

Blind multi-microphone noise reduction and dereverberation algorithms for speech communication applications

Prof. Dr. Simon Doclo

University of Oldenburg, Germany

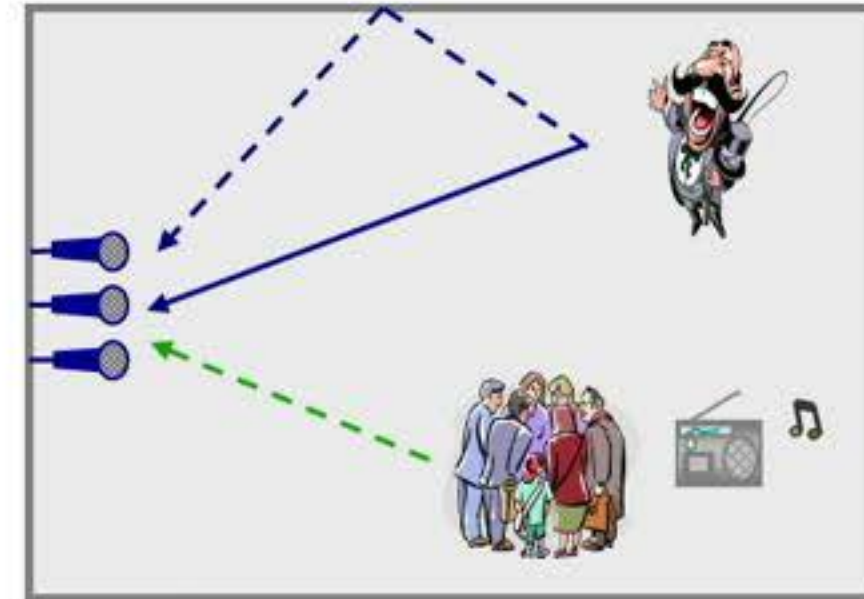
Dept. of Medical Physics and Acoustics, Cluster of Excellence Hearing4all

<http://www.sigproc.uni-oldenburg.de/>

Microsoft Research, 29.10.2019

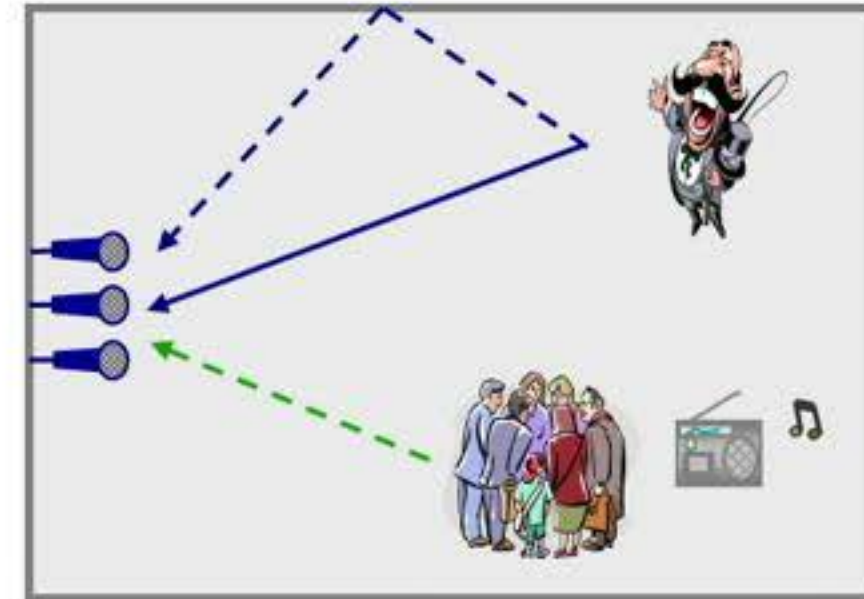
Introduction

- **Problem**
 - **Ambient noise and reverberation** jointly present in typical acoustic environments
 - **Speech quality and intelligibility** degradation for speech communication applications
 - Performance degradation of voice-controlled systems



Introduction

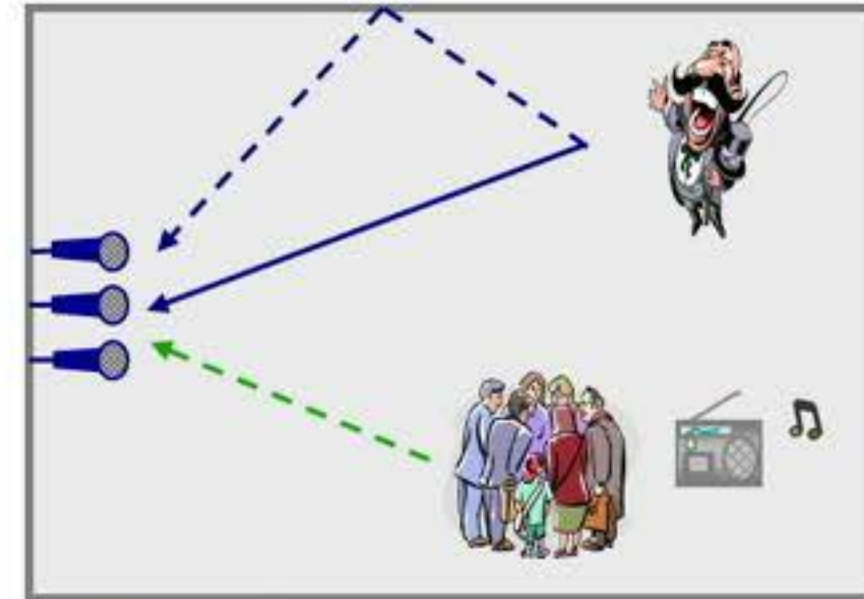
- **Problem**
 - **Ambient noise and reverberation** jointly present in typical acoustic environments
 - **Speech quality and intelligibility** degradation for speech communication applications
 - Performance degradation of voice-controlled systems



Introduction

- **Problem**

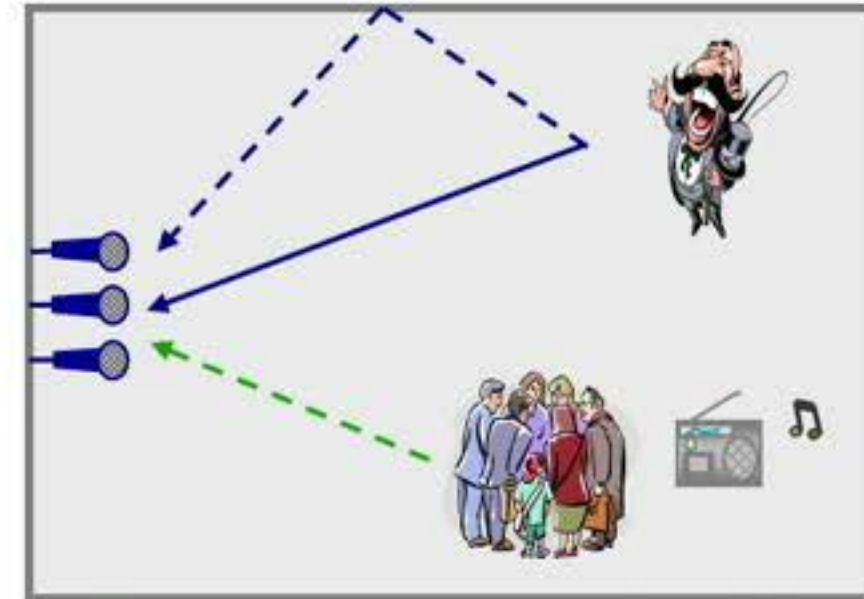
- **Ambient noise and reverberation** jointly present in typical acoustic environments
- **Speech quality and intelligibility** degradation for speech communication applications
- Performance degradation of voice-controlled systems



Introduction

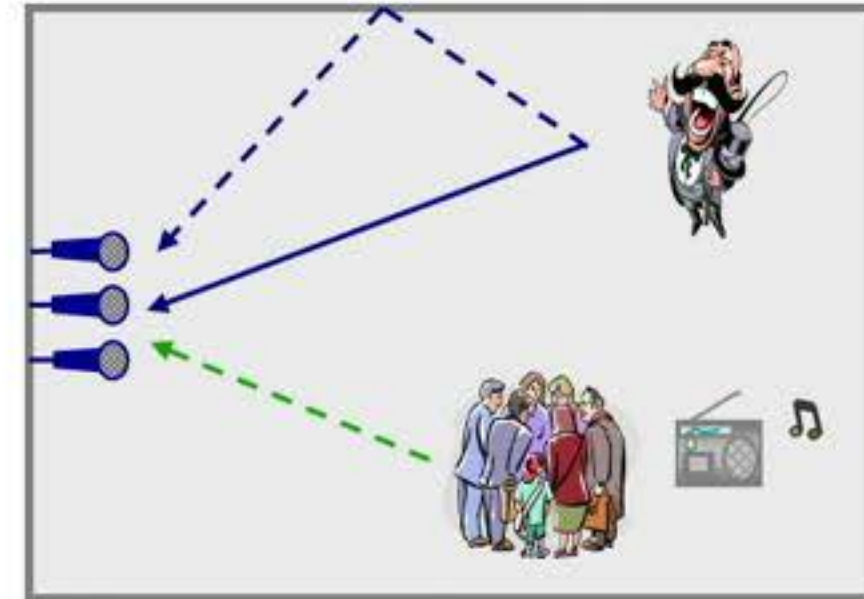
- **Problem**

- **Ambient noise and reverberation** jointly present in typical acoustic environments
- **Speech quality and intelligibility** degradation for speech communication applications
- Performance degradation of voice-controlled systems



Introduction

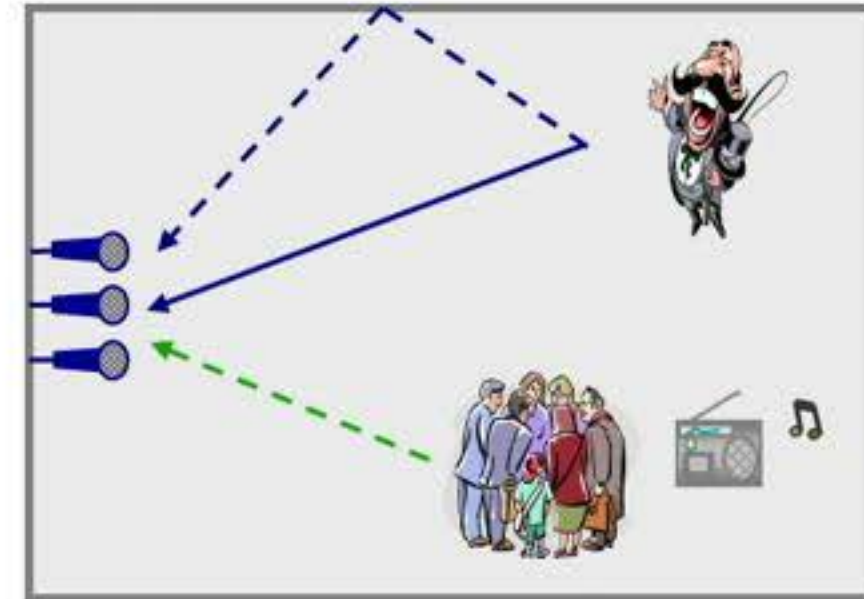
- **Objectives**
 - Single- and **multi-microphone joint noise reduction and dereverberation** algorithms
 - Speech communication applications: **blind and on-line processing** for time-varying dynamic acoustic scenarios
 - Exploit knowledge or (statistical) models of speech signals and room acoustics



Introduction

- **Objectives**

- Single- and **multi-microphone joint noise reduction and dereverberation** algorithms
- Speech communication applications: **blind and on-line processing** for time-varying dynamic acoustic scenarios
- Exploit knowledge or (statistical) models of speech signals and room acoustics



- **This presentation**

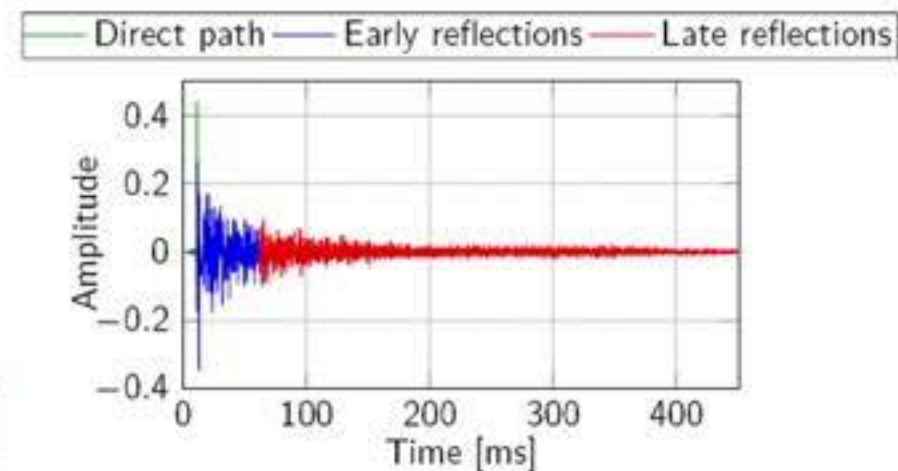
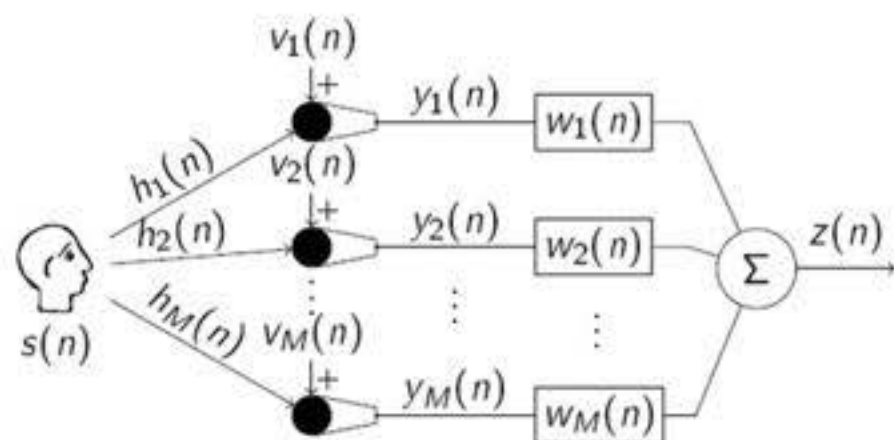
1. **Joint estimation of (time-varying) spatial and spectral variables** for multi-microphone speech enhancement
2. **Binaural hearing devices:** combination of speech enhancement and preservation of auditory scene
3. **Extension to acoustic sensor networks** with spatially distributed microphones

1. Joint dereverberation and noise reduction

Signal model

- **Scenario:** speech source in noisy and reverberant environment, M microphones
- **Model in Short-Time Fourier Transform (STFT) domain:**

$$\begin{aligned}
 y_m(k, l) &= h_m(k, l) \star s(k, l) + v_m(k, l) \\
 &= \underset{\substack{\uparrow \\ \text{direct and early} \\ \text{reverberation}}}{a_m(k, l)} x_1(k, l) + \underset{\substack{\uparrow \\ \text{late} \\ \text{reverberation}}}{x_{r,m}(k, l)} + \underset{\substack{\uparrow \\ \text{ambient} \\ \text{noise}}}{v_m(k, l)}
 \end{aligned}$$



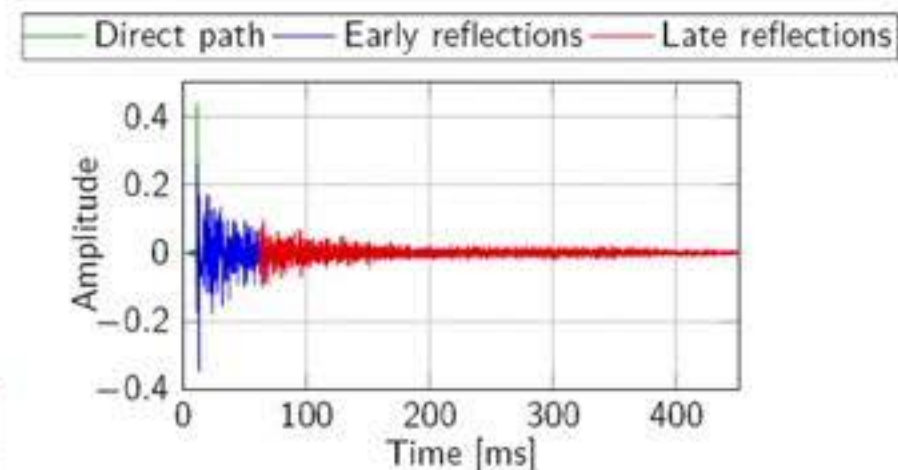
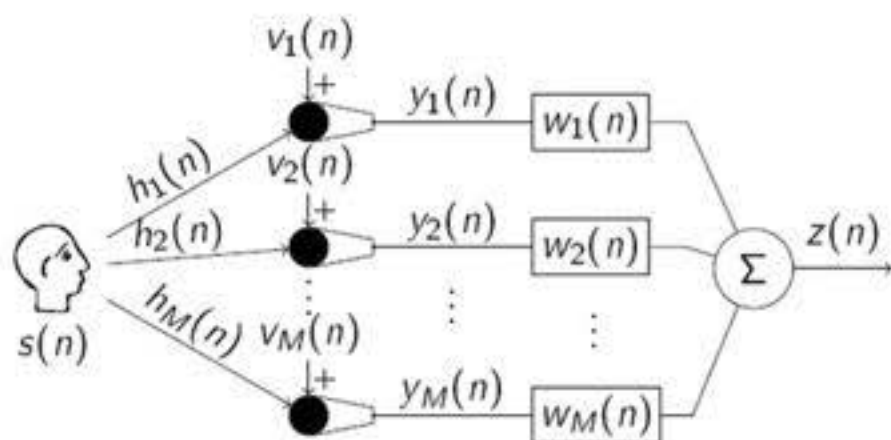
Signal model

- **Scenario:** speech source in noisy and reverberant environment, M microphones
- **Model in Short-Time Fourier Transform (STFT) domain:**

$$\begin{aligned} y_m(k, l) &= h_m(k, l) \star s(k, l) + v_m(k, l) \\ &= \underset{\substack{\uparrow \\ \text{direct and early} \\ \text{reverberation}}}{a_m(k, l)} x_1(k, l) + \underset{\substack{\uparrow \\ \text{late} \\ \text{reverberation}}}{x_{r,m}(k, l)} + \underset{\substack{\uparrow \\ \text{ambient} \\ \text{noise}}}{v_m(k, l)} \end{aligned}$$

$$\mathbf{y}(k, l) = \mathbf{a}(k, l) x_1(k, l) + \mathbf{x}_r(k, l) + \mathbf{v}(k, l)$$

$\mathbf{a}(k, l)$ = vector of **relative early transfer functions (RETFs)** of target source



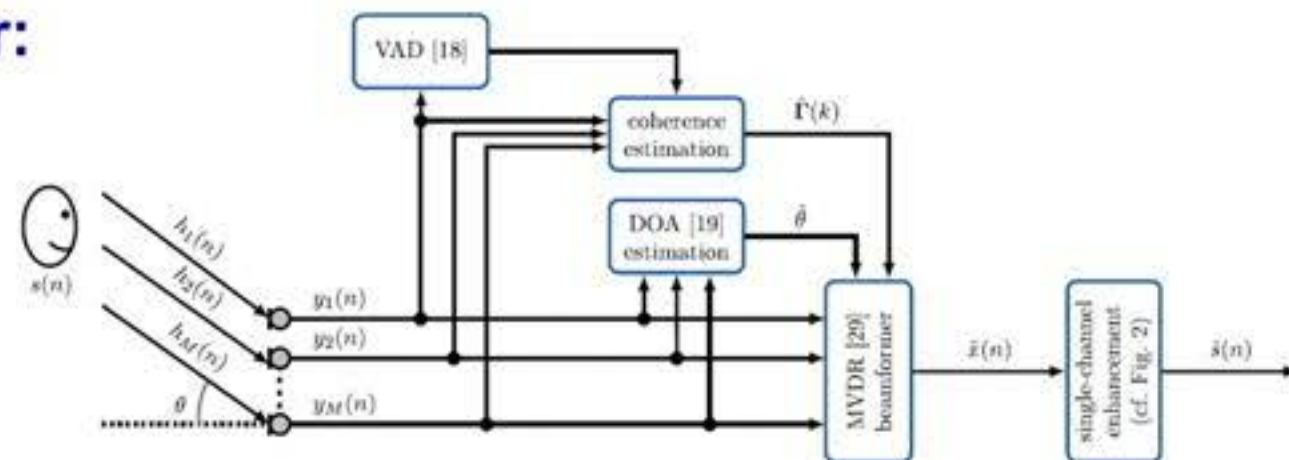
Multi-microphone dereverberation and noise reduction

1. Beamforming + spectral postfilter:

multiply each time-frequency bin with real-valued gain

$$\mathbf{y}(l) = \mathbf{a}(l)x_1(l) + \mathbf{x}_r(l) + \mathbf{v}(l)$$

$$\mathbf{W}_{MWF} = \frac{(\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a}}{\mathbf{a}^H (\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + (\mathbf{a}^H (\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a})^{-1}}$$

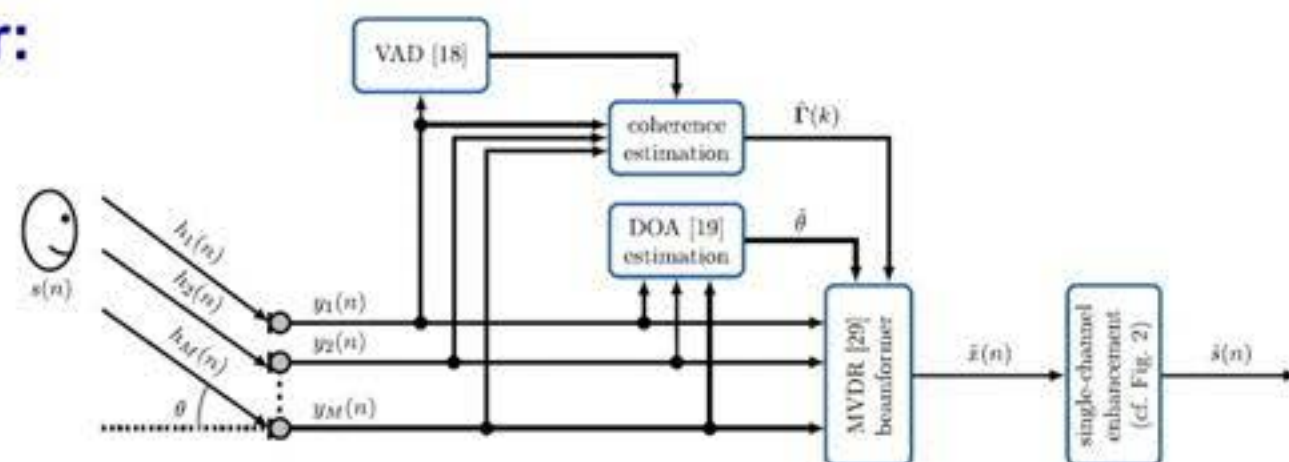


Multi-microphone dereverberation and noise reduction

1. Beamforming + spectral postfilter:

multiply each time-frequency bin with real-valued gain

$$\mathbf{y}(l) = \mathbf{a}(l)x_1(l) + \mathbf{x}_r(l) + \mathbf{v}(l)$$

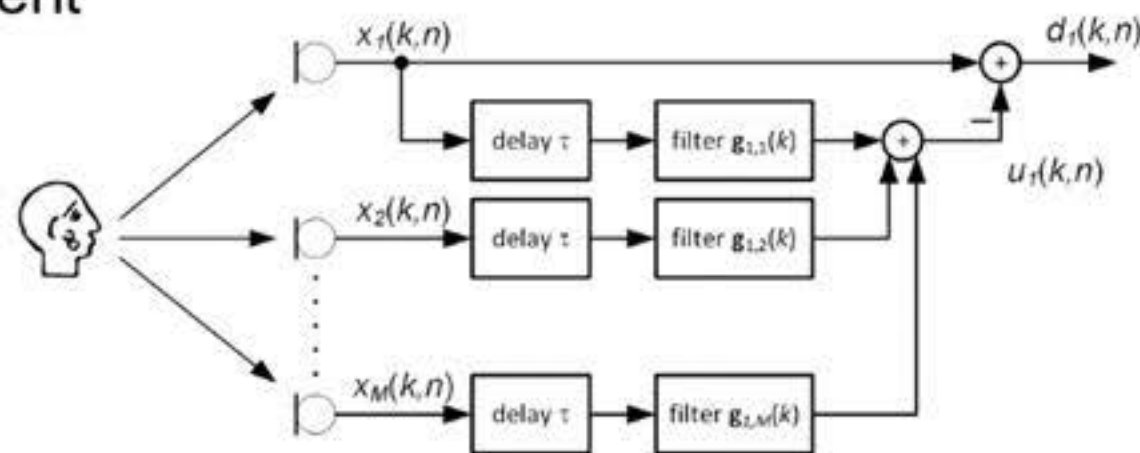


$$\mathbf{W}_{MWF} = \frac{(\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a}}{\mathbf{a}^H (\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + (\mathbf{a}^H (\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a})^{-1}}$$

