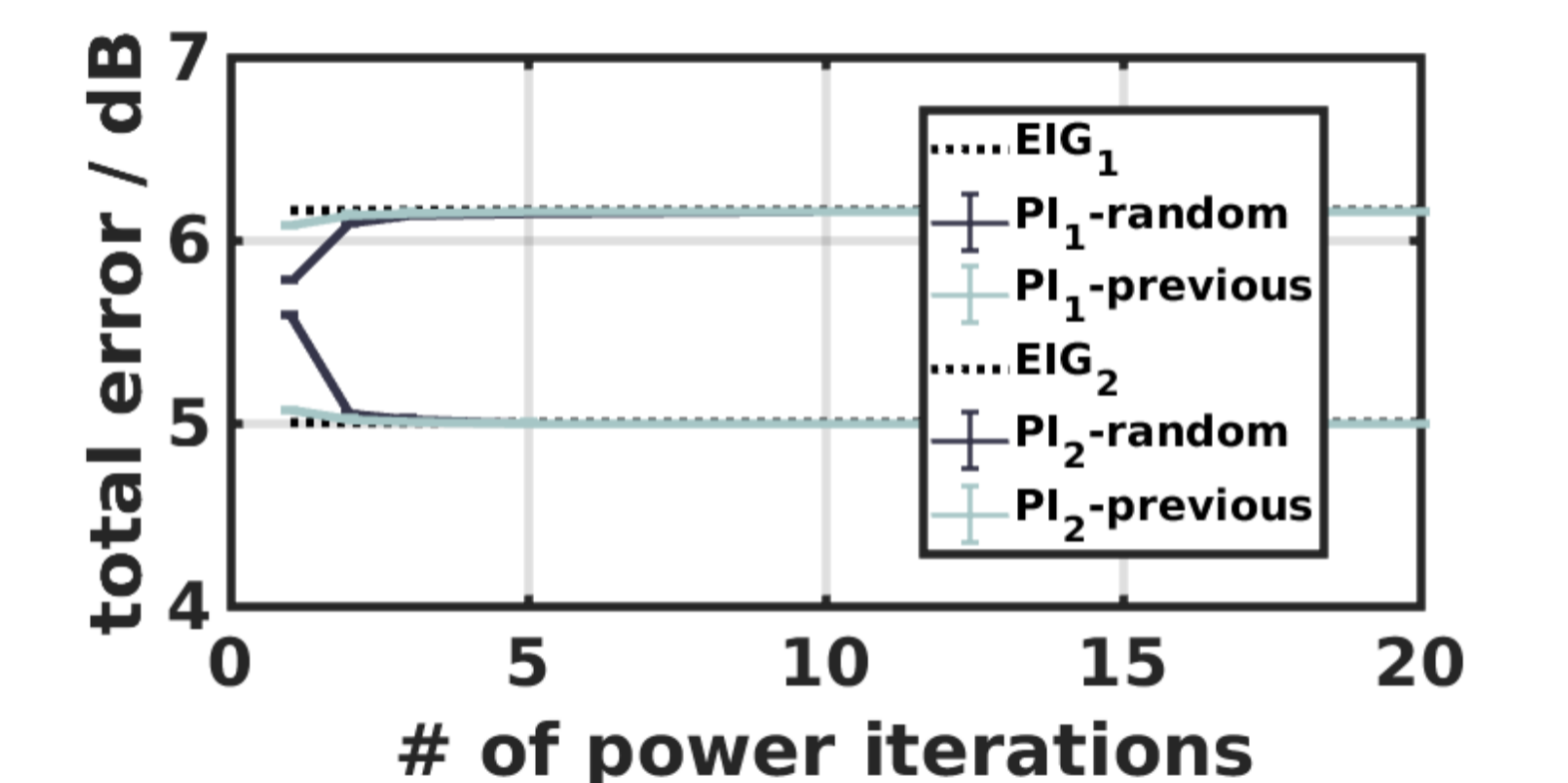
- $M = 4$ microphones, $f_s = 16$ kHz
- diffuse babble noise [2] at $\text{SNR}_{\text{in}} = \{10, \dots, 40\}$ dB, no sensor noise
- STFT: 64 ms frame length ($N_{\text{FFT}} = 1024$), 16 ms shift
- $\Phi_{\mathbf{y}}(l)$ estimated using recursive averaging, 40 ms smoothing constant
- speech PSD $\Phi_s(l)$ estimated using decision-directed approach[3]
- performance measures:
 - PSD estimation error (averaged over frames and frequencies)
 - speech quality of processed signal $\hat{S}(l)$ using fwsSNR and PESQ

	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

RESULTS

- PSD estimation accuracy:** different initialization, $\text{SNR}_{\text{in}} = 10$ dB, AS₁

- convergence to full EVD after few iterations
- best initialization utilizing estimate of previous frame

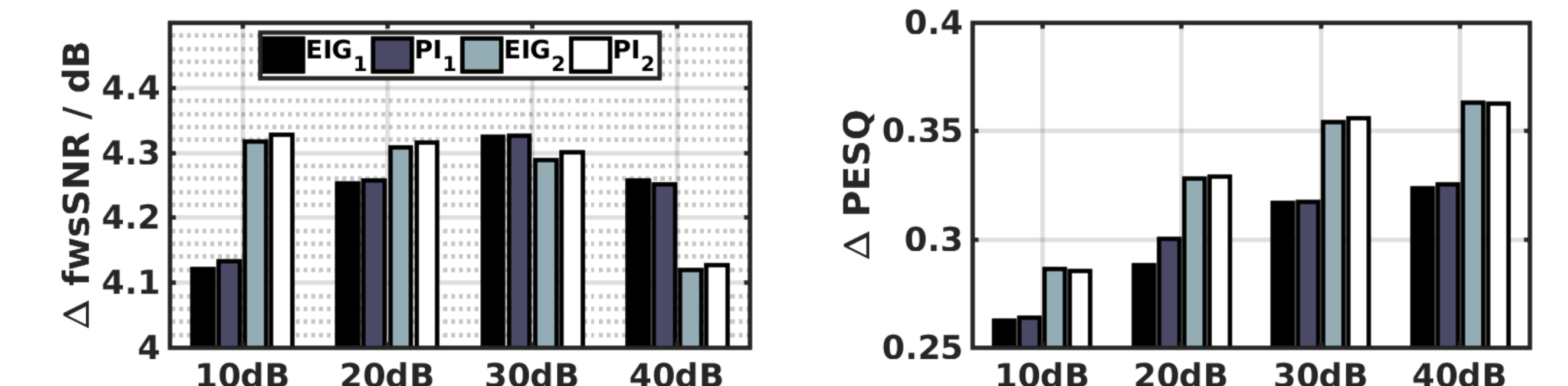


- computational complexity:** non-optimized runtimes, $N = 2$

- power method significantly faster
- no noticeable difference in accuracy

Method	Total Error / dB	Av. Duration / 10 ⁻⁵ s
Power Method, λ_1	6.09	0.35
QR Method, λ_1	6.11	1.03
QZ Method (MATLAB), λ_1	6.16	0.55

- speech quality** of processed signal using different PSD estimates (average over AS₁₋₃, $N = 2$)



- no significant difference between power method and full EVD