2. Reverberation and noise suppression: *subtract complex-valued estimate of late reverberant and noise component*

$$y_m(l) = h_m(l) \star s(l) + v_m(l)$$

$$\hat{x}_{e,1}(l) = y_1(l) - \mathbf{Y}_\tau(l) \mathbf{g}(l)$$



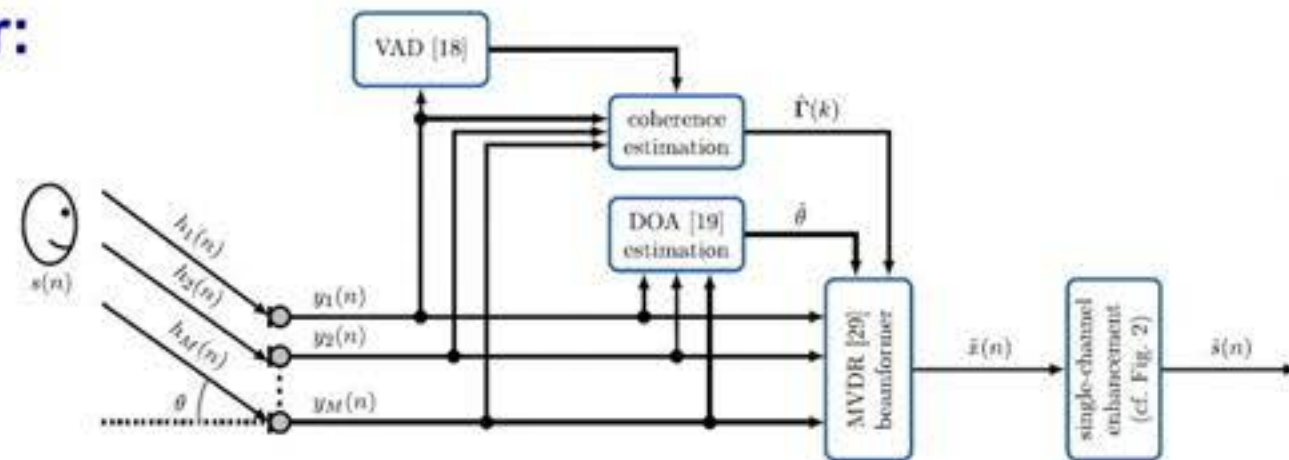
Multi-microphone dereverberation and noise reduction

1. Beamforming + spectral postfilter:

multiply each time-frequency bin with real-valued gain

$$\mathbf{y}(l) = \mathbf{a}(l)x_1(l) + \mathbf{x}_r(l) + \mathbf{v}(l)$$

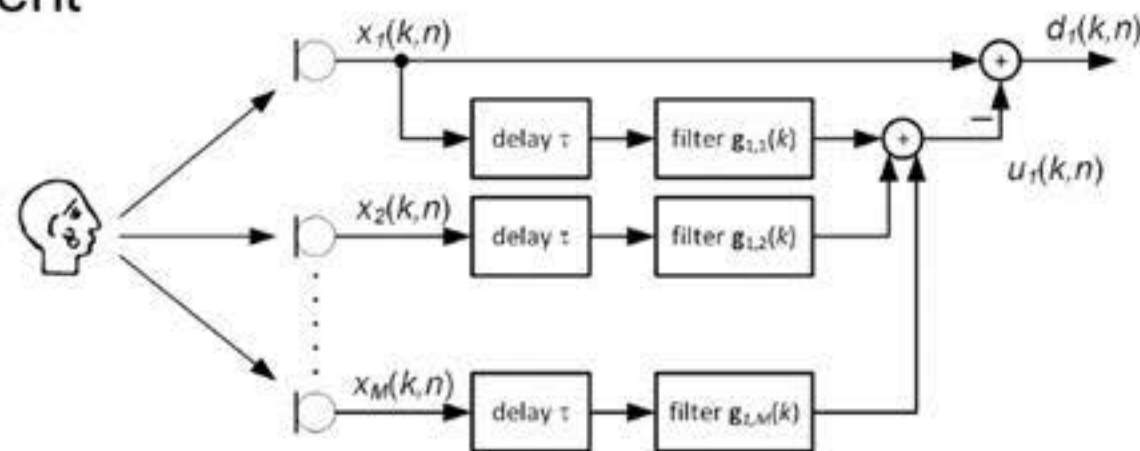
$$\mathbf{W}_{MWF} = \frac{(\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a}}{\mathbf{a}^H (\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + (\mathbf{a}^H (\phi_d \mathbf{\Gamma} + \mathbf{\Phi}_v)^{-1} \mathbf{a})^{-1}}$$



2. Reverberation and noise suppression: *subtract complex-valued estimate of late reverberant and noise component*

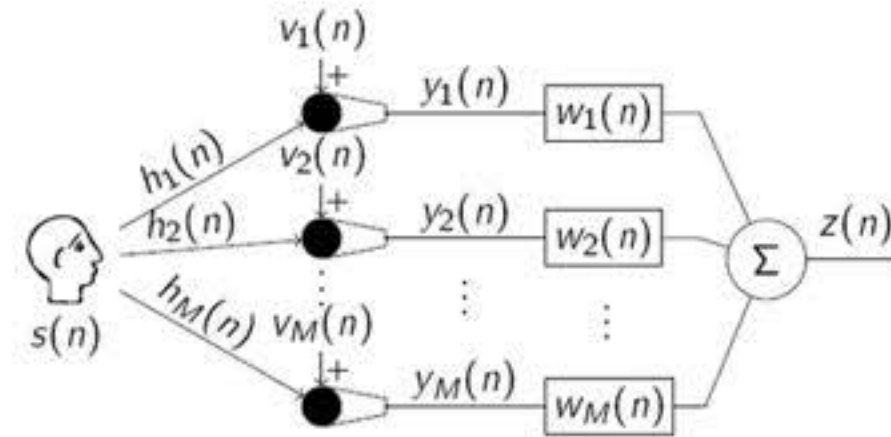
$$y_m(l) = h_m(l) \star s(l) + v_m(l)$$

$$\hat{x}_{e,1}(l) = y_1(l) - \mathbf{Y}_\tau(l) \mathbf{g}(l)$$



Beamforming + spectral postfilter

- Filter-and-sum structure : $z = \mathbf{w}^H \mathbf{y}$



Beamforming + spectral postfilter

- **Filter-and-sum structure :** $z = \mathbf{w}^H \mathbf{y}$ $\mathbf{y} = \mathbf{a}x_1 + \mathbf{n}$
- **“Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)**
Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Beamforming + spectral postfilter

- **Filter-and-sum structure :**

$$z = \mathbf{w}^H \mathbf{y}$$

$$\mathbf{y} = \mathbf{a}x_1 + \mathbf{n}$$

- **“Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)**

Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Beamforming + spectral postfilter

- **Filter-and-sum structure :** $z = \mathbf{w}^H \mathbf{y}$ $\mathbf{y} = \mathbf{a}x_1 + \mathbf{n}$
- **“Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)**
Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Beamforming + spectral postfilter

- **Filter-and-sum structure :**

$$z = \mathbf{w}^H \mathbf{y}$$

$$\mathbf{y} = \mathbf{a}x_1 + \mathbf{n}$$

- **“Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)**

Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Can be decomposed as **MVDR beamformer and spectral postfilter**

$$\mathbf{W}_{MWF} = \frac{\Phi_n^{-1} \mathbf{a}}{\mathbf{a}^H \Phi_n^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + \mu (\mathbf{a}^H \Phi_n^{-1} \mathbf{a})^{-1}}$$

Requires estimate/model of interference **covariance matrix** Φ_n ,
estimate/model of **relative (early) transfer function vector** \mathbf{a} of desired
source, and **PSDs** of speech and interference components at MVDR output

Beamforming + spectral postfilter

- Filter-and-sum structure :

$$z = \mathbf{w}^H \mathbf{y}$$

$$\mathbf{y} = \mathbf{a}x_1 + \mathbf{n}$$

- “Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)

Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Can be decomposed as **MVDR beamformer and spectral postfilter**

$$\mathbf{w}_{MWF} = \frac{\Phi_n^{-1} \mathbf{a}}{\mathbf{a}^H \Phi_n^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + \mu (\mathbf{a}^H \Phi_n^{-1} \mathbf{a})^{-1}}$$

Requires estimate/model of interference **covariance matrix** Φ_n , estimate/model of **relative (early) transfer function vector** \mathbf{a} of desired source, and **PSDs** of speech and interference components at MVDR output

Beamforming + spectral postfilter

- Filter-and-sum structure :

$$z = \mathbf{w}^H \mathbf{y}$$

$$\mathbf{y} = \mathbf{a}x_1 + \mathbf{n}$$

- “Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)

Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Can be decomposed as **MVDR beamformer and spectral postfilter**

$$\mathbf{W}_{MWF} = \frac{\Phi_n^{-1} \mathbf{a}}{\mathbf{a}^H \Phi_n^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + \mu (\mathbf{a}^H \Phi_n^{-1} \mathbf{a})^{-1}}$$

Requires estimate/model of interference **covariance matrix** Φ_n ,
estimate/model of **relative (early) transfer function vector** \mathbf{a} of desired
source, and **PSDs** of speech and interference components at MVDR output

Beamforming + spectral postfilter

- **Filter-and-sum structure :** $z = \mathbf{w}^H \mathbf{y}$ $y = \mathbf{a}x_1 + \mathbf{n}$
- **“Workhorse algorithm”: parametric Multi-channel Wiener filter (MWF)**
Goal: estimate desired speech component in reference microphone + trade off interference (*noise and/or reverberation*) reduction and speech distortion

$$\min_{\mathbf{w}} \mathcal{E}\{|\mathbf{w}^H \mathbf{x} - x_1|^2\} + \mu \mathcal{E}\{|\mathbf{w}^H \mathbf{n}|^2\} \Rightarrow \mathbf{w}_{MWF} = (\Phi_x + \mu \Phi_n)^{-1} \Phi_x \mathbf{e}$$

Requires estimate of covariance matrices, e.g., based on speech presence probability (SPP)

Can be decomposed as **MVDR beamformer and spectral postfilter**

$$\mathbf{W}_{MWF} = \frac{\Phi_n^{-1} \mathbf{a}}{\mathbf{a}^H \Phi_n^{-1} \mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + \mu (\mathbf{a}^H \Phi_n^{-1} \mathbf{a})^{-1}}$$

Requires estimate/model of interference **covariance matrix** Φ_n ,
 estimate/model of **relative (early) transfer function vector** \mathbf{a} of desired source, and **PSDs** of speech and interference components at MVDR output

Beamforming + spectral postfilter

- Signal model

$$\mathbf{y}(l) = \mathbf{a}(l)x_1(l) + \mathbf{x}_r(l) + \mathbf{v}(l)$$

$$\Phi_y(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \Phi_{x_r}(l) + \Phi_v(l)$$

Late reverberation: model as diffuse sound field $\Phi_{x_r}(l) = \phi_d(l)\Gamma$

with $\phi_d(l)$ *time-varying* diffuse PSD and Γ *time-invariant* coherence matrix
(also incorporating diffuse noise !)

$$\mathbf{W}_{MWF} = \frac{(\phi_d\Gamma + \Phi_v)^{-1}\mathbf{a}}{\mathbf{a}^H(\phi_d\Gamma + \Phi_v)^{-1}\mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + (\mathbf{a}^H(\phi_d\Gamma + \Phi_v)^{-1}\mathbf{a})^{-1}}$$

Beamforming + spectral postfilter

- Signal model

$$\mathbf{y}(l) = \mathbf{a}(l)x_1(l) + \mathbf{x}_r(l) + \mathbf{v}(l)$$

$$\Phi_y(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \Phi_{x_r}(l) + \Phi_v(l)$$

Late reverberation: model as diffuse sound field $\Phi_{x_r}(l) = \phi_d(l)\Gamma$

with $\phi_d(l)$ *time-varying* diffuse PSD and Γ *time-invariant* coherence matrix (also incorporating diffuse noise !)

$$\mathbf{W}_{MWF} = \frac{(\phi_d\Gamma + \Phi_v)^{-1}\mathbf{a}}{\mathbf{a}^H(\phi_d\Gamma + \Phi_v)^{-1}\mathbf{a}} \cdot \frac{\phi_{x_1}}{\phi_{x_1} + (\mathbf{a}^H(\phi_d\Gamma + \Phi_v)^{-1}\mathbf{a})^{-1}}$$

- Key estimation tasks:

- RETF vector $\mathbf{a}(l)$:** *anechoic* (based on DOA estimate) or *reverberant*
- Diffuse/late reverberant PSD $\phi_d(l)$:** using single-channel *temporal model* (exponential decay) or based on multi-channel *diffuse sound field model*
- Noise covariance matrix $\Phi_v(l)$:** *estimate* (based on SPP) or *model* (e.g., spatially white noise)

Estimation of PSDs

- **Requiring estimate of RETF vector and noise covariance matrix**

$$\hat{\Phi}_x(l) = \hat{\Phi}_y(l) - \hat{\Phi}_v(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\mathbf{\Gamma}$$

- *Maximum-likelihood estimators*, requiring iterative optimisation procedure
- *Closed-form least-squares estimators*, based on Frobenius norm

$$\min_{\phi_{x_1}(l), \phi_d(l)} \|\hat{\Phi}_x(l) - \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) - \phi_d(l)\mathbf{\Gamma}\|_F^2$$

Estimation of PSDs

- Requiring estimate of RETF vector and noise covariance matrix

$$\hat{\Phi}_x(l) = \hat{\Phi}_y(l) - \hat{\Phi}_v(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\mathbf{\Gamma}$$

- Maximum-likelihood estimators, requiring iterative optimisation procedure
- Closed-form least-squares estimators, based on Frobenius norm

$$\min_{\phi_{x_1}(l), \phi_d(l)} \|\hat{\Phi}_x(l) - \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) - \phi_d(l)\mathbf{\Gamma}\|_F^2$$

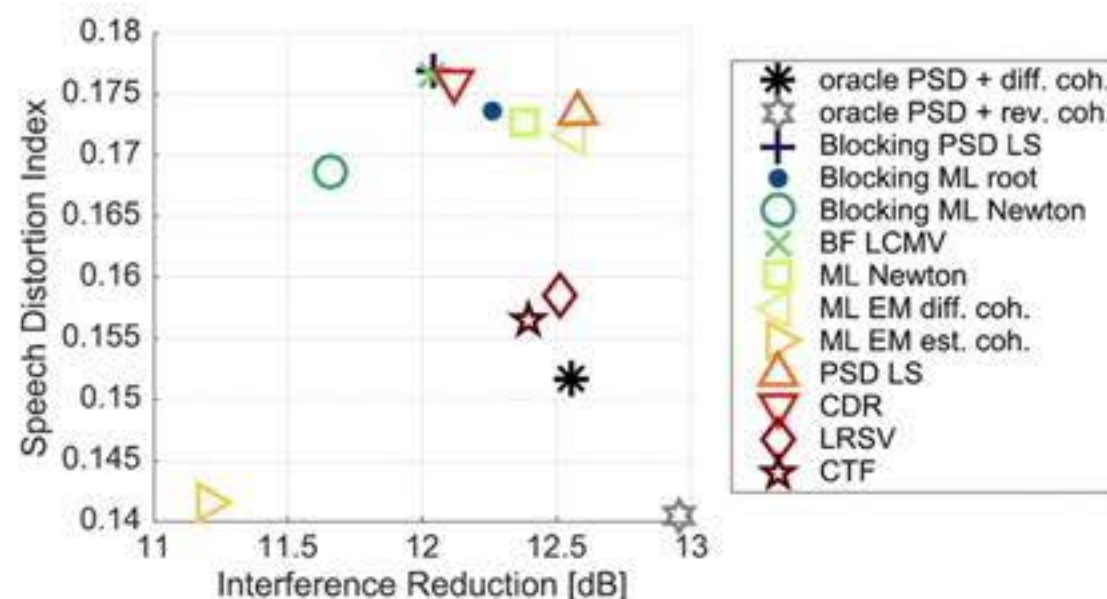


Fig. 9. Speech distortion vs. interference reduction for RSNR = 15 dB.

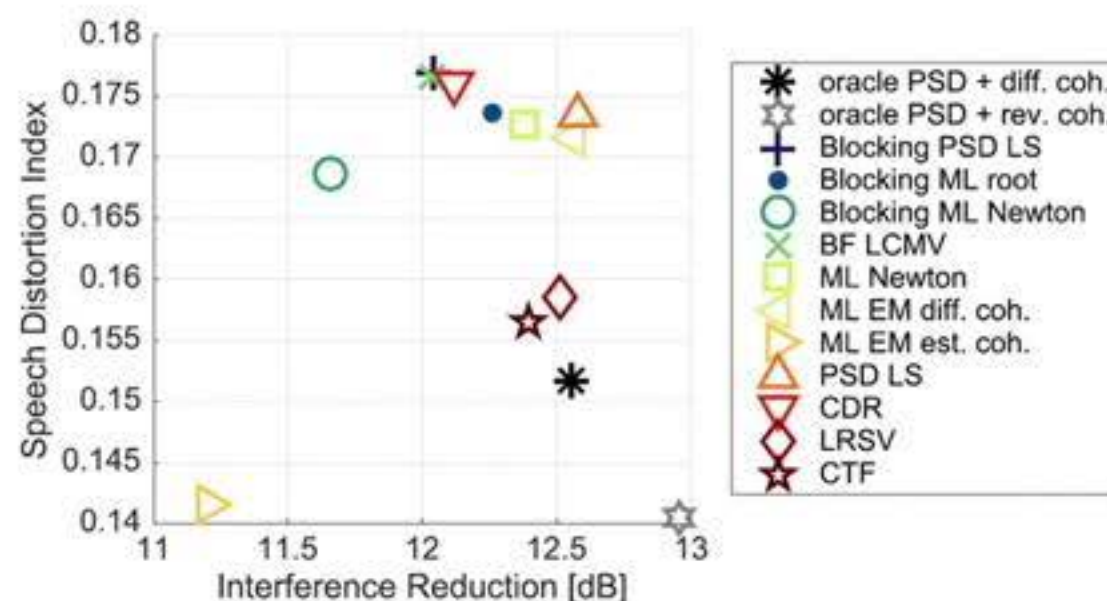
Estimation of PSDs

- Requiring estimate of RETF vector and noise covariance matrix

$$\hat{\Phi}_x(l) = \hat{\Phi}_y(l) - \hat{\Phi}_v(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\mathbf{\Gamma}$$

- Maximum-likelihood estimators, requiring iterative optimisation procedure
- Closed-form least-squares estimators, based on Frobenius norm

$$\min_{\phi_{x_1}(l), \phi_d(l)} \|\hat{\Phi}_x(l) - \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) - \phi_d(l)\mathbf{\Gamma}\|_F^2$$



Similar performance for most methods...

Fig. 9. Speech distortion vs. interference reduction for RSNR = 15 dB.

Joint Estimation of RETF vector and PSDs

1. Covariance whitening (CW) method:

- Requires estimate of noise covariance matrix

$$\hat{\Phi}_x(l) = \hat{\Phi}_y(l) - \hat{\Phi}_v(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\mathbf{\Gamma}$$

Joint Estimation of RETF vector and PSDs

1. Covariance whitening (CW) method:

- Requires estimate of noise covariance matrix

$$\hat{\Phi}_x(l) = \hat{\Phi}_y(l) - \hat{\Phi}_v(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\mathbf{\Gamma}$$

- Eigenvalue decomposition of prewhitened signal correlation matrix

$$\hat{\Phi}_x^w(l) = \mathbf{\Gamma}^{-1/2}\hat{\Phi}_x(l)\mathbf{\Gamma}^{-H/2} = \phi_{x_1}(l)\mathbf{b}(l)\mathbf{b}^H(l) + \phi_d(l)\mathbf{I}$$

Joint Estimation of RETF vector and PSDs

1. Covariance whitening (CW) method:

- Requires estimate of noise covariance matrix

$$\hat{\Phi}_x(l) = \hat{\Phi}_y(l) - \hat{\Phi}_v(l) = \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) + \phi_d(l)\mathbf{\Gamma}$$

- Eigenvalue decomposition of prewhitened signal correlation matrix

$$\hat{\Phi}_x^w(l) = \mathbf{\Gamma}^{-1/2}\hat{\Phi}_x(l)\mathbf{\Gamma}^{-H/2} = \phi_{x_1}(l)\mathbf{b}(l)\mathbf{b}^H(l) + \phi_d(l)\mathbf{I}$$

- Principal eigenvector $\mathbf{u}(l)$: estimate of RETF vector

$$\hat{\mathbf{a}}(l) = \frac{\mathbf{\Gamma}^{1/2}\mathbf{u}(l)}{\mathbf{e}^T\mathbf{\Gamma}^{1/2}\mathbf{u}(l)}$$

- Eigenvalues: estimate of PSDs

$$\hat{\phi}_d(l) = \lambda_2\{\hat{\Phi}_x^w(l)\} \quad \hat{\phi}_{d,\mu}(l) = \frac{1}{M-1}(\text{tr}\{\hat{\Phi}_x^w(l)\} - \lambda_1\{\hat{\Phi}_x^w(l)\})$$

$$\hat{\phi}_{x_1}(l) = \lambda_1\{\hat{\Phi}_x^w(l)\}/\|\hat{\mathbf{b}}\|_2^2$$

Joint Estimation of RETF vector and PSDs

2. **Alternating least squares (ALS) method**, minimizing Frobenius norm
 - **Model noise covariance matrix + estimate noise PSD**

$$\min_{\phi_{x_1}(l), \phi_d(l), \phi_v(l), \mathbf{a}(l)} \|\hat{\Phi}_y(l) - \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) - \phi_d(l)\mathbf{\Gamma} - \phi_v(l)\mathbf{\Psi}\|_F^2$$

Joint Estimation of RETF vector and PSDs

2. Alternating least squares (ALS) method, minimizing Frobenius norm

- Model noise covariance matrix + estimate noise PSD

$$\min_{\phi_{x_1}(l), \phi_d(l), \phi_v(l), \mathbf{a}(l)} \|\hat{\Phi}_y(l) - \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) - \phi_d(l)\mathbf{\Gamma} - \phi_v(l)\mathbf{\Psi}\|_F^2$$

- No closed-form solution → two-step alternating procedure
(least-squares problem for PSDs, eigenvalue problem for RETF vector)

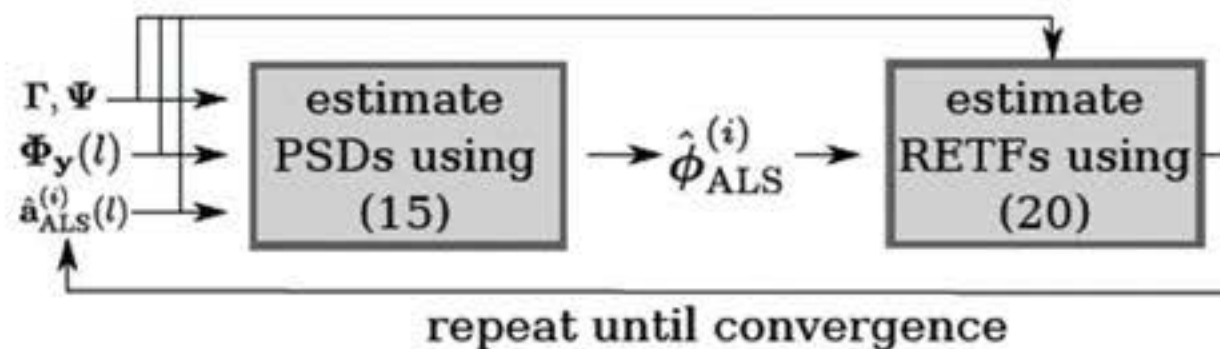


Fig. 1: Block diagram of ALS-based RETF vector and PSD estimation.

Joint Estimation of RETF vector and PSDs

2. Alternating least squares (ALS) method, minimizing Frobenius norm

- Model noise covariance matrix + estimate noise PSD

$$\min_{\phi_{x_1}(l), \phi_d(l), \phi_v(l), \mathbf{a}(l)} \|\hat{\Phi}_y(l) - \phi_{x_1}(l)\mathbf{a}(l)\mathbf{a}^H(l) - \phi_d(l)\mathbf{\Gamma} - \phi_v(l)\mathbf{\Psi}\|_F^2$$

- No closed-form solution → two-step alternating procedure
(least-squares problem for PSDs, eigenvalue problem for RETF vector)

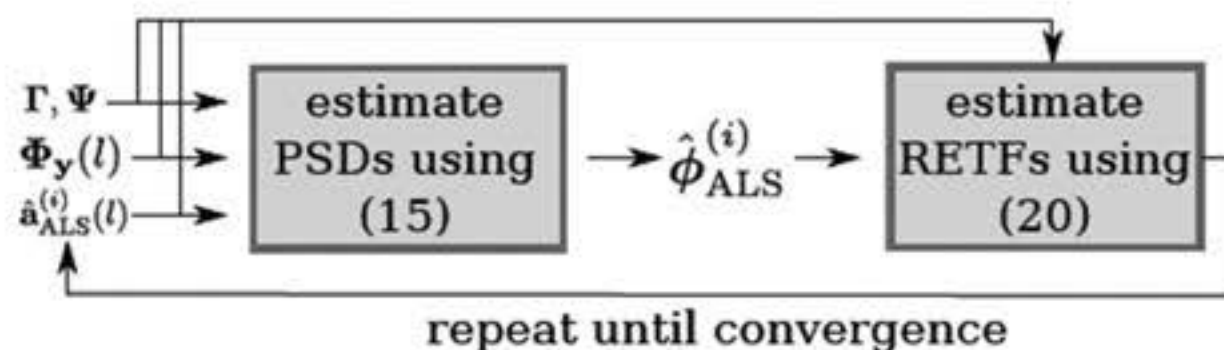


Fig. 1: Block diagram of ALS-based RETF vector and PSD estimation.

Algorithm 1: ALS method to jointly estimate the RETF vector and PSDs.

Input: $\mathbf{\Gamma}(k)$, $\mathbf{\Psi}(k)$, $\hat{\Phi}_y(k, l)$, iterations N , init. $\hat{\mathbf{a}}^{(0)}(k, 1)$

Output: $\hat{\mathbf{a}}_{\text{ALS}}(k, l)$, $\hat{\phi}_{\text{ALS}}(k, l)$

```

for all (k, l) do
  for i = 1 : N do
    compute  $\mathbf{A}^{(i-1)}(k, l)$  using (16) and  $\mathbf{b}^{(i-1)}(k, l)$  using (17)
     $\hat{\phi}_{\text{ALS}}^{(i)}(k, l) = (\mathbf{A}^{(i-1)}(k, l))^{-1} \mathbf{b}^{(i-1)}(k, l)$  (15)
    constrain  $\hat{\phi}_{\text{ALS}}^{(i)}(k, l)$  using (25)
     $\hat{\Phi}_x^{(i)}(k, l) = \hat{\Phi}_y(k, l) - (\hat{\phi}_{d, \text{ALS}}^{(i)}(k, l)\mathbf{\Gamma}(k) + \hat{\phi}_{v, \text{ALS}}^{(i)}(k, l)\mathbf{\Psi}(k))$ 
     $\hat{\Phi}_x^{(i)}(k, l) = \hat{\mathbf{N}}^{(i)}(k, l)\hat{\mathbf{\Lambda}}^{(i)}(k, l)\hat{\mathbf{N}}^{(i), H}(k, l)$  (EVD)
     $\hat{\mathbf{a}}_{\text{ALS}}^{(i)}(k, l) = \sqrt{\hat{\lambda}_1^{(i)}(k, l) / \hat{\phi}_{s, \text{ALS}}^{(i)}(k, l)} \hat{\nu}_1^{(i)}(k, l)$  (20)
  end
   $\hat{\mathbf{a}}_{\text{ALS}}^{(1)}(k, l+1) = \hat{\mathbf{a}}_{\text{ALS}}^{(N)}(k, l) / (\mathbf{e}^T \hat{\mathbf{a}}_{\text{ALS}}^{(N)}(k, l))$  (for next frame)
end
  
```

Simulation results

1. Simulated stationary source (ACE)

- Linear microphone array (M=6, d=6cm)
- Target source at 15° (measured room impulse responses, $T_{60} \approx 1.25$ s)
- Simulated diffuse babble noise (SDR=10 dB)



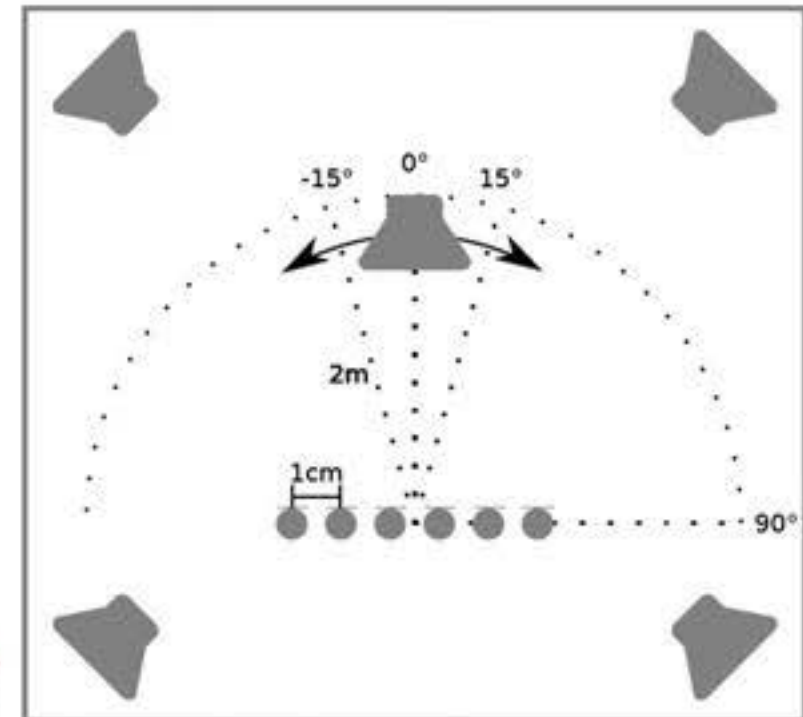
Simulation results

1. Simulated stationary source (ACE)

- Linear microphone array ($M=6$, $d=6\text{cm}$)
- Target source at 15° (measured room impulse responses, $T_{60} \approx 1.25\text{ s}$)
- Simulated diffuse babble noise ($\text{SDR}=10\text{ dB}$)

2. Recorded moving source (varechoic lab)

- Linear microphone array ($M=6$, $d=1\text{cm}$)
- Moving target source ($T_{60} \approx 0.35\text{ s}$)
- Recorded pseudo-diffuse babble noise ($\text{SDR}=10\text{ dB}$)



Simulation results

1. Simulated stationary source (ACE)

- Linear microphone array ($M=6$, $d=6\text{cm}$)
- Target source at 15° (measured room impulse responses, $T_{60} \approx 1.25\text{ s}$)
- Simulated diffuse babble noise ($\text{SDR}=10\text{ dB}$)

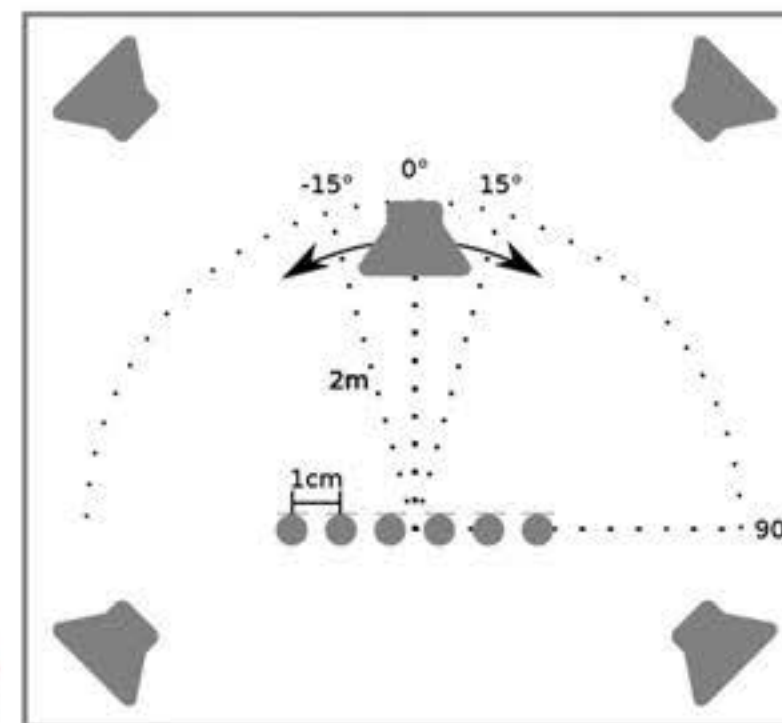
2. Recorded moving source (varechoic lab)

- Linear microphone array ($M=6$, $d=1\text{cm}$)
- Moving target source ($T_{60} \approx 0.35\text{ s}$)
- Recorded pseudo-diffuse babble noise ($\text{SDR}=10\text{ dB}$)



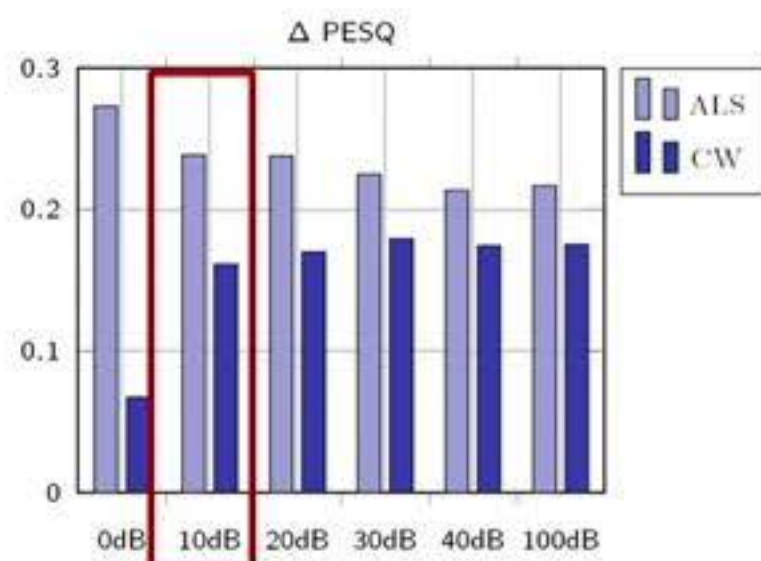
Simulation parameters:

- $f_s = 16\text{ kHz}$, STFT: 64 ms, 75% overlap, Hamming window
- Γ : spherically diffuse; smoothing: 40 ms; speech PSD estimated using decision-directed approach, $G_{\min} = -10\text{ dB}$
- CW: noise covariance matrix estimated during first second; ALS: 5 iterations



Simulation results (PESQ improvement)

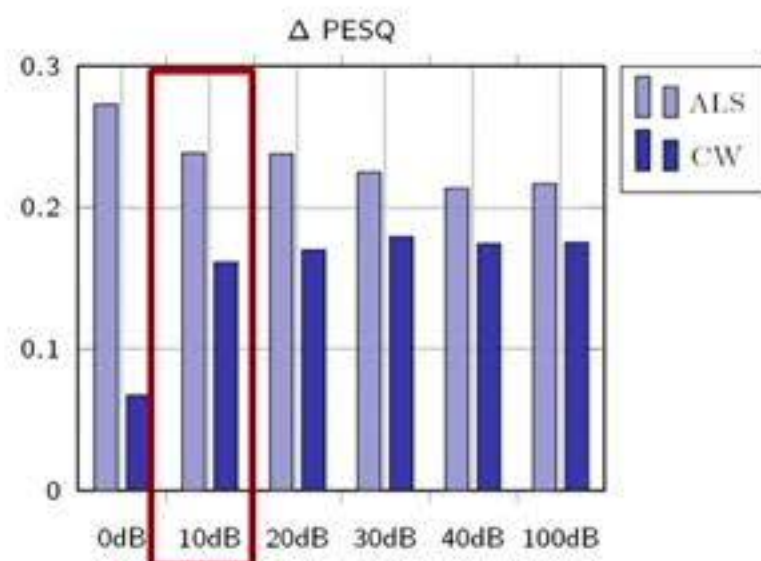
1. Simulated stationary source



Linear array ($M=6$, $d=6\text{cm}$), $f_s=16\text{kHz}$, stationary source at $\theta=15^\circ$, perfectly diffuse babble noise ($\text{SDR}=10\text{dB}$), sensor noise ($\text{DNR}=10\text{dB}$)

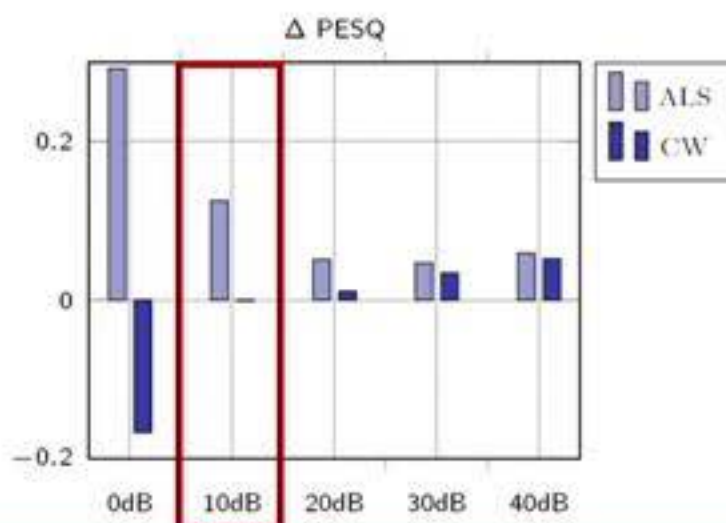
Simulation results (PESQ improvement)

1. Simulated stationary source



Linear array ($M=6$, $d=6\text{cm}$), $f_s=16\text{kHz}$, stationary source at $\theta=15^\circ$, perfectly diffuse babble noise (SDR=10dB), sensor noise (DNR=10dB)

2. Recorded moving source

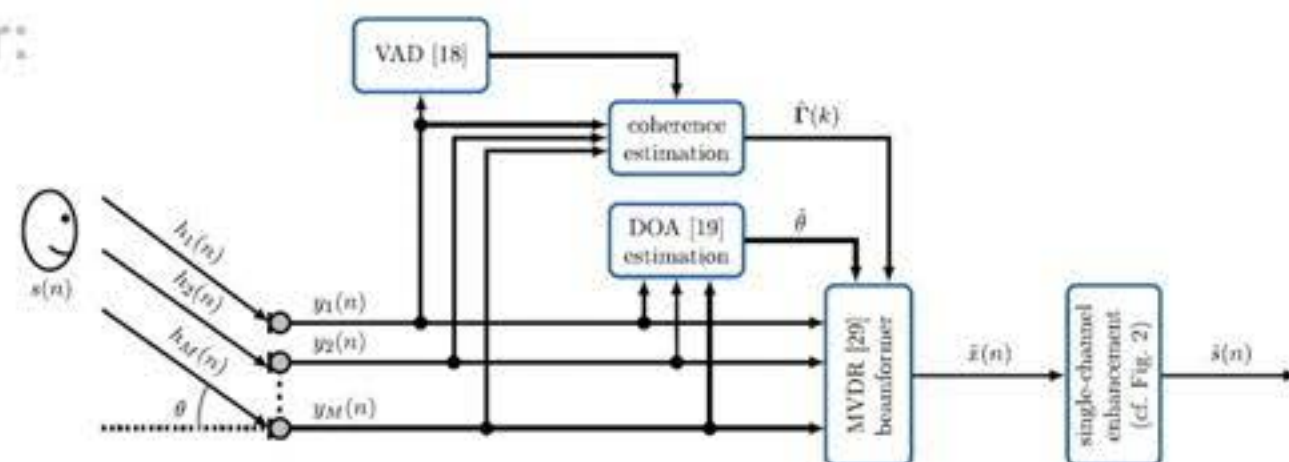


Linear array ($M=6$, $d=1\text{cm}$), $f_s=16\text{kHz}$, moving source $\theta=0^\circ$ to $\theta=90^\circ$ pseudo-diffuse babble noise (SDR=10dB), sensor noise (DNR=10dB)

Multi-microphone dereverberation and noise reduction

1. **Beamforming + spectral postfilter:** *multiply each time-frequency bin with real-valued gain*

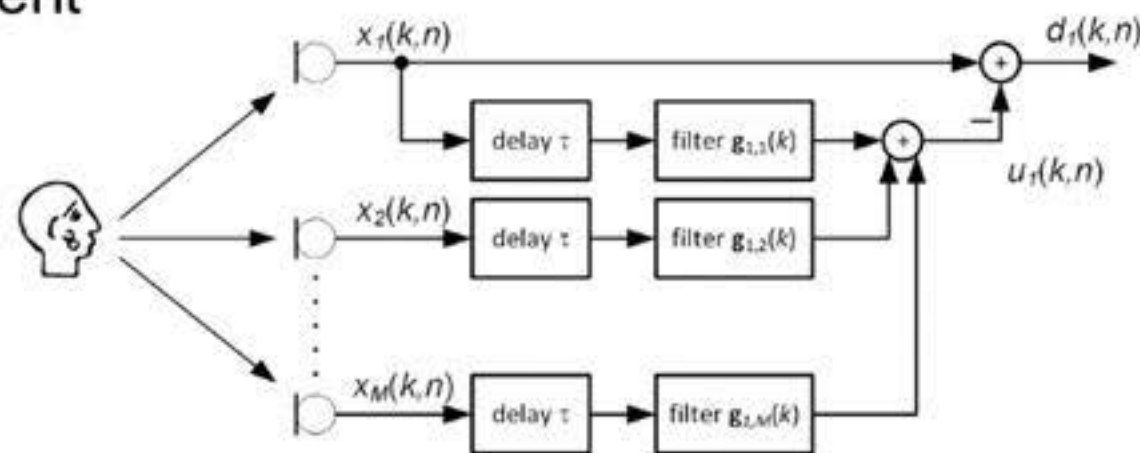
$$\mathbf{y}(l) = \mathbf{a}(l)x_1(l) + \mathbf{x}_r(l) + \mathbf{v}(l)$$



2. **Reverberation and noise suppression:** *subtract complex-valued estimate of late reverberant and noise component*

$$y_m(l) = h_m(l) \star s(l) + v_m(l)$$

$$\hat{x}_{e,1}(l) = y_1(l) - \mathbf{Y}_\tau(l)\mathbf{g}(l)$$

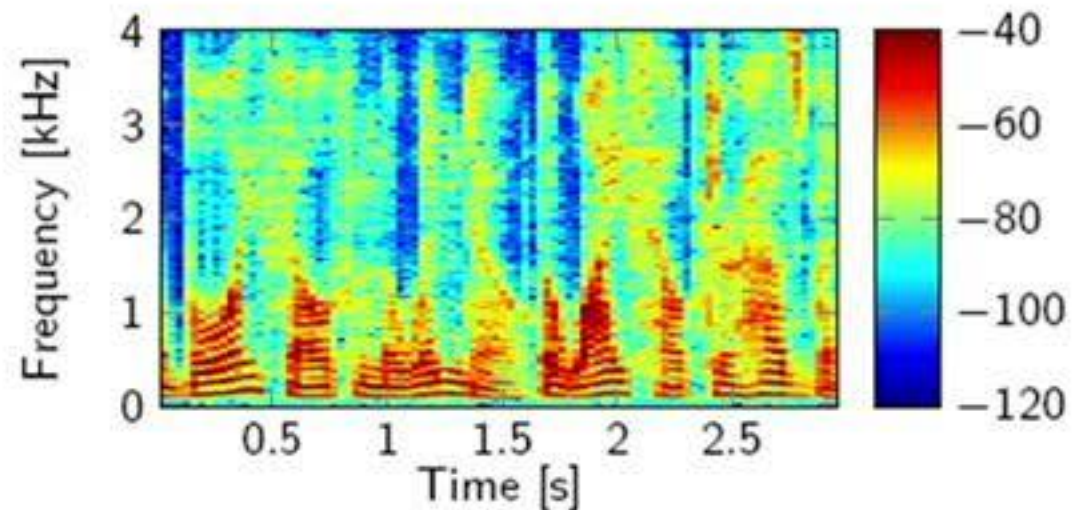


Reverberation suppression

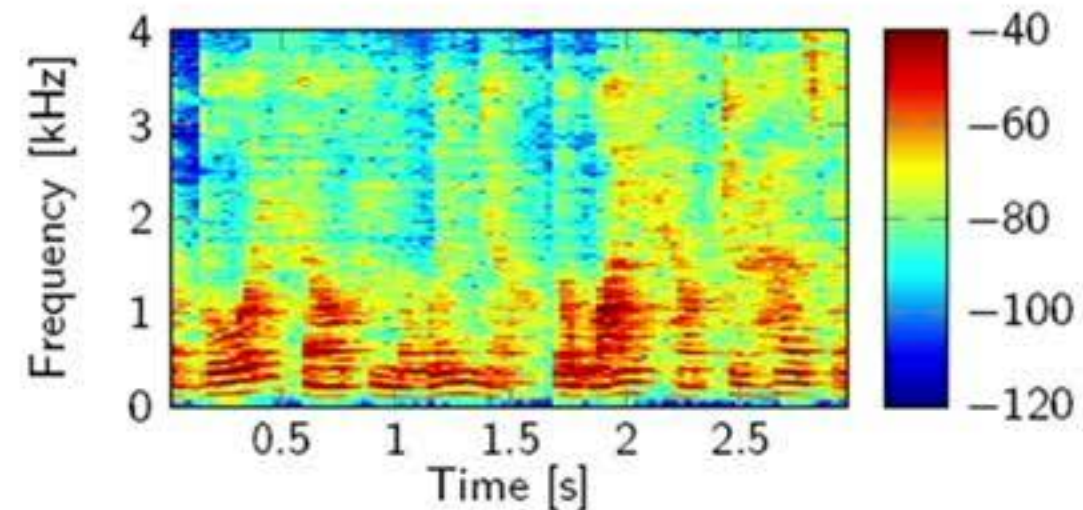
- **Goal:** estimate clean speech STFT coefficients $s(k, l)$ from reverberant (and noisy) STFT coefficients $y_m(k, l)$ by subtracting late reverberant component

$$y_m(k, l) = \underbrace{h_m(k, l) \star s(k, l)}_{x_m(k, l)} + v_m(k, l)$$

- Probabilistic estimation using (statistical) **models of desired speech signal and reverberation**
- Exploit **sparsity properties** of speech in STFT-domain



Clean



Reverberant

Reverberation suppression

- **Goal:** estimate clean speech STFT coefficients $s(k, l)$ from reverberant (and noisy) STFT coefficients $y_m(k, l)$ by subtracting late reverberant component

$$y_m(k, l) = \underbrace{h_m(k, l) \star s(k, l)}_{x_m(k, l)} + v_m(k, l)$$

- Probabilistic estimation using (statistical) **models of desired speech signal and reverberation**
 - Exploit **sparsity properties** of speech in STFT-domain
- **Approach:** transform to equivalent AR model \rightarrow sparse **multi-channel linear prediction (MCLP)**

$$x_1(k, l) = d(k, l) + \sum_{m=1}^M \sum_{n=0}^{L_g-1} g_m(k, n) x_m(k, l - \tau - n)$$

Reverberation suppression

- **Goal:** estimate clean speech STFT coefficients $s(k, l)$ from reverberant (and noisy) STFT coefficients $y_m(k, l)$ by subtracting late reverberant component

$$y_m(k, l) = \underbrace{h_m(k, l) \star s(k, l)}_{x_m(k, l)} + v_m(k, l)$$

- Probabilistic estimation using (statistical) **models of desired speech signal and reverberation**
 - Exploit **sparsity properties** of speech in STFT-domain
- **Approach:** transform to equivalent AR model \rightarrow sparse **multi-channel linear prediction (MCLP)**

$$x_1(k, l) = \underbrace{d(k, l)}_{\substack{\uparrow \\ \text{clean signal} \\ \text{(incl. early reflections)}}} + \sum_{m=1}^M \sum_{n=0}^{L_g-1} g_m(k, n) \underbrace{x_m(k, l - \tau - n)}_{\substack{\uparrow \\ \text{delay} \\ \text{(early reflections)}}$$

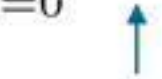
Reverberation suppression

- **Goal:** estimate clean speech STFT coefficients $s(k, l)$ from reverberant (and noisy) STFT coefficients $y_m(k, l)$ by subtracting late reverberant component

$$y_m(k, l) = \underbrace{h_m(k, l) \star s(k, l)}_{x_m(k, l)} + v_m(k, l)$$

- Probabilistic estimation using (statistical) **models of desired speech signal and reverberation**
 - Exploit **sparsity properties** of speech in STFT-domain
- **Approach:** transform to equivalent AR model \rightarrow sparse **multi-channel linear prediction (MCLP)**

$$x_1(k, l) = d(k, l) + \sum_{m=1}^M \sum_{n=0}^{L_g-1} g_m(k, n) x_m(k, l - \tau - n)$$



prediction
filters

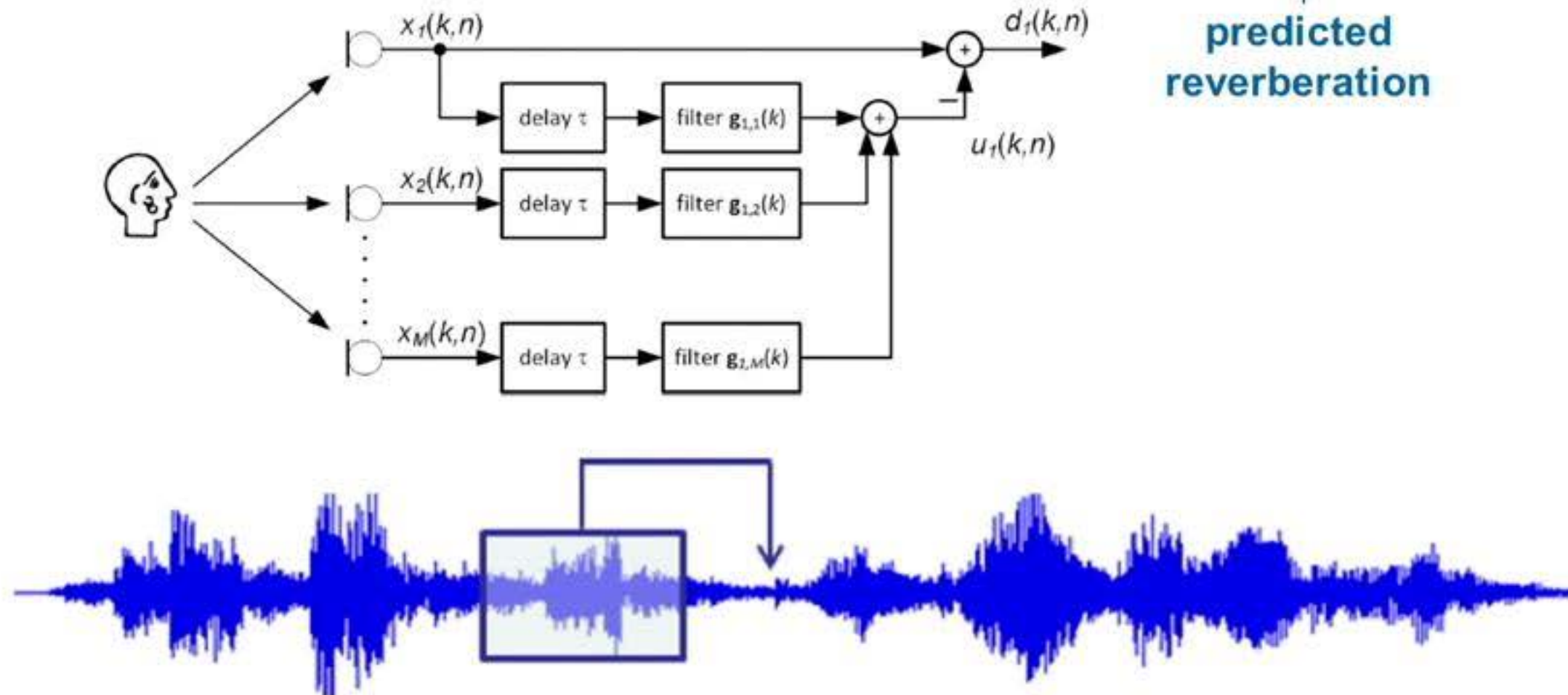
Multi-channel linear prediction

- AR model of reverberant speech

$$\mathbf{x}_1(k) = \mathbf{d}(k) + \mathbf{X}_\tau(k)\mathbf{g}(k)$$

$$\hat{\mathbf{d}}(k) = \mathbf{x}_1(k) - \mathbf{X}_\tau(k)\hat{\mathbf{g}}(k)$$

↑
predicted
reverberation

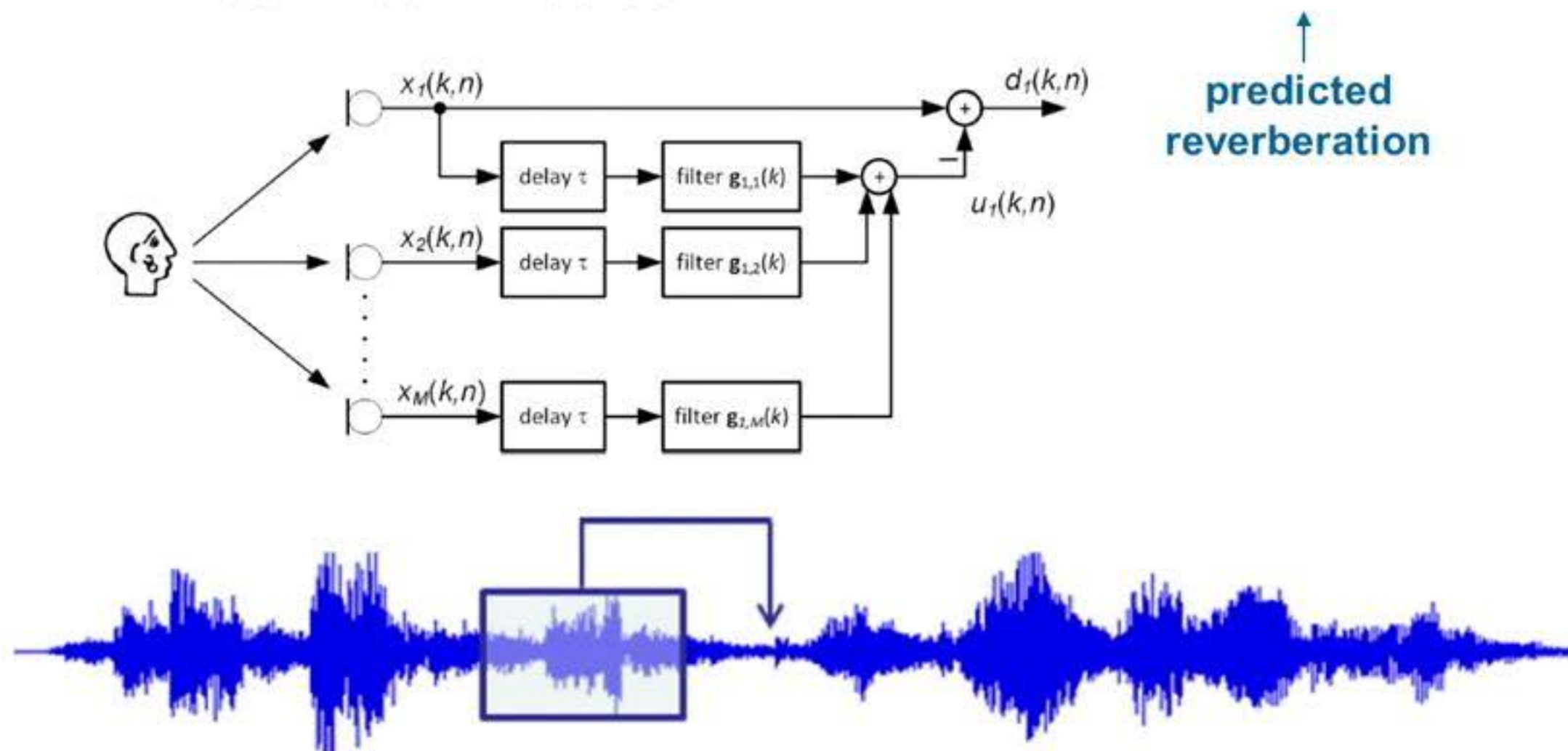


Multi-channel linear prediction

- AR model of reverberant speech

$$\mathbf{x}_1(k) = \mathbf{d}(k) + \mathbf{X}_\tau(k)\mathbf{g}(k)$$

$$\hat{\mathbf{d}}(k) = \mathbf{x}_1(k) - \mathbf{X}_\tau(k)\hat{\mathbf{g}}(k)$$



How to select suitable cost function for prediction filters ?

Multi-channel linear prediction

- **Approach:**

- STFT coefficients of desired signal are modelled using **circular sparse/super-Gaussian prior with time-varying variance** $\lambda(n)$

$$\rho(d(n)) = \max_{\lambda(n) > 0} \mathcal{N}_C(d(n); 0, \lambda(n)) \boxed{\psi(\lambda(n))}$$

Scaling function $\psi(\cdot)$ can be interpreted as **hyper-prior on variance**

Multi-channel linear prediction

- **Approach:**

- STFT coefficients of desired signal are modelled using **circular sparse/super-Gaussian prior with time-varying variance** $\lambda(n)$

$$\rho(d(n)) = \max_{\lambda(n) > 0} \mathcal{N}_C(d(n); 0, \lambda(n)) \boxed{\psi(\lambda(n))}$$

Scaling function $\psi(\cdot)$ can be interpreted as **hyper-prior on variance**

- **Maximum-Likelihood Estimation** (batch, per frequency bin)

$$\mathcal{L}(\mathbf{g}) = \prod_{n=1}^N \rho(d(n)) \Rightarrow \min_{\lambda > 0, \mathbf{g}} \sum_{n=1}^N \left(\frac{|d(n)|^2}{\lambda(n)} + \log \pi \lambda(n) \boxed{-\log \psi(\lambda(n))} \right)$$

Multi-channel linear prediction

- **Approach:**

- STFT coefficients of desired signal are modelled using **circular sparse/super-Gaussian prior with time-varying variance** $\lambda(n)$

$$\rho(d(n)) = \max_{\lambda(n) > 0} \mathcal{N}_C(d(n); 0, \lambda(n)) \boxed{\psi(\lambda(n))}$$

Scaling function $\psi(\cdot)$ can be interpreted as **hyper-prior on variance**

- **Maximum-Likelihood Estimation** (batch, per frequency bin)

$$\mathcal{L}(\mathbf{g}) = \prod_{n=1}^N \rho(d(n)) \Rightarrow \min_{\lambda > 0, \mathbf{g}} \sum_{n=1}^N \left(\frac{|d(n)|^2}{\lambda(n)} + \log \pi \lambda(n) \boxed{-\log \psi(\lambda(n))} \right)$$

- **Alternating optimization procedure**

1. Estimate **prediction vector** (assuming fixed variances)

$$\hat{\mathbf{g}}^{(i+1)} = \left(\mathbf{X}_\tau^H \mathcal{D}_{\hat{\lambda}^{(i)}}^{-1} \mathbf{X}_\tau \right)^{-1} \mathbf{X}_\tau^H \mathcal{D}_{\hat{\lambda}^{(i)}}^{-1} \mathbf{x}_1$$

2. Estimate **variances** (assuming fixed prediction vector)

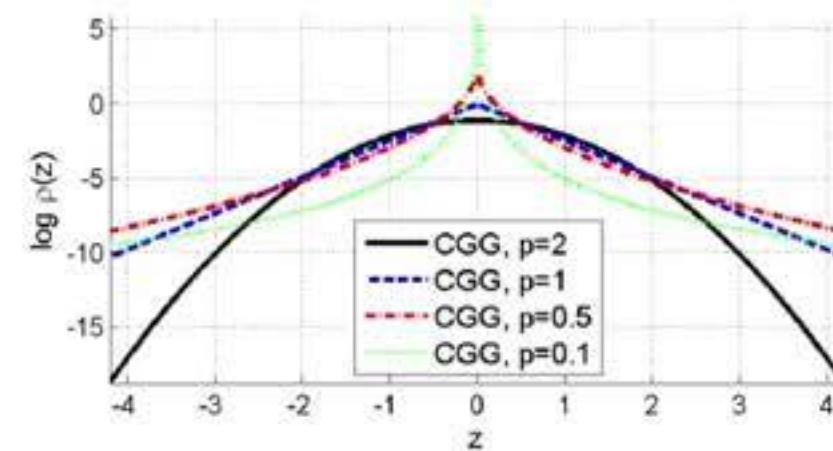
$$\hat{\lambda}^{(i+1)}(n) = \arg \min_{\lambda(n) > 0} \frac{|\hat{d}^{(i+1)}(n)|^2}{\lambda(n)} + \log \pi \lambda(n) \boxed{-\log \psi(\lambda(n))}$$

Multi-channel linear prediction

- Example:** complex generalized Gaussian (CGG) prior with shape parameter p

$$\rho(z) = \frac{p}{2\pi\gamma\Gamma(2/p)} e^{-\frac{|z|^p}{\gamma^{p/2}}}$$

$$\hat{\lambda}^{(i+1)}(n) = |\hat{d}^{(i+1)}(n)|^{2-p},$$

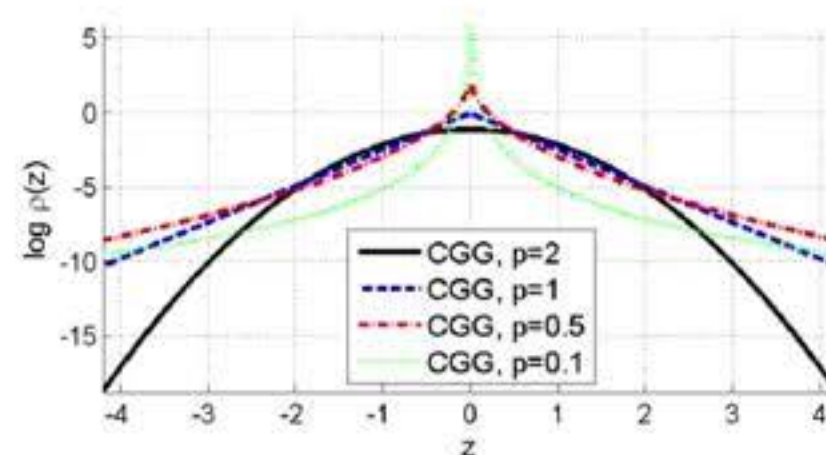


Multi-channel linear prediction

- Example:** complex generalized Gaussian (CGG) prior with shape parameter p

$$\rho(z) = \frac{p}{2\pi\gamma\Gamma(2/p)} e^{-\frac{|z|^p}{\gamma^{p/2}}}$$

$$\hat{\lambda}^{(i+1)}(n) = |\hat{d}^{(i+1)}(n)|^{2-p},$$

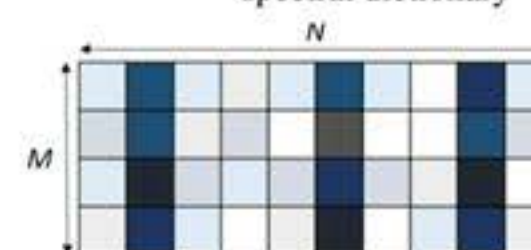


- Remarks:**

- ML estimation using CGG prior is equivalent to l_p -norm minimization
→ **promotes sparsity of TF-coefficients across time** (for $p < 2$)
- Incorporate additional knowledge of speech signal,
e.g. **low-rank structure** (NMF)
- Group sparsity** for MIMO speech dereverberation
→ mixed norms
- Recursive version** by constraining MCLP-based
estimate of undesired component

$$\min_{\mathbf{g}} \|\mathbf{d}\|_p^p,$$

$$|\mathbf{D}|^2 \approx \underbrace{\mathbf{W}}_{\text{spectral dictionary}} \mathbf{H}$$



ℓ_2 norm of the columns

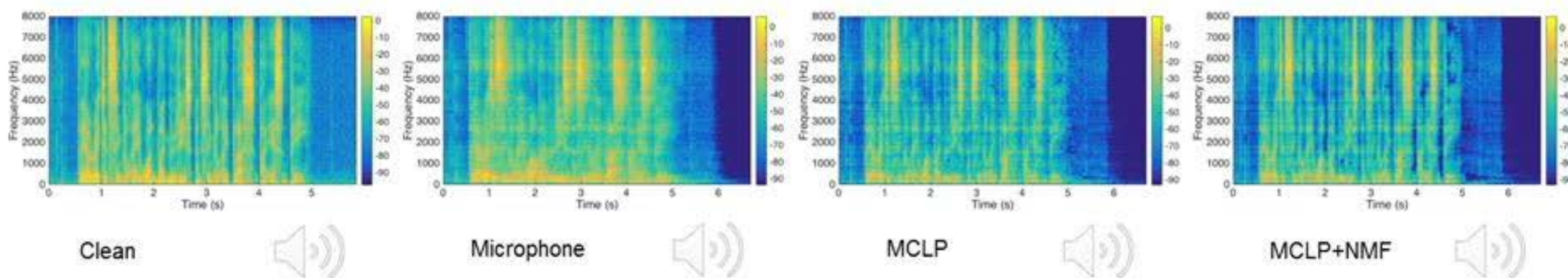
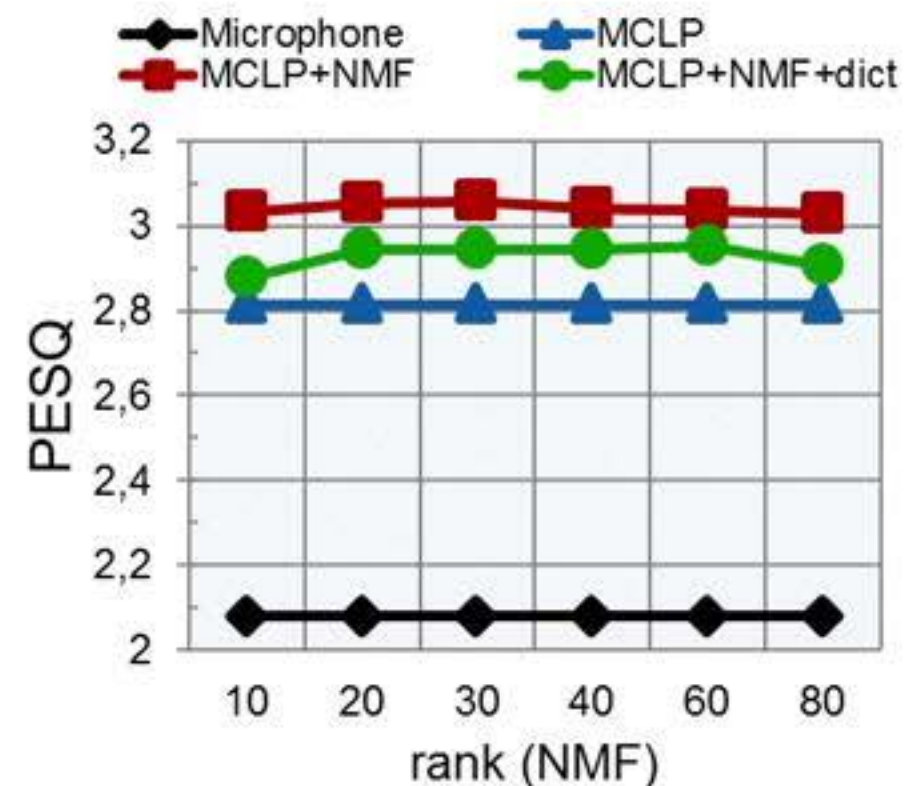


$\|\mathbf{D}\|_{2,p} = \ell_p$ norm of the vector

Multi-channel linear prediction

- Instrumental validation (noiseless, batch)**

- MCLP exploits sparsity
- NMF introduces speech structure (unsupervised vs. supervised NMF)



$T_{60} \sim 700\text{ms}$, $M=4$, $f_s=16\text{ kHz}$; STFT: 64ms (overlap 16ms); MCLP: $L_0=8$, $\tau=2$, $p=0$

Current/future work

- Estimation of RETF vectors and PSDs for **multi-speaker scenarios** (e.g. based on Procrustes problem)
- **Joint noise reduction and dereverberation:** integration of multi-channel linear prediction and generalized sidelobe canceller

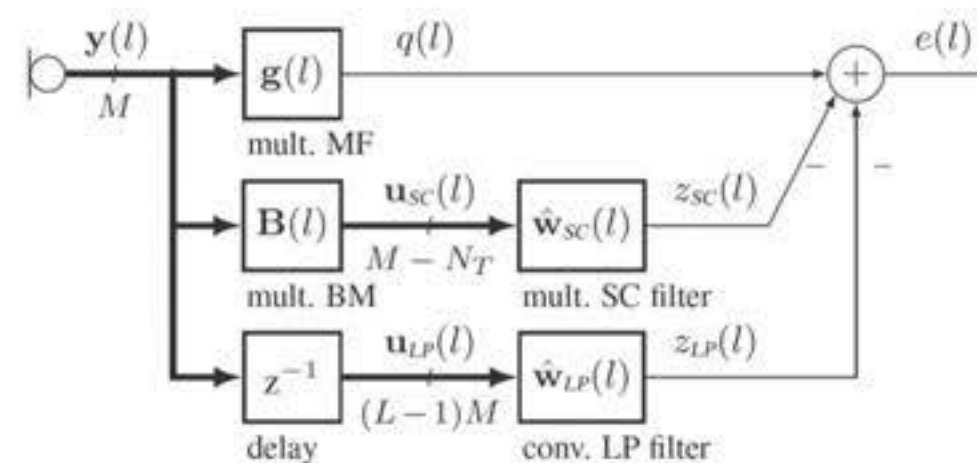
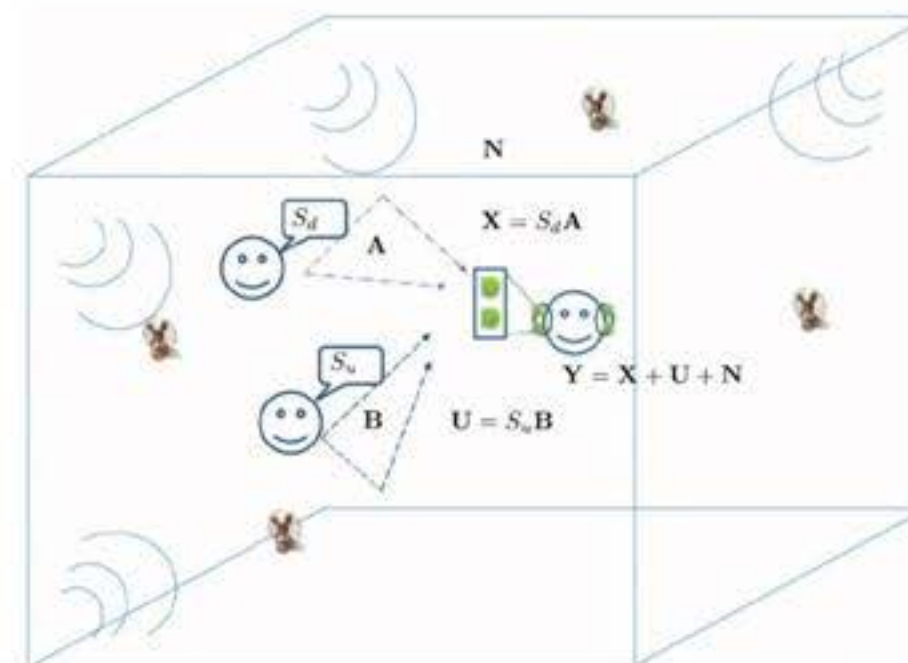


Fig. 1. The integrated sidelobe cancellation and linear prediction (ISCLP) architecture.

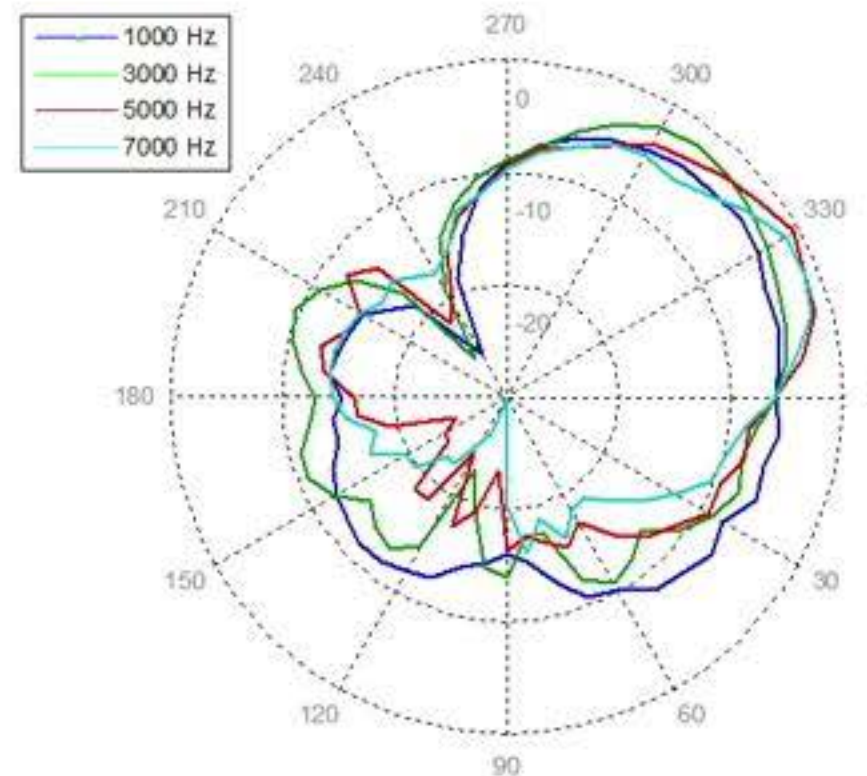
2. Acoustic signal processing for binaural hearing devices

Hearing devices / hearables

- ❑ Hearing devices generally have multiple microphones available and allow for advanced acoustical signal pre-processing



Monaural (2-3)



Hearing devices / hearables

- ❑ Hearing devices generally have multiple microphones available and allow for advanced acoustical signal pre-processing

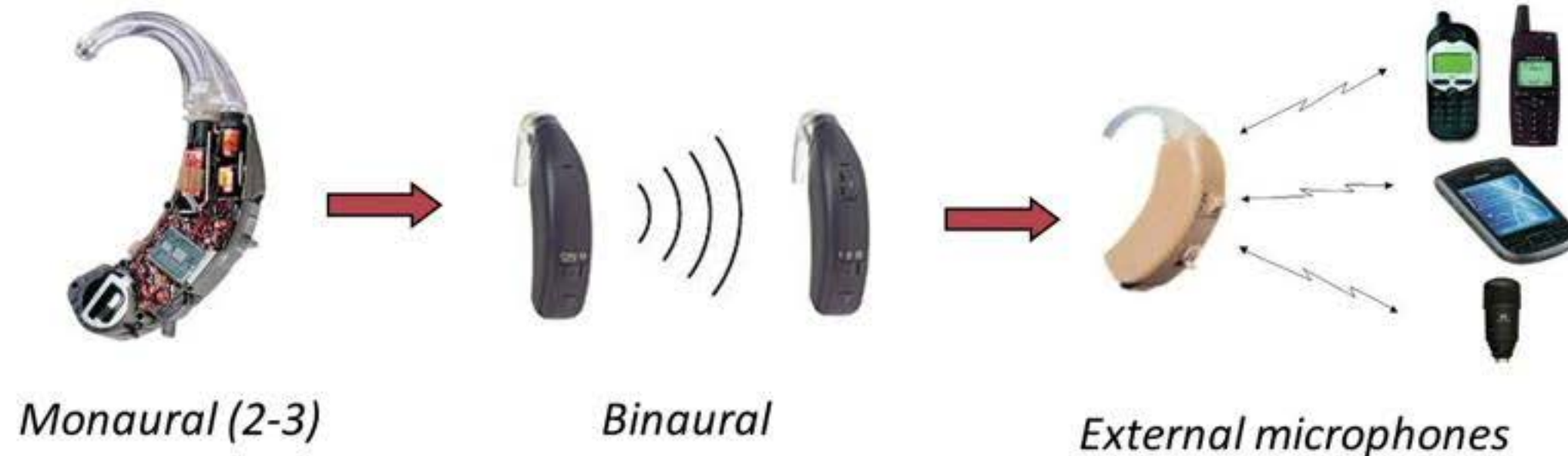


Monaural (2-3)

Binaural

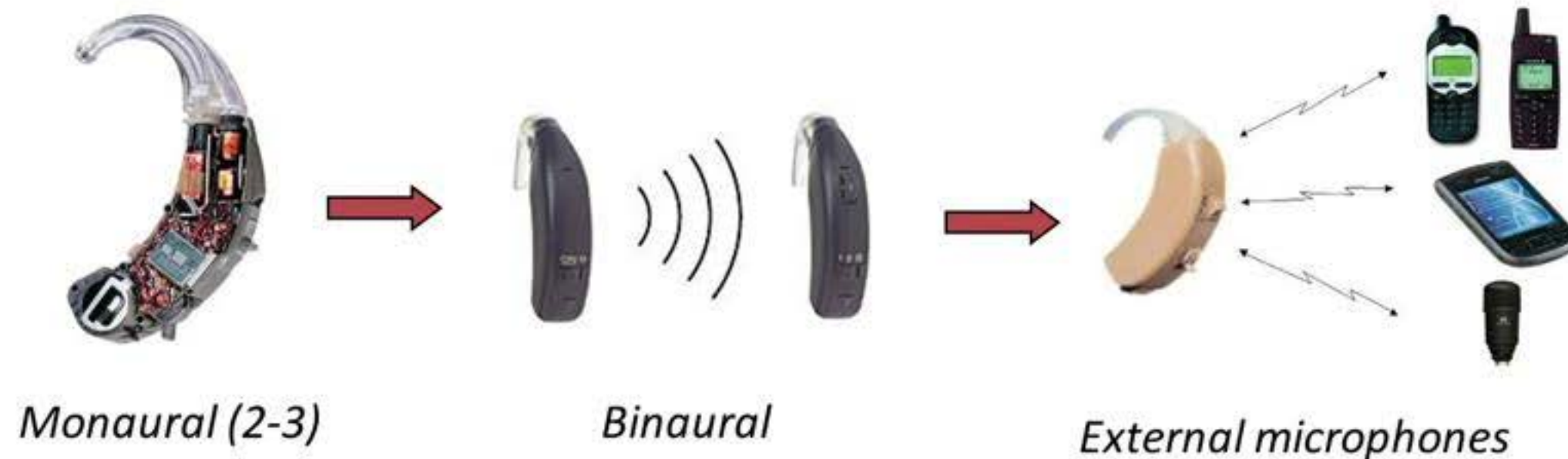
Hearing devices / hearables

- ❑ Hearing devices generally have multiple microphones available and allow for advanced acoustical signal pre-processing



Hearing devices / hearables

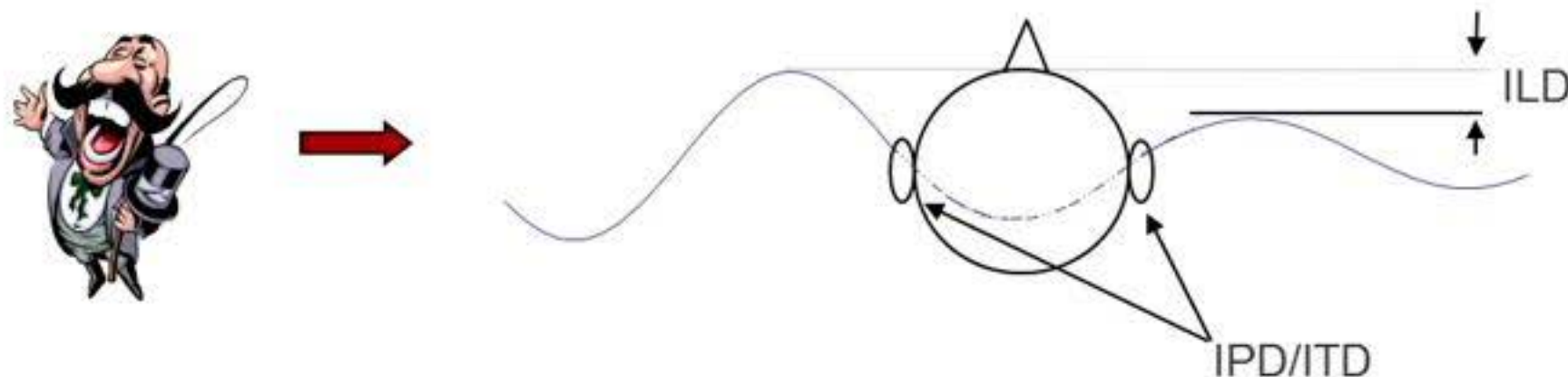
- ❑ Hearing devices generally have multiple microphones available and allow for **advanced acoustical signal pre-processing**



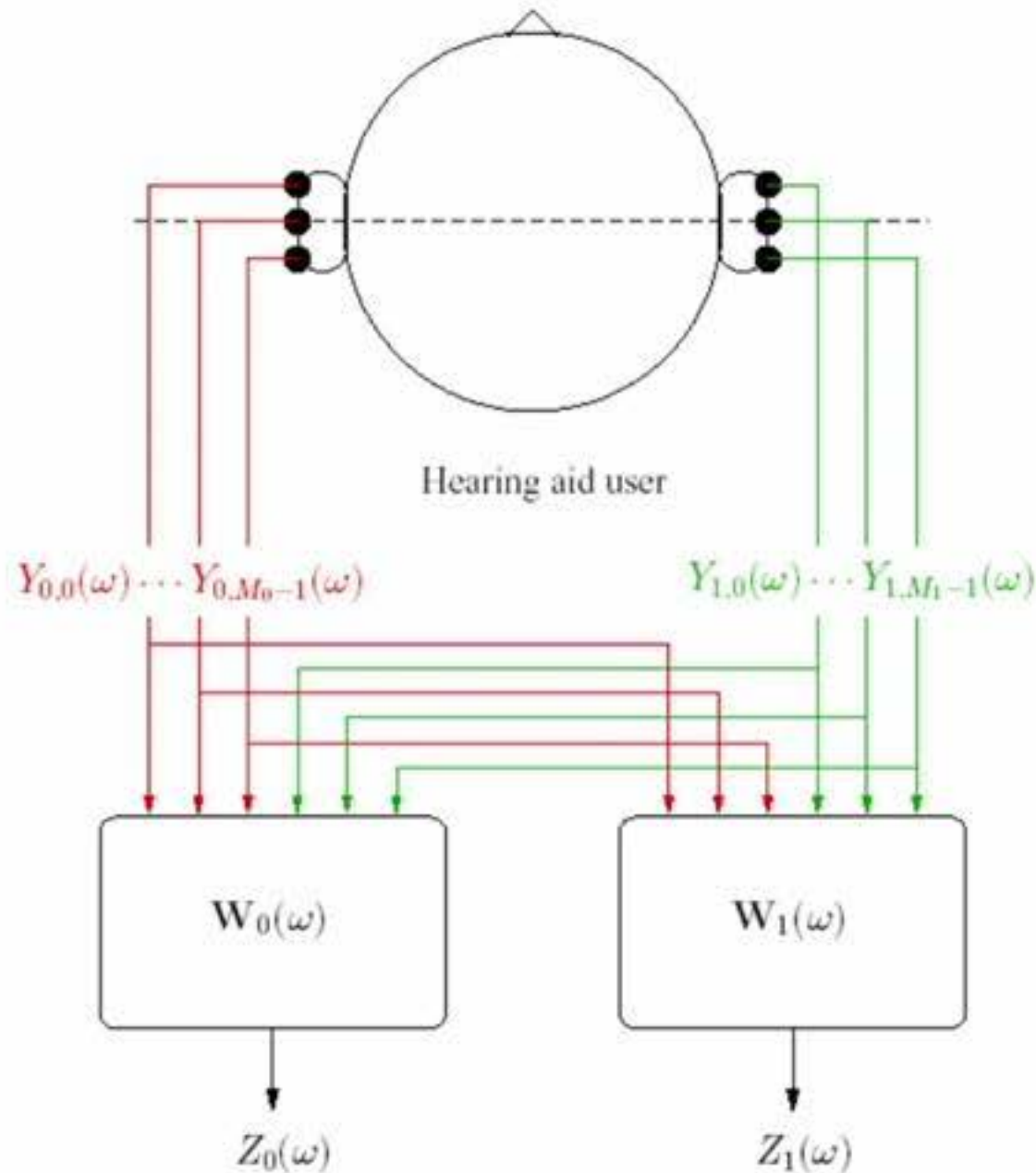
- ❑ **Main objectives of binaural speech enhancement algorithms:** improve speech intelligibility + **preserve spatial awareness (binaural cues)**

Binaural auditory cues

- ❑ **Interaural Time/Phase Difference (ITD/IPD)**
Interaural Level Difference (ILD)
Interaural Coherence (IC)
 - ❑ ITD: $f < 1500$ Hz, ILD: $f > 2000$ Hz
 - ❑ IC: describes spatial characteristics, e.g. perceived width, of diffuse noise, and determines when ITD/ILD cues are *reliable*
- ❑ Binaural cues, in addition to spectro-temporal cues, play an important role in auditory scene analysis (source segregation) and speech intelligibility

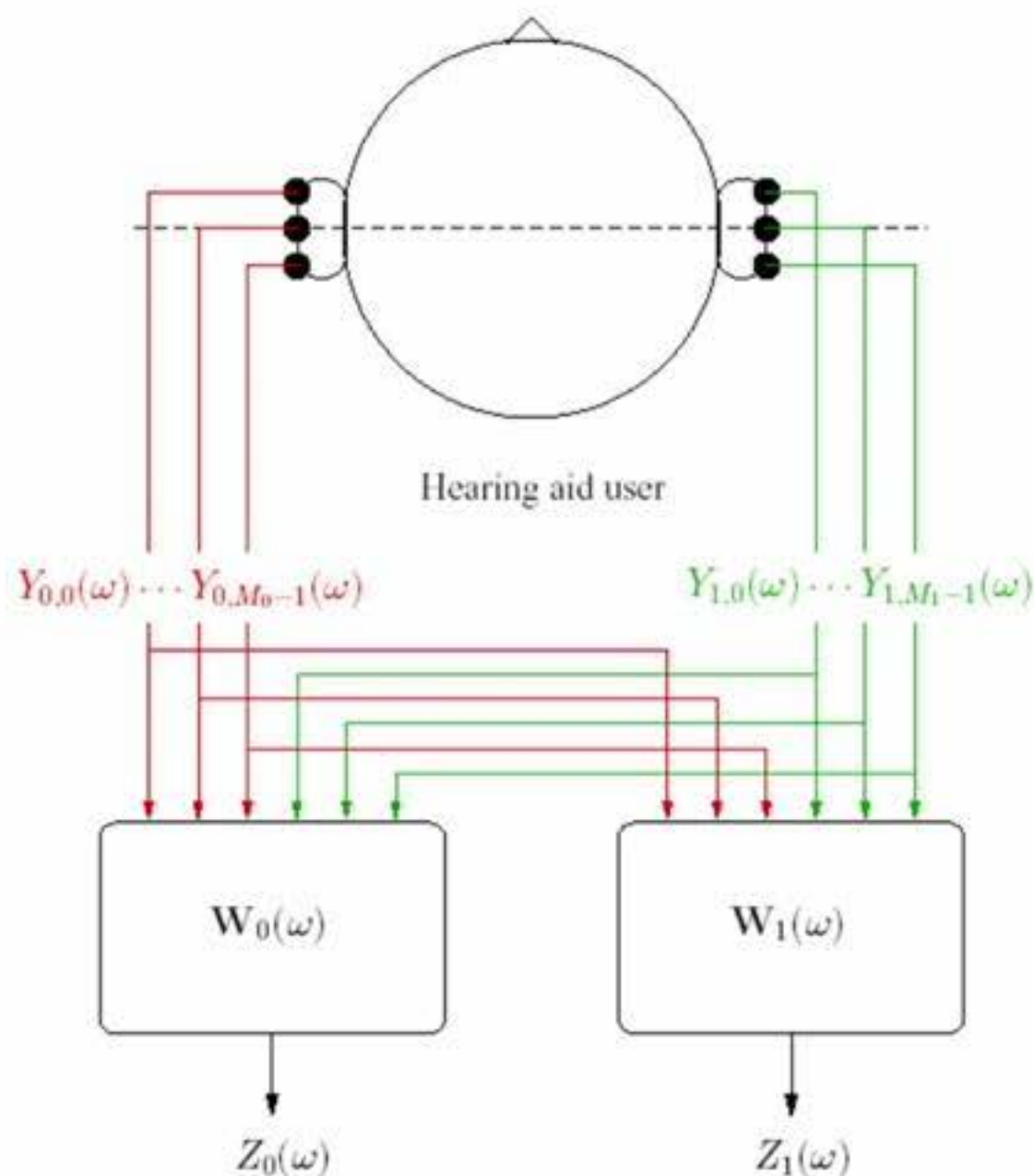


Binaural noise reduction: Configuration



- Binaural hearing aid configuration:
 - Two hearing aids with in total M microphones
 - All microphone signals \mathbf{Y} are assumed to be available at both hearing aids (perfect wireless link)

Binaural noise reduction: Configuration



□ Binaural hearing aid configuration:

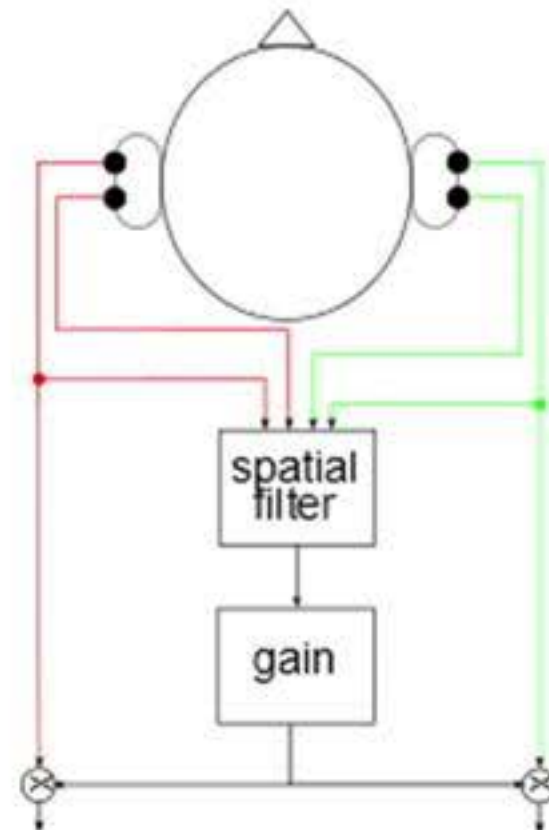
- Two hearing aids with in total M microphones
- All microphone signals \mathbf{Y} are assumed to be available at both hearing aids (perfect wireless link)
- Apply a filter \mathbf{W}_0 and \mathbf{W}_1 at the left and the right hearing aid, generating binaural output signals Z_0 and Z_1

$$Z_0(\omega) = \mathbf{W}_0^H(\omega)\mathbf{Y}(\omega), \quad Z_1(\omega) = \mathbf{W}_1^H(\omega)\mathbf{Y}(\omega)$$

Binaural noise reduction: Two main paradigms

Spectral post-filtering (based on multi-microphone noise reduction)

[Wittkop 2003, Lotter 2006, Rohdenburg 2008, Grimm 2009, Kamkar-Parsi 2011, Reindl 2013, Baumgärtel 2015, Enzner 2016]

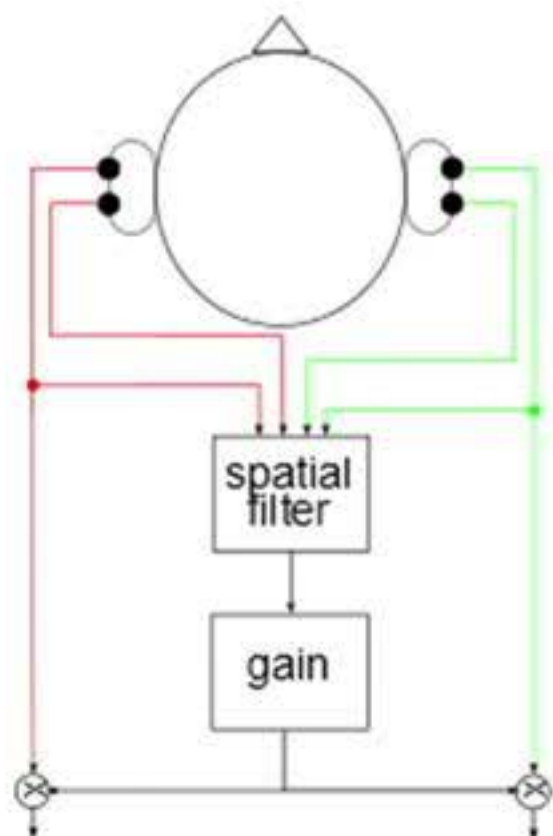


- ⊕ Binaural cue preservation
- ⊖ Possible single-channel artifacts

Binaural noise reduction: Two main paradigms

Spectral post-filtering (based on multi-microphone noise reduction)

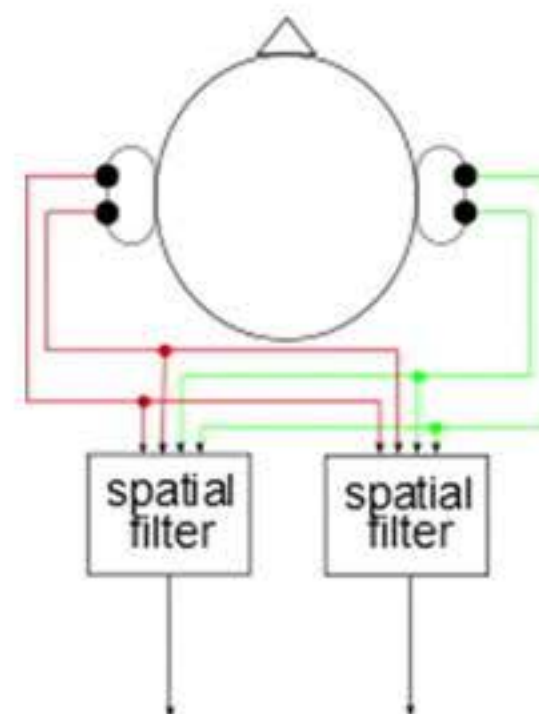
[Wittkop 2003, Lotter 2006, Rohdenburg 2008, Grimm 2009, Kamkar-Parsi 2011, Reindl 2013, Baumgärtel 2015, Enzner 2016]



- ⊕ Binaural cue preservation
- ⊖ Possible single-channel artifacts

Binaural spatial filtering techniques

[Welker 1997, Aichner 2007, Doclo 2010, Cornelis 2012, Hadad 2015-2016, Marquardt 2015-2018, Koutrouvelis 2017-2019]



- ⊕ Larger noise reduction performance
- ⊕ Merge spatial and spectral post-filtering
- ⊖ Binaural cue preservation not guaranteed

Binaural MVDR and MWF

Minimum-Variance-Distortionless-Response (MVDR) beamformer

Goal: minimize output noise power without distorting speech component in reference microphone signals

$$\begin{array}{ll} \min_{\mathbf{W}_0} \mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_0 & \text{subject to } \mathbf{W}_0^H \mathbf{A} = A_0 \\ \min_{\mathbf{W}_1} \mathbf{W}_1^H \mathbf{R}_v \mathbf{W}_1 & \text{subject to } \mathbf{W}_1^H \mathbf{A} = A_1 \end{array}$$

↑
↑
 noise reduction distortionless constraint

Requires estimate/model of noise coherence matrix (e.g. diffuse) and estimate/model of relative transfer function (RTF) of target speech source

Multi-channel Wiener Filter (MWF)

Goal: estimate speech component in reference microphone signals + trade off noise reduction and speech distortion

$$J_{\text{MWF}}(\mathbf{W}) = \mathcal{E} \left\{ \left\| \begin{bmatrix} X_0 - \mathbf{W}_0^H \mathbf{X} \\ X_1 - \mathbf{W}_1^H \mathbf{X} \end{bmatrix} \right\|^2 + \mu \left\| \begin{bmatrix} \mathbf{W}_0^H \mathbf{V} \\ \mathbf{W}_1^H \mathbf{V} \end{bmatrix} \right\|^2 \right\}$$

↑
↑
 speech distortion noise reduction

Requires estimate of speech and noise covariance matrices, e.g. based on SPP

Can be decomposed as binaural MVDR beamformer and spectral postfilter

Binaural MVDR and MWF

Minimum-Variance-Distortionless-Response (MVDR) beamformer

Goal: minimize output noise power without distorting speech component in reference microphone signals

$$\begin{aligned} \min_{\mathbf{W}_0} \mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_0 \quad & \text{subject to} \quad \mathbf{W}_0^H \mathbf{A} = A_0 \\ \min_{\mathbf{W}_1} \mathbf{W}_1^H \mathbf{R}_v \mathbf{W}_1 \quad & \text{subject to} \quad \mathbf{W}_1^H \mathbf{A} = A_1 \end{aligned}$$

↑
↑
 noise reduction distortionless constraint

Requires estimate/model of noise coherence matrix (e.g. diffuse) and estimate/model of relative transfer function (RTF) of target speech source

Multi-channel Wiener Filter (MWF)

Goal: estimate speech component in reference microphone signals + trade off noise reduction and speech distortion

$$J_{\text{MWF}}(\mathbf{W}) = \mathcal{E} \left\{ \left\| \begin{bmatrix} X_0 - \mathbf{W}_0^H \mathbf{X} \\ X_1 - \mathbf{W}_1^H \mathbf{X} \end{bmatrix} \right\|^2 + \mu \left\| \begin{bmatrix} \mathbf{W}_0^H \mathbf{V} \\ \mathbf{W}_1^H \mathbf{V} \end{bmatrix} \right\|^2 \right\}$$

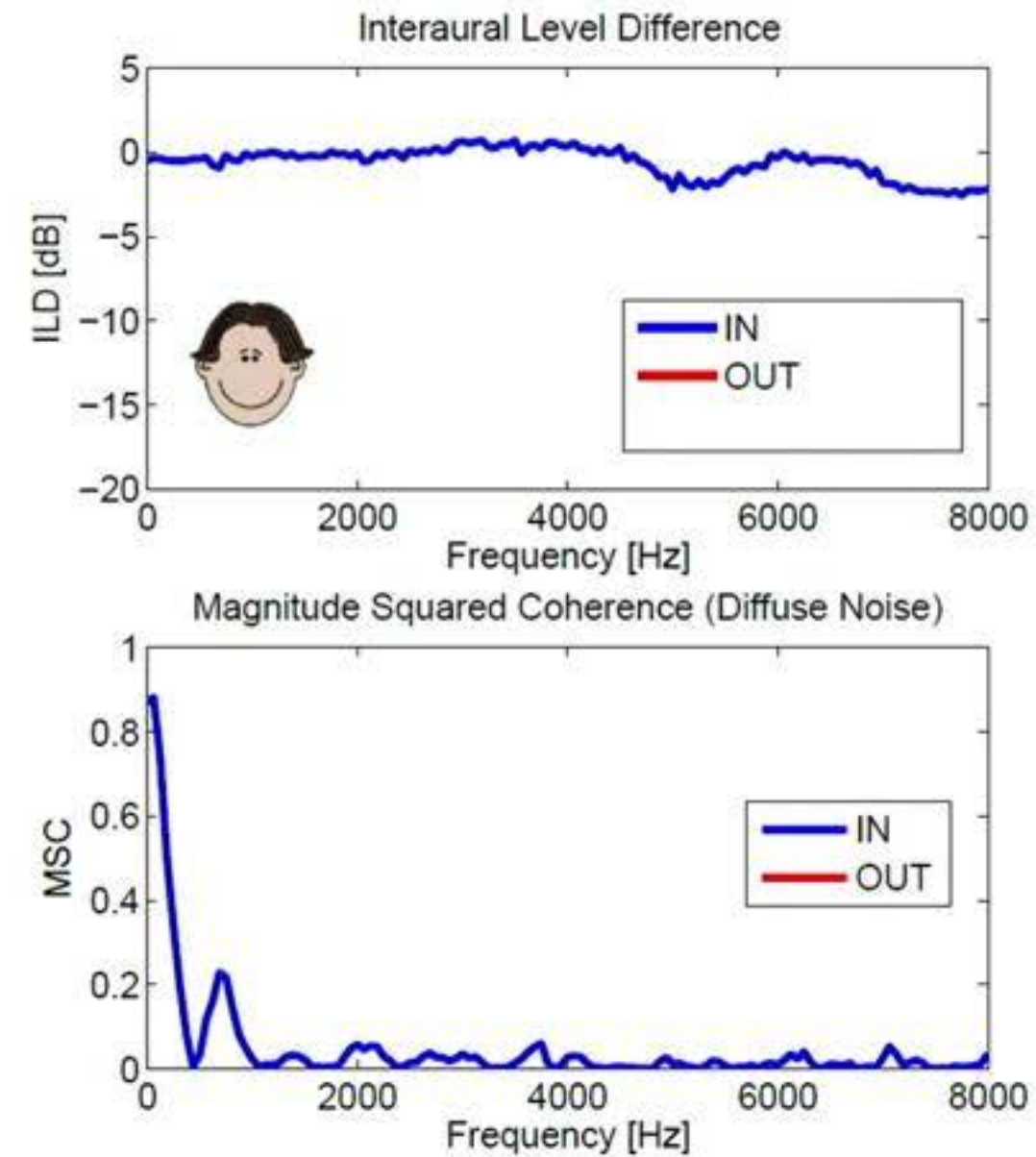
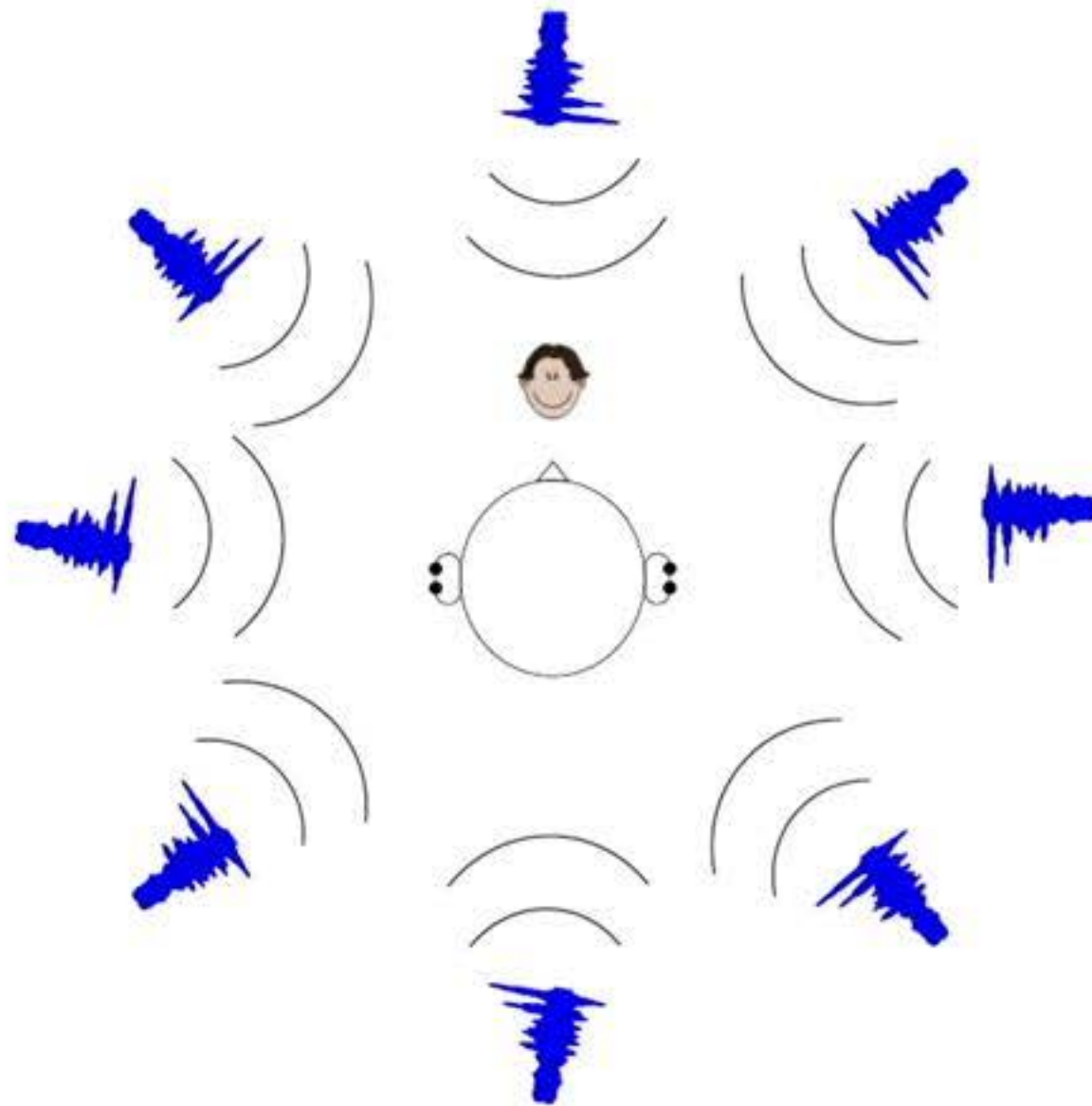
↑
↑
 speech distortion noise reduction

Requires estimate of speech and noise covariance matrices, e.g. based on SPP

Can be decomposed as binaural MVDR beamformer and spectral postfilter

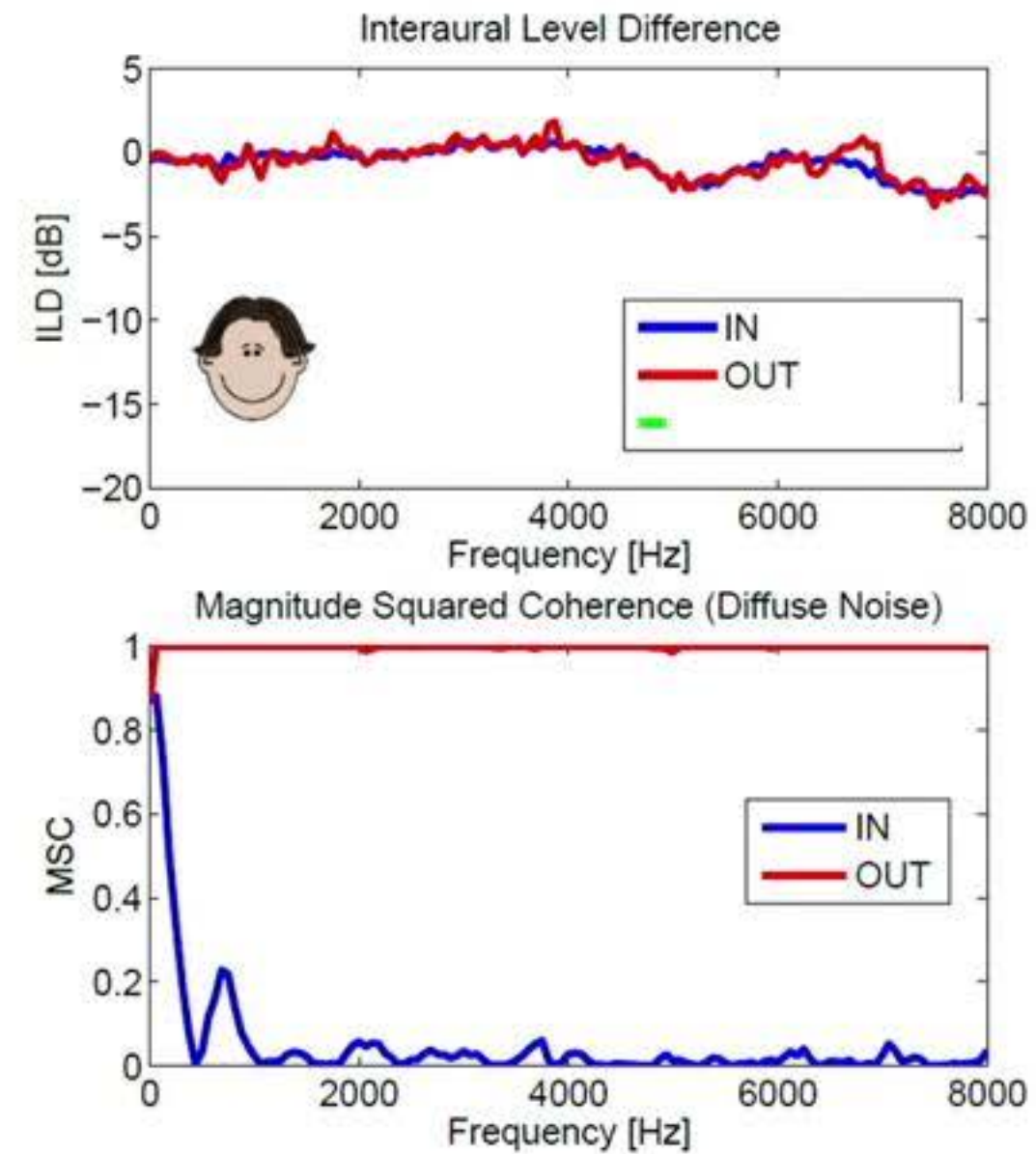
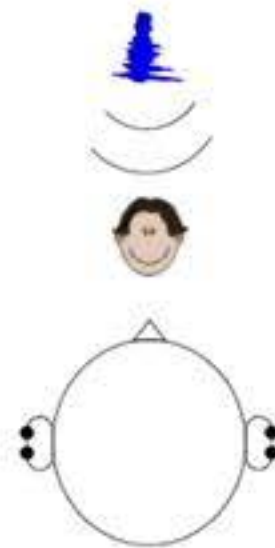
Good noise reduction performance, what about binaural cues ?

Binaural MVDR/MWF: binaural cues

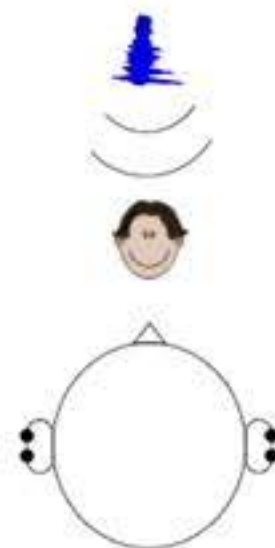


Note: MSC = Magnitude Squared Coherence

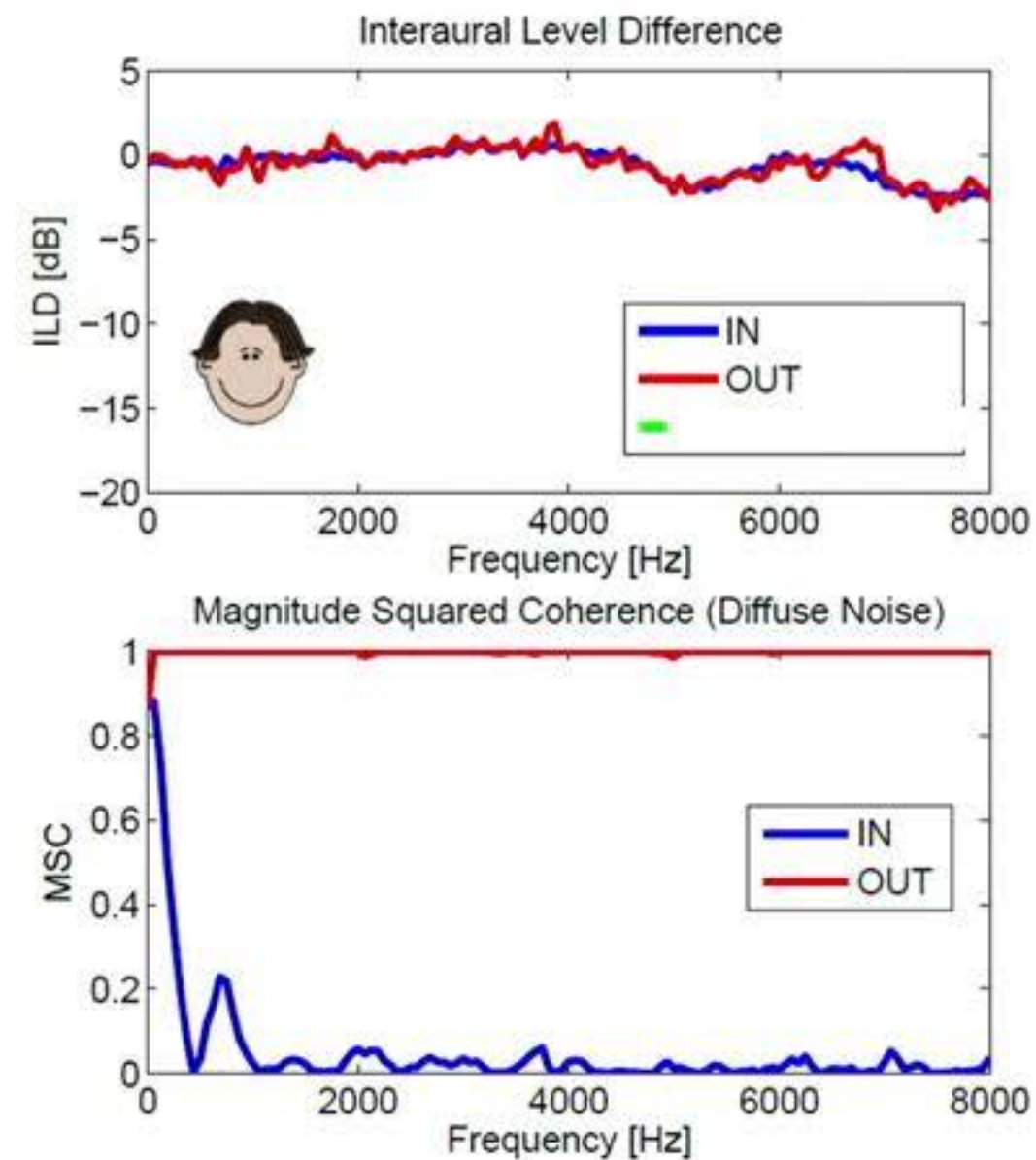
Binaural MVDR/MWF: binaural cues



Binaural MVDR/MWF: binaural cues



**Binaural cues for residual
noise/interference in
binaural MVDR/MWF
not preserved**



Binaural MWF: Extensions for diffuse noise

Binaural MWF

- ⊕ SNR improvement
- ⊕ Binaural cues of speech source
- ⊖ Binaural cues of noise

Binaural MWF: Extensions for diffuse noise

Binaural MWF

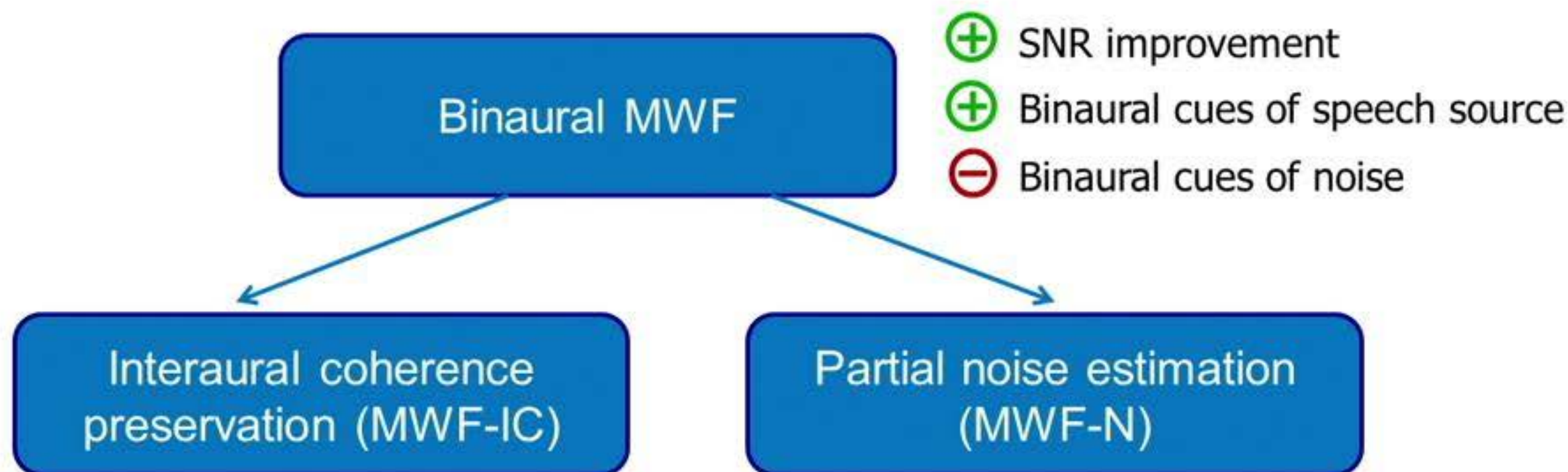
- ⊕ SNR improvement
- ⊕ Binaural cues of speech source
- ⊖ Binaural cues of noise

Interaural coherence
preservation (MWF-IC)

$$J_{MWF-IC}(\mathbf{W}) = J_{MWF}(\mathbf{W}) + \lambda \left| \frac{\mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_1}{\sqrt{\mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_0 \mathbf{W}_1^H \mathbf{R}_v \mathbf{W}_1}} - IC_v^{des} \right|^2$$

- ⊖ No closed-form solution, iterative optimization procedures required

Binaural MWF: Extensions for diffuse noise



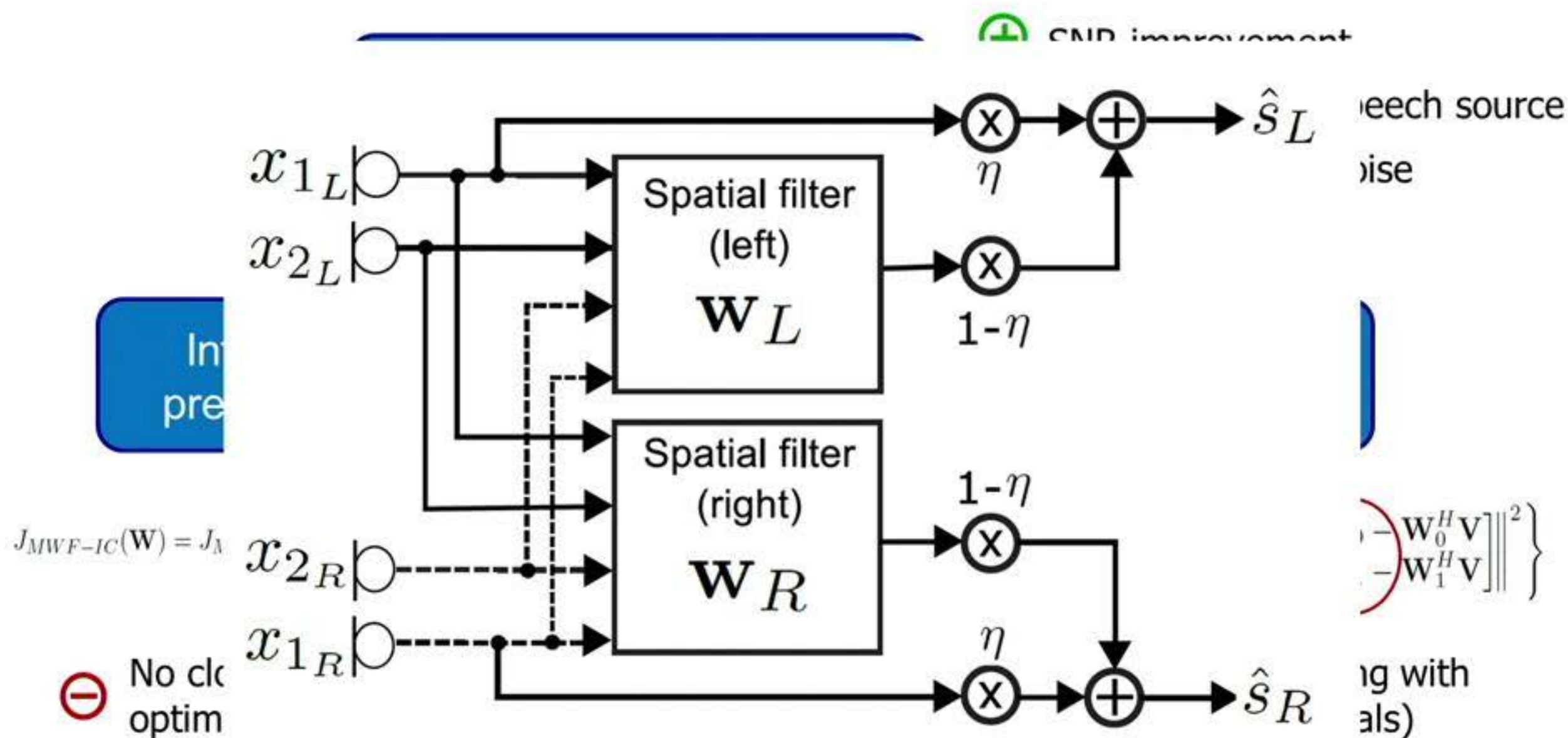
$$J_{MWF-IC}(\mathbf{W}) = J_{MWF}(\mathbf{W}) + \lambda \left| \frac{\mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_1}{\sqrt{\mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_0 \mathbf{W}_1^H \mathbf{R}_v \mathbf{W}_1}} - IC_v^{des} \right|^2$$

⊖ No closed-form solution, iterative optimization procedures required

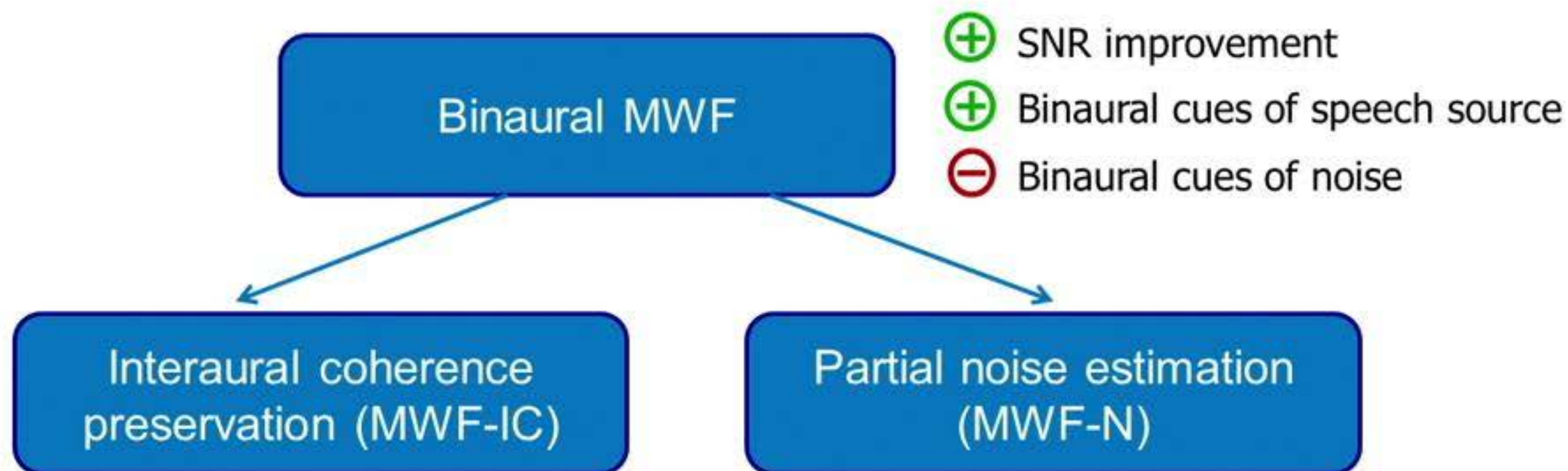
$$J_{MWF-N}(\mathbf{W}) = \mathcal{E} \left\{ \left\| \begin{bmatrix} X_0 - \mathbf{W}_0^H \mathbf{X} \\ X_1 - \mathbf{W}_1^H \mathbf{X} \end{bmatrix} \right\|^2 + \mu \left\| \begin{bmatrix} \eta V_0 - \mathbf{W}_0^H \mathbf{V} \\ \eta V_1 - \mathbf{W}_1^H \mathbf{V} \end{bmatrix} \right\|^2 \right\}$$

⊕ Closed-form solution (mixing with reference microphone signals)

Binaural MWF: Extensions for diffuse noise



Binaural MWF: Extensions for diffuse noise



$$J_{MWF-IC}(\mathbf{W}) = J_{MWF}(\mathbf{W}) + \lambda \left| \frac{\mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_1}{\sqrt{\mathbf{W}_0^H \mathbf{R}_v \mathbf{W}_0 \mathbf{W}_1^H \mathbf{R}_v \mathbf{W}_1}} - IC_v^{des} \right|^2$$

⊖ No closed-form solution, iterative optimization procedures required

$$J_{MWF-N}(\mathbf{W}) = \mathcal{E} \left\{ \left\| \begin{bmatrix} X_0 - \mathbf{W}_0^H \mathbf{X} \\ X_1 - \mathbf{W}_1^H \mathbf{X} \end{bmatrix} \right\|^2 + \mu \left\| \begin{bmatrix} \eta V_0 - \mathbf{W}_0^H \mathbf{V} \\ \eta V_1 - \mathbf{W}_1^H \mathbf{V} \end{bmatrix} \right\|^2 \right\}$$

⊕ Closed-form solution (mixing with reference microphone signals)

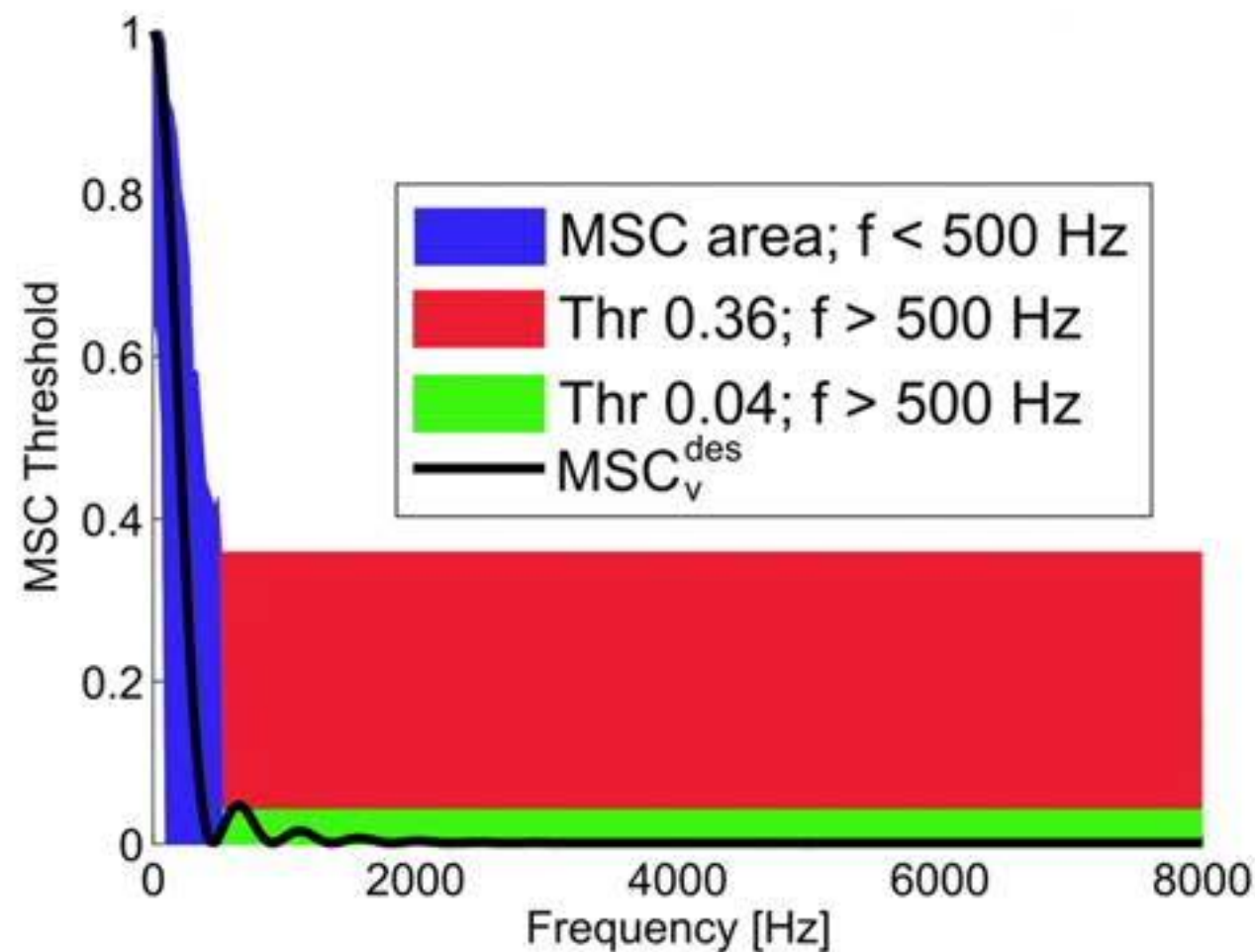
⚖ **Trade-off** between SNR improvement and binaural cue preservation, depending on **parameters** (η and λ)

Trade-off parameters for binaural MVDR/MWF

- ❑ **Fixed broadband values** ($\eta = 0.1 \dots 0.3$)

Trade-off parameters for binaural MVDR/MWF

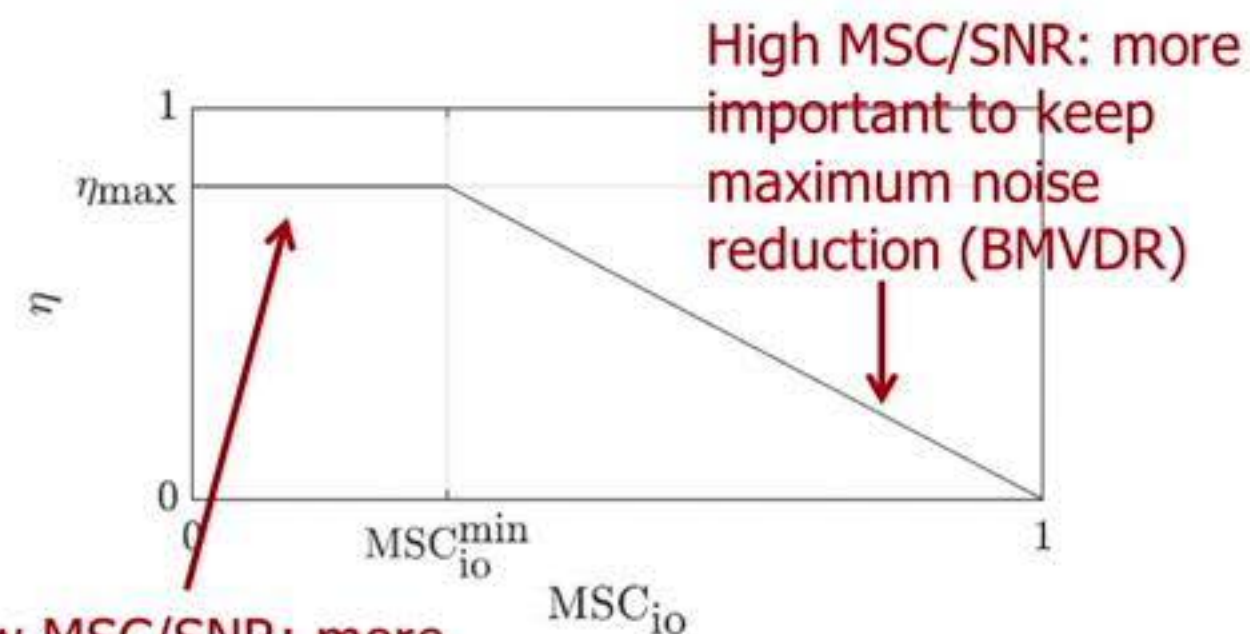
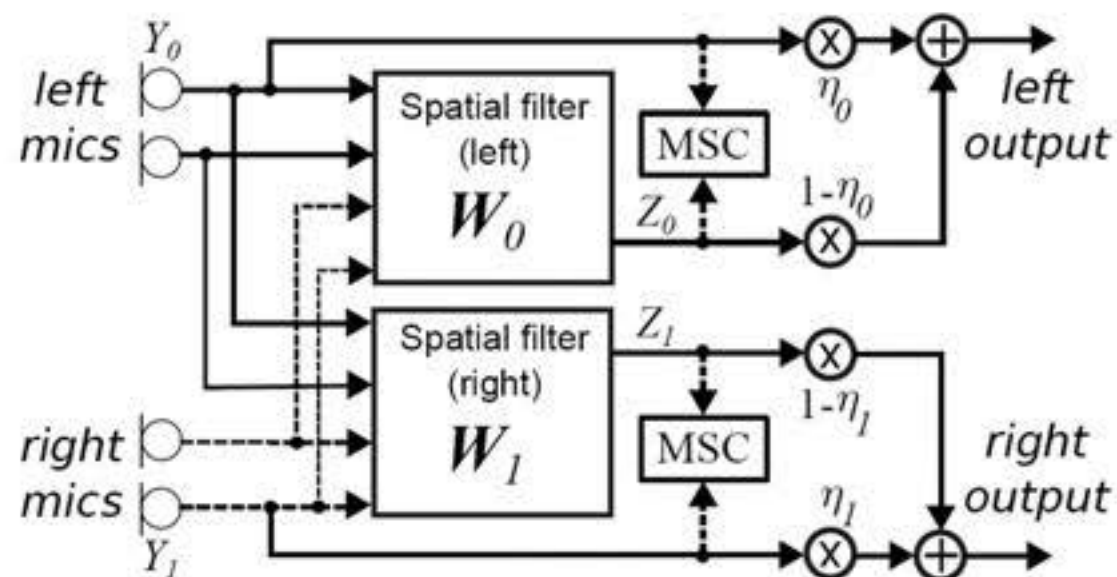
- ❑ **Fixed broadband values** ($\eta = 0.1 \dots 0.3$)
- ❑ **Frequency-dependent values** based on IC discrimination ability of human auditory system



- IC discrimination ability depends on magnitude of reference IC
- **Boundaries on Magnitude Squared Coherence** ($MSC = |IC|^2$) :
 - For $f < 500$ Hz ("large" IC): frequency-dependent MSC boundaries (**blue**)
 - For $f > 500$ Hz ("small" IC): fixed MSC boundary, e.g. 0.36 (**red**) or 0.04 (**green**)

Trade-off parameters for binaural MVDR/MWF

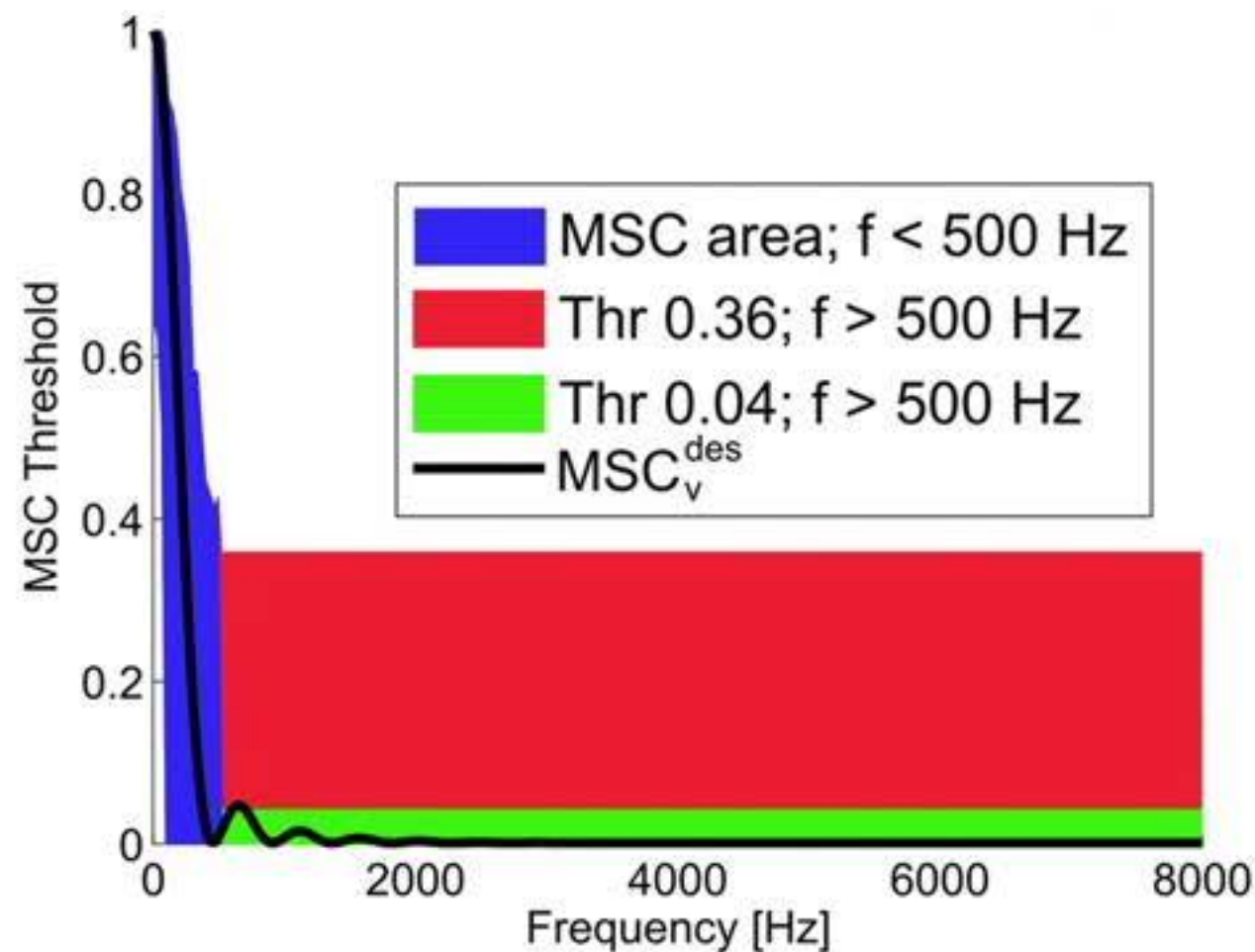
- ❑ **Fixed broadband values** ($\eta = 0.1 \dots 0.3$)
- ❑ **Frequency-dependent values** based on IC discrimination ability of human auditory system
- ❑ **Frequency-dependent function** of MSC between noisy reference microphone signals and output signals of BMVDR beamformer



Low MSC/SNR: more important to preserve binaural cues (scaled input signals)

Trade-off parameters for binaural MVDR/MWF

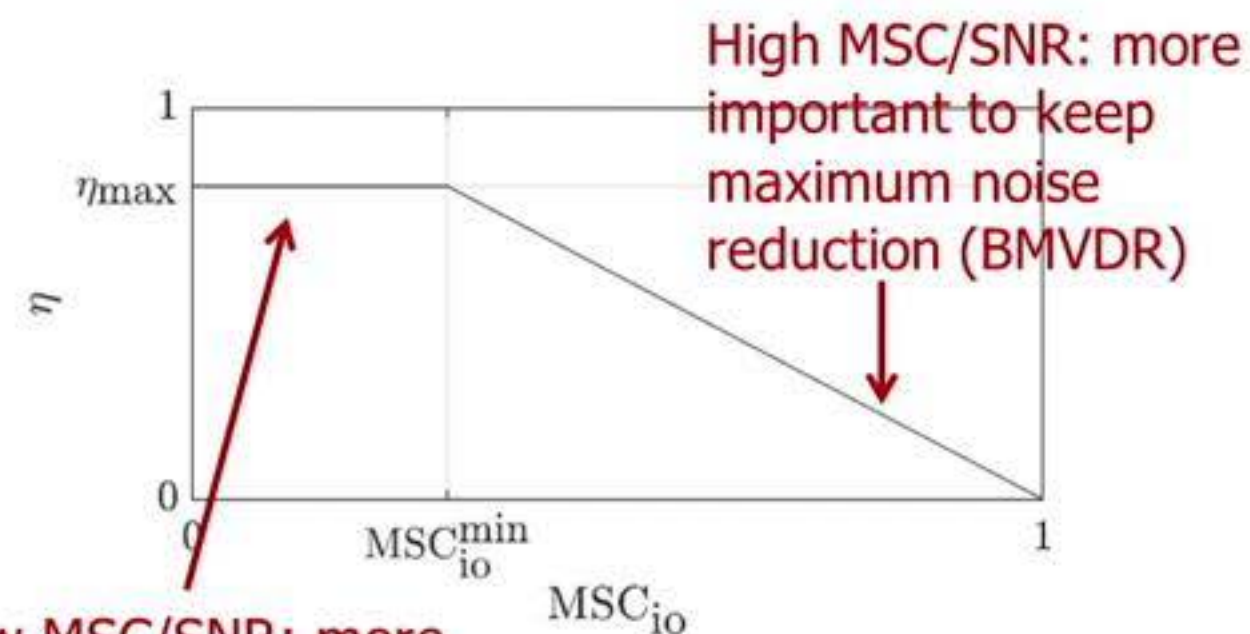
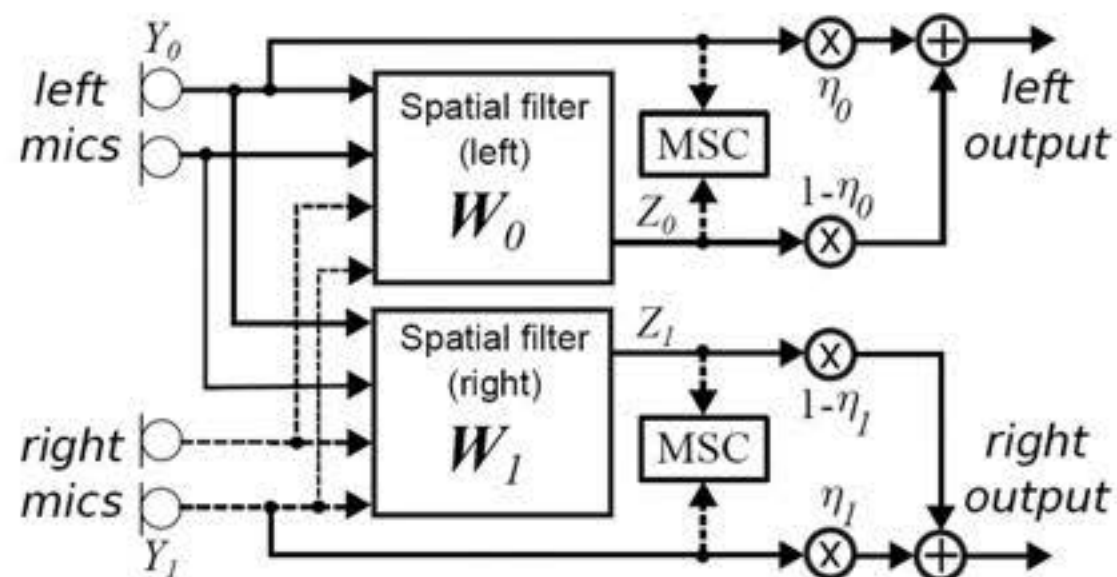
- ❑ **Fixed broadband values** ($\eta = 0.1 \dots 0.3$)
- ❑ **Frequency-dependent values** based on IC discrimination ability of human auditory system



- IC discrimination ability depends on magnitude of reference IC
- **Boundaries on Magnitude Squared Coherence** ($MSC = |IC|^2$) :
 - For $f < 500$ Hz ("large" IC): frequency-dependent MSC boundaries (**blue**)
 - For $f > 500$ Hz ("small" IC): fixed MSC boundary, e.g. 0.36 (**red**) or 0.04 (**green**)

Trade-off parameters for binaural MVDR/MWF

- ❑ **Fixed broadband values** ($\eta = 0.1 \dots 0.3$)
- ❑ **Frequency-dependent values** based on IC discrimination ability of human auditory system
- ❑ **Frequency-dependent function** of MSC between noisy reference microphone signals and output signals of BMVDR beamformer

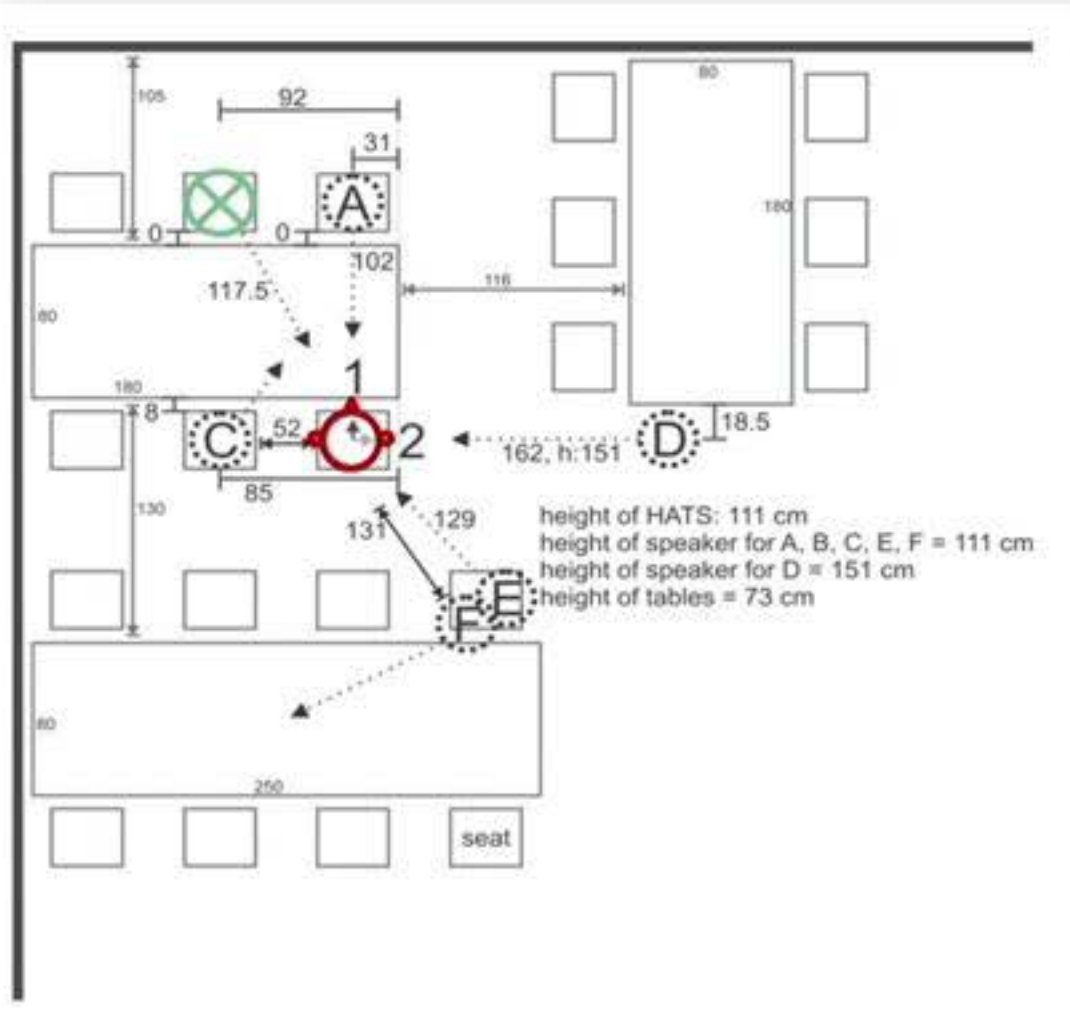


Evaluation: Test setup



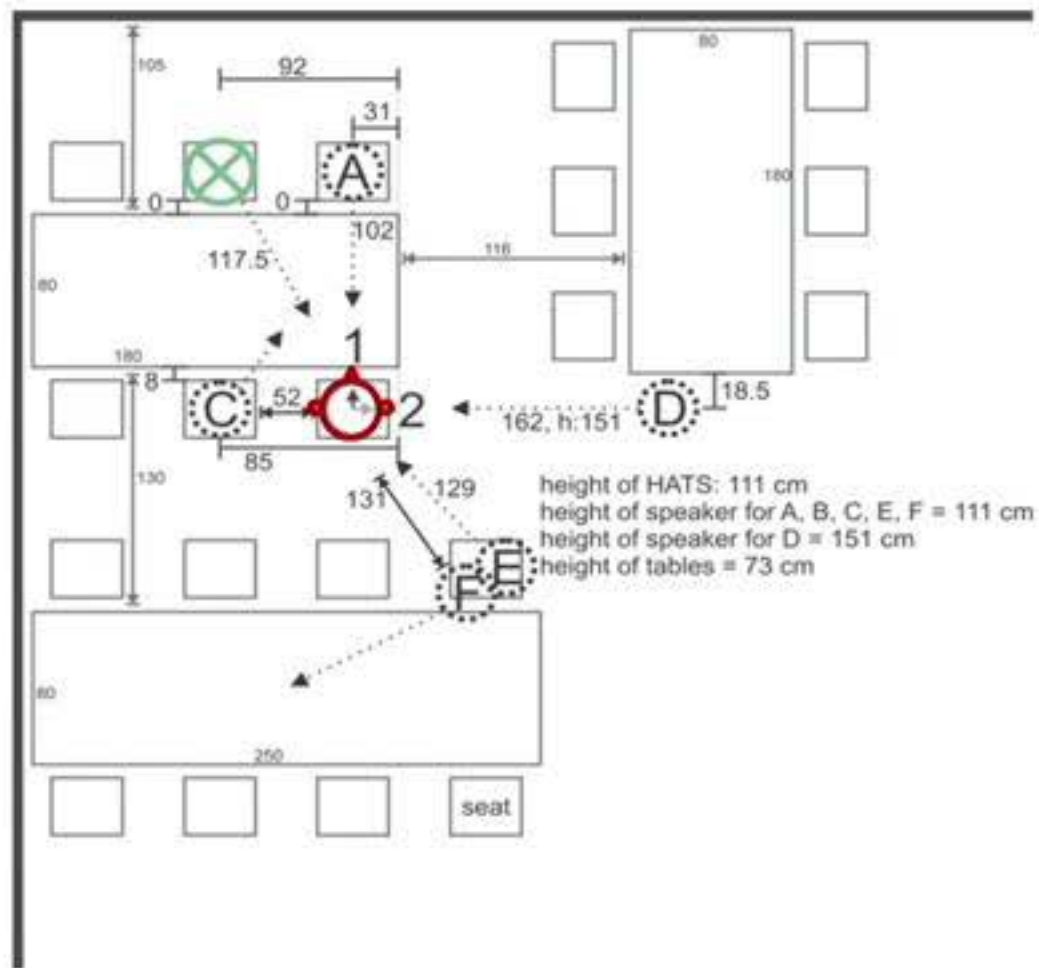
- Binaural hearing aid recordings ($M=4$ mics) in **cafeteria** ($T_{60} \approx 1250$ ms)
 - Target speaker at -35°
 - Realistic cafeteria ambient noise
- **Algorithms:** binaural MVDR and binaural MVDR-N with different **trade-off parameters**:
 - MVDR-IC
 - MVDR-MS1: $\eta_{\max}=0.7$, $MSC_{\min}=0$
 - MVDR-MS2: $\eta_{\max}=1.0$, $MSC_{\min}=0.1$

Evaluation: Test setup



- Binaural hearing aid recordings ($M=4$ mics) in **cafeteria** ($T_{60} \approx 1250$ ms)
 - Target speaker at -35°
 - Realistic cafeteria ambient noise
- **Algorithms:** binaural MVDR and binaural MVDR-N with different **trade-off parameters**:
 - MVDR-IC
 - MVDR-MS1: $\eta_{\max}=0.7$, $MSC_{\min}=0$
 - MVDR-MS2: $\eta_{\max}=1.0$, $MSC_{\min}=0.1$
- **Subjective listening experiments:**
 - 11 normal-hearing subjects
 - **SRT** using Oldenburg Sentence Test (OLSA)
 - **Spatial quality (diffuseness)** using MUSHRA

Evaluation: Test setup

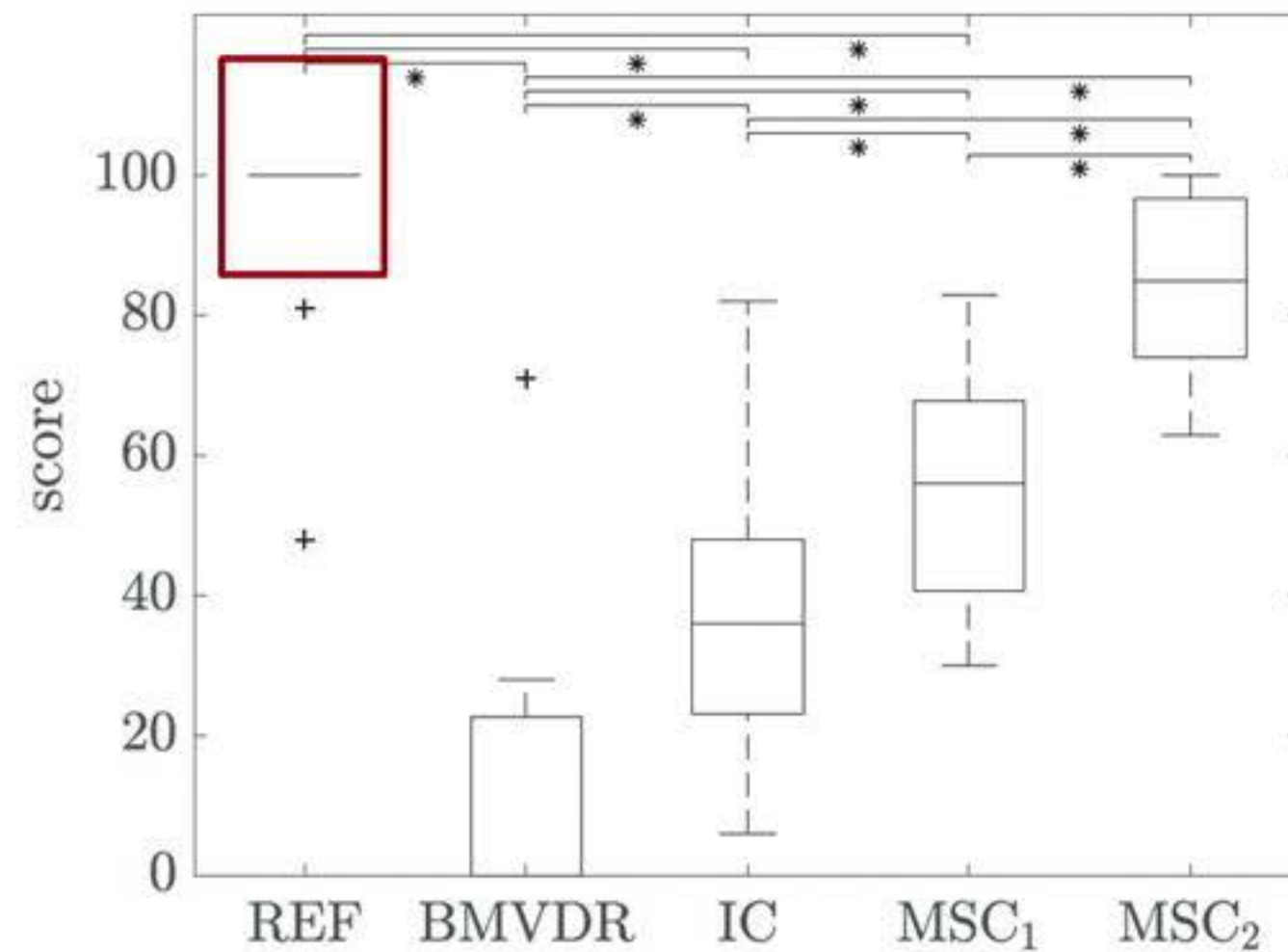


- Binaural hearing aid recordings ($M=4$ mics) in **cafeteria** ($T_{60} \approx 1250$ ms)
 - Target speaker at -35°
 - Realistic cafeteria ambient noise
- **Algorithms:** binaural MVDR and binaural MVDR-N with different **trade-off parameters**:
 - MVDR-IC
 - MVDR-MS1: $\eta_{\max}=0.7$, $MSC_{\min}=0$
 - MVDR-MS2: $\eta_{\max}=1.0$, $MSC_{\min}=0.1$
- **Subjective listening experiments:**
 - 11 normal-hearing subjects
 - **SRT** using Oldenburg Sentence Test (OLSA)
 - **Spatial quality (diffuseness)** using MUSHRA

Does binaural unmasking compensate for SNR decrease of cue preservation algorithms (MVDR-N) ?

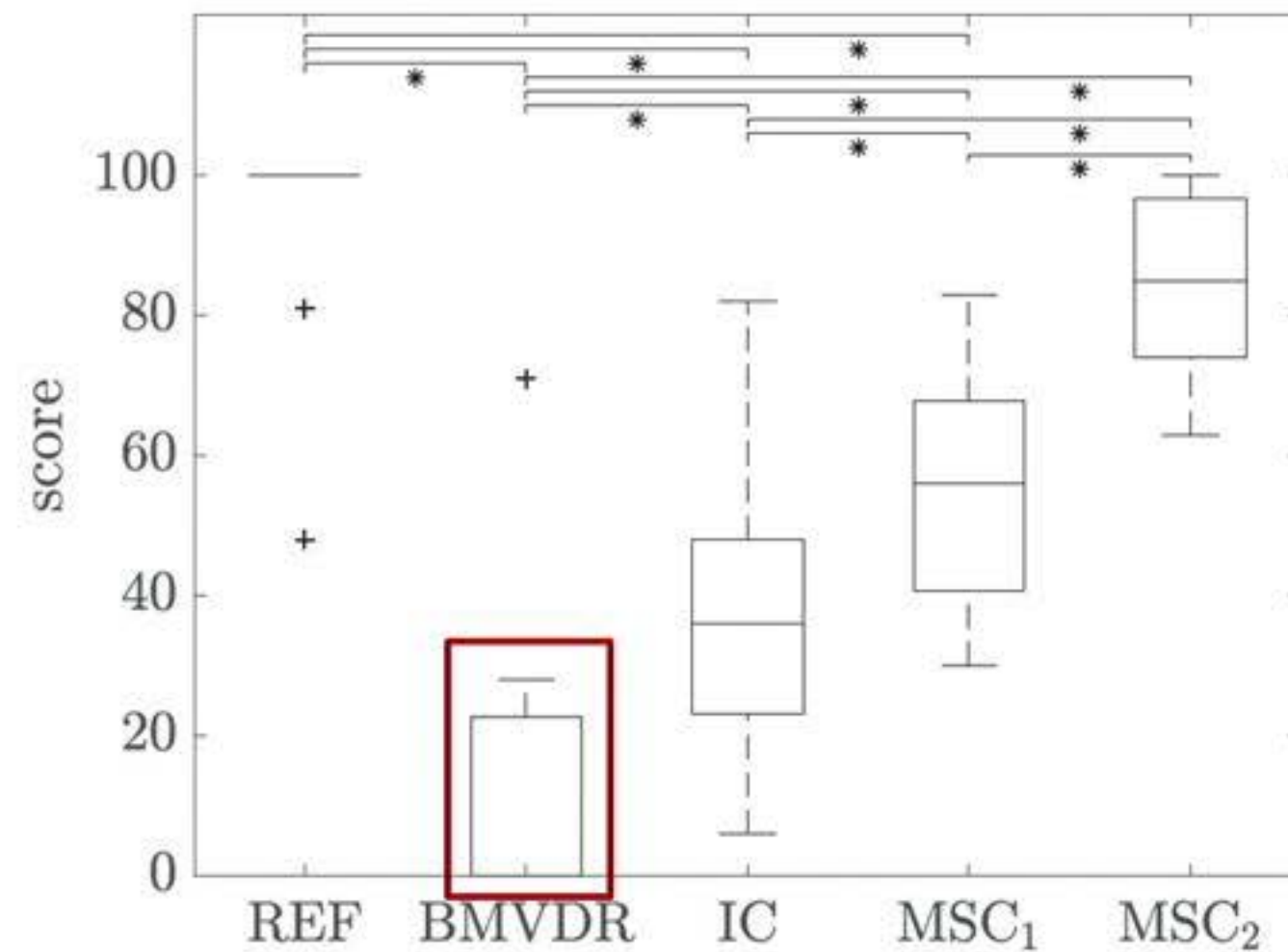
Evaluation: Spatial quality (MUSHRA)

- Evaluate spatial difference between reference microphone signals and binaural output signals



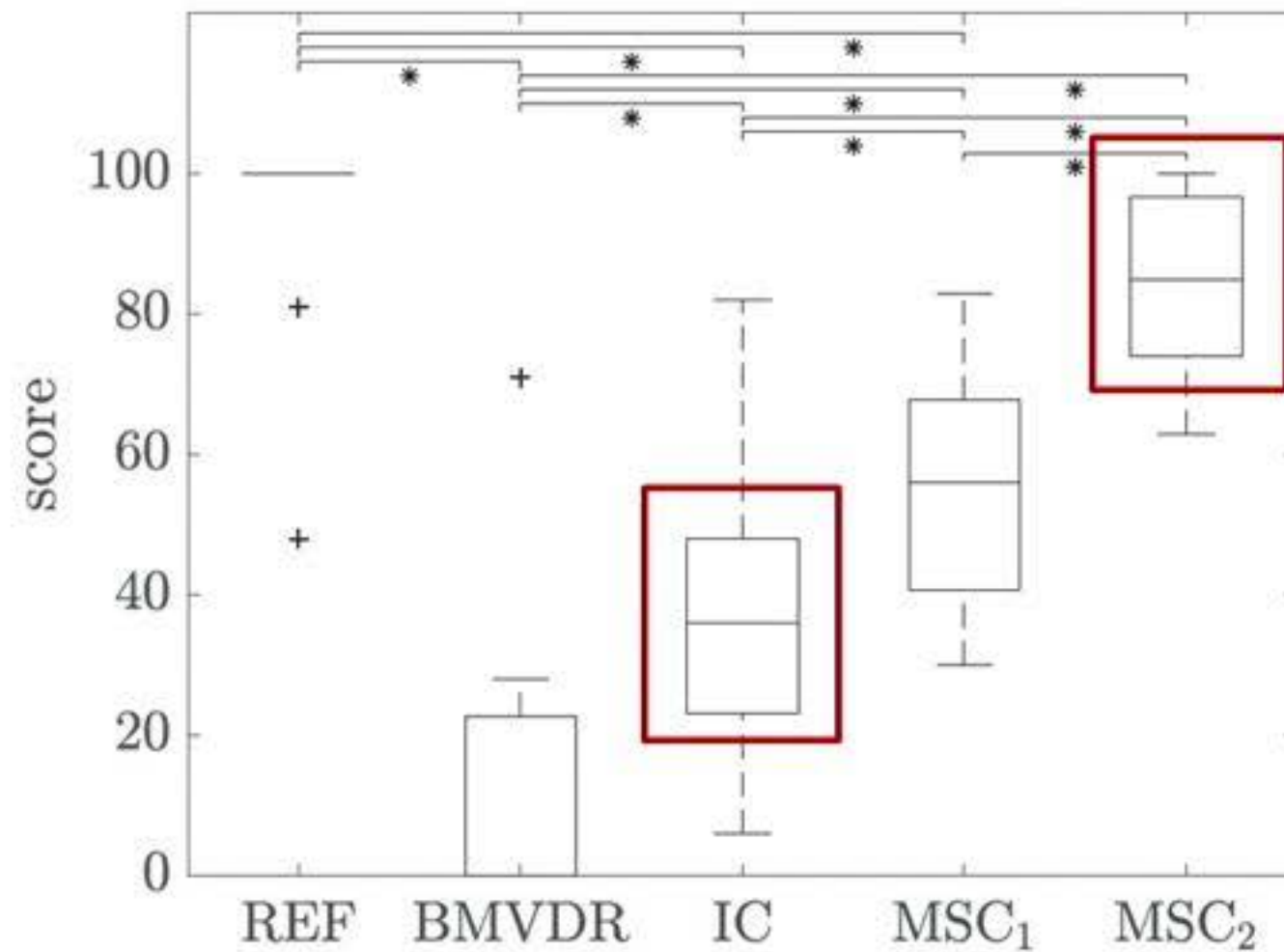
Evaluation: Spatial quality (MUSHRA)

- Evaluate spatial difference between reference microphone signals and binaural output signals



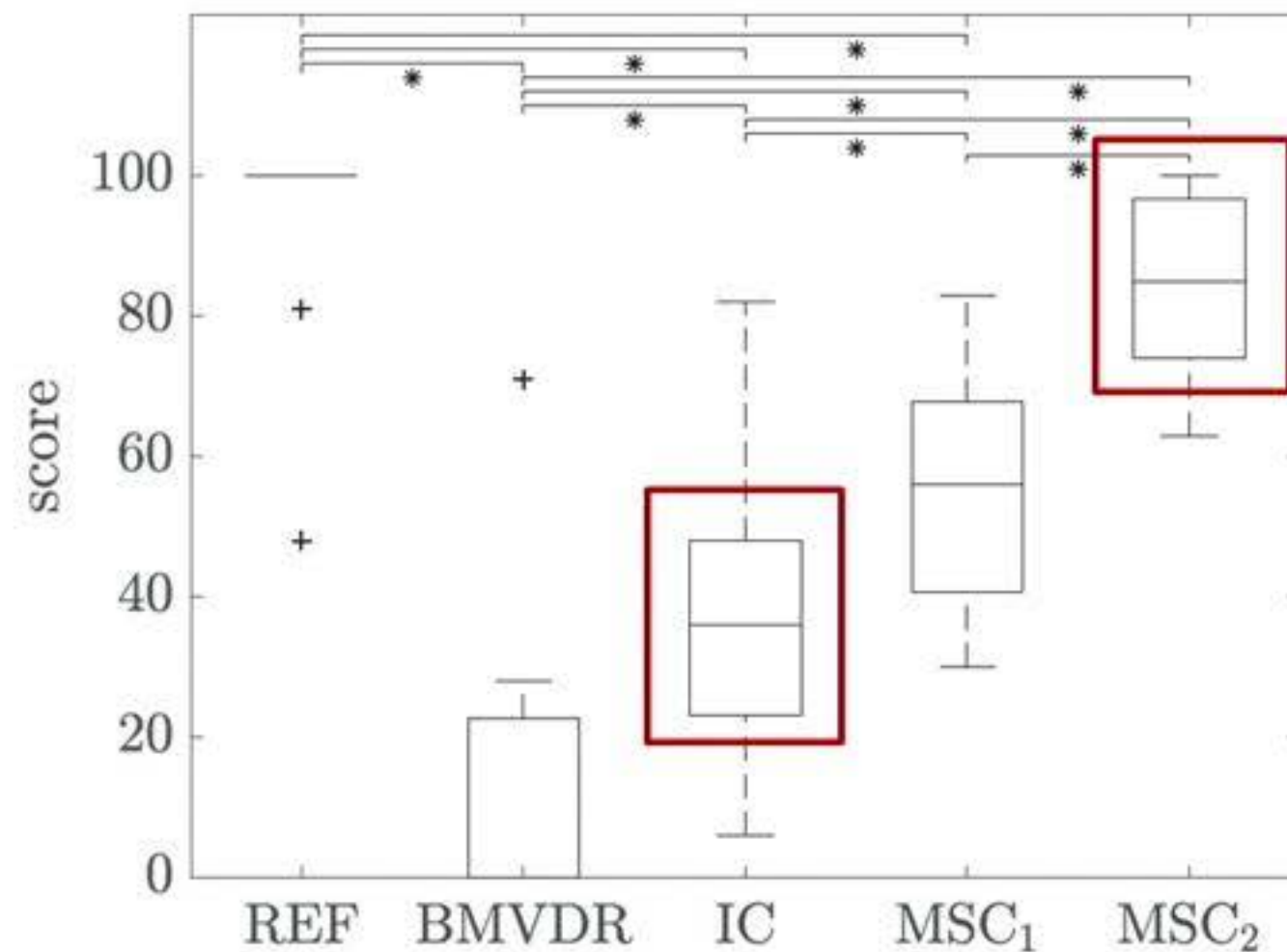
Evaluation: Spatial quality (MUSHRA)

- Evaluate spatial difference between reference microphone signals and binaural output signals
- **MVDR-N outperforms BMVDR**
 - Trade-off parameters: MSC-based better than IC-based
 - Using MSC2 hardly any difference to input !



Evaluation: Spatial quality (MUSHRA)

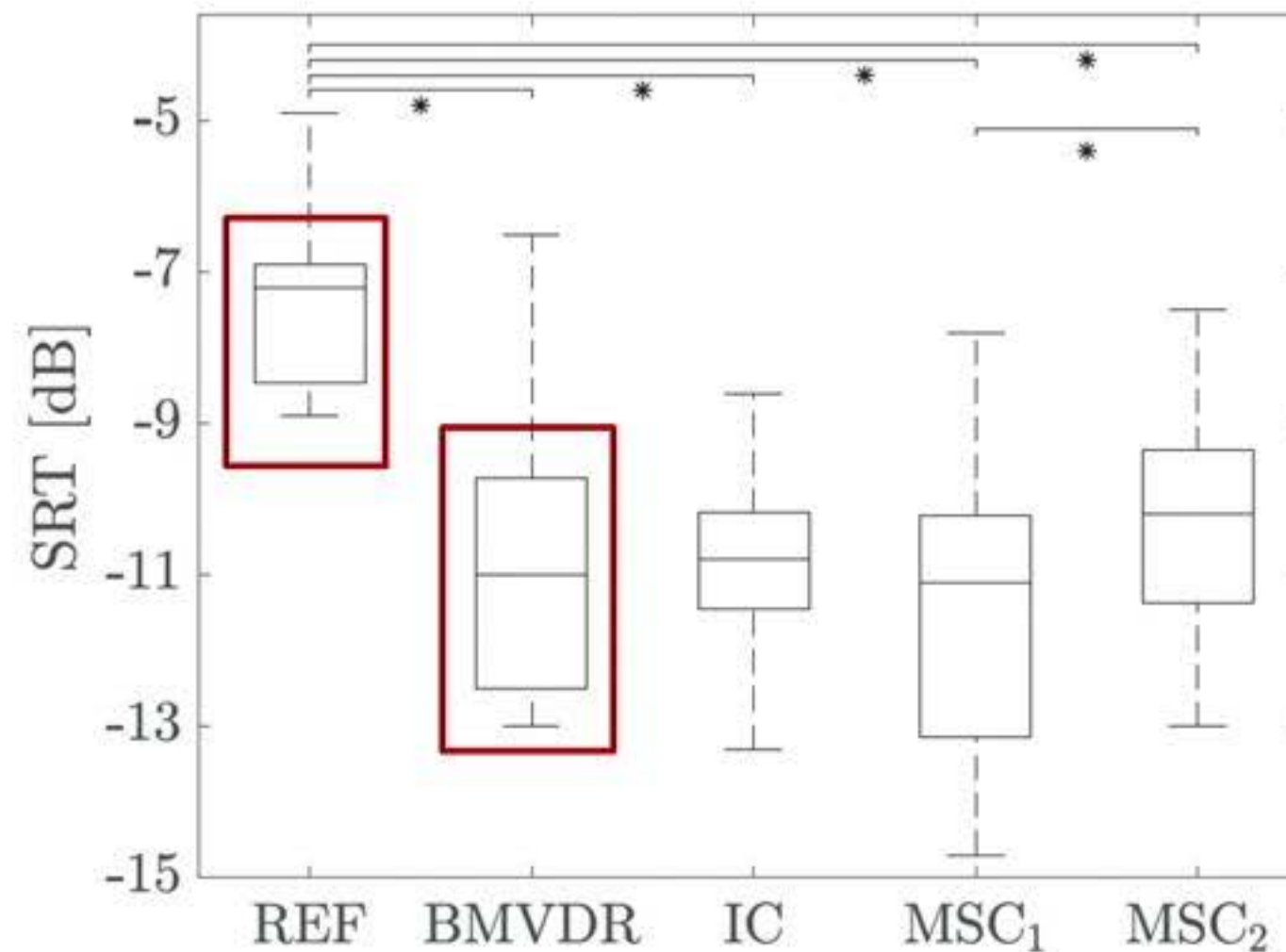
- Evaluate spatial difference between reference microphone signals and binaural output signals
- **MVDR-N outperforms BMVDR**
 - Trade-off parameters: MSC-based better than IC-based
 - Using MSC2 hardly any difference to input !



**Binaural cue preservation for diffuse noise
significantly improves spatial quality**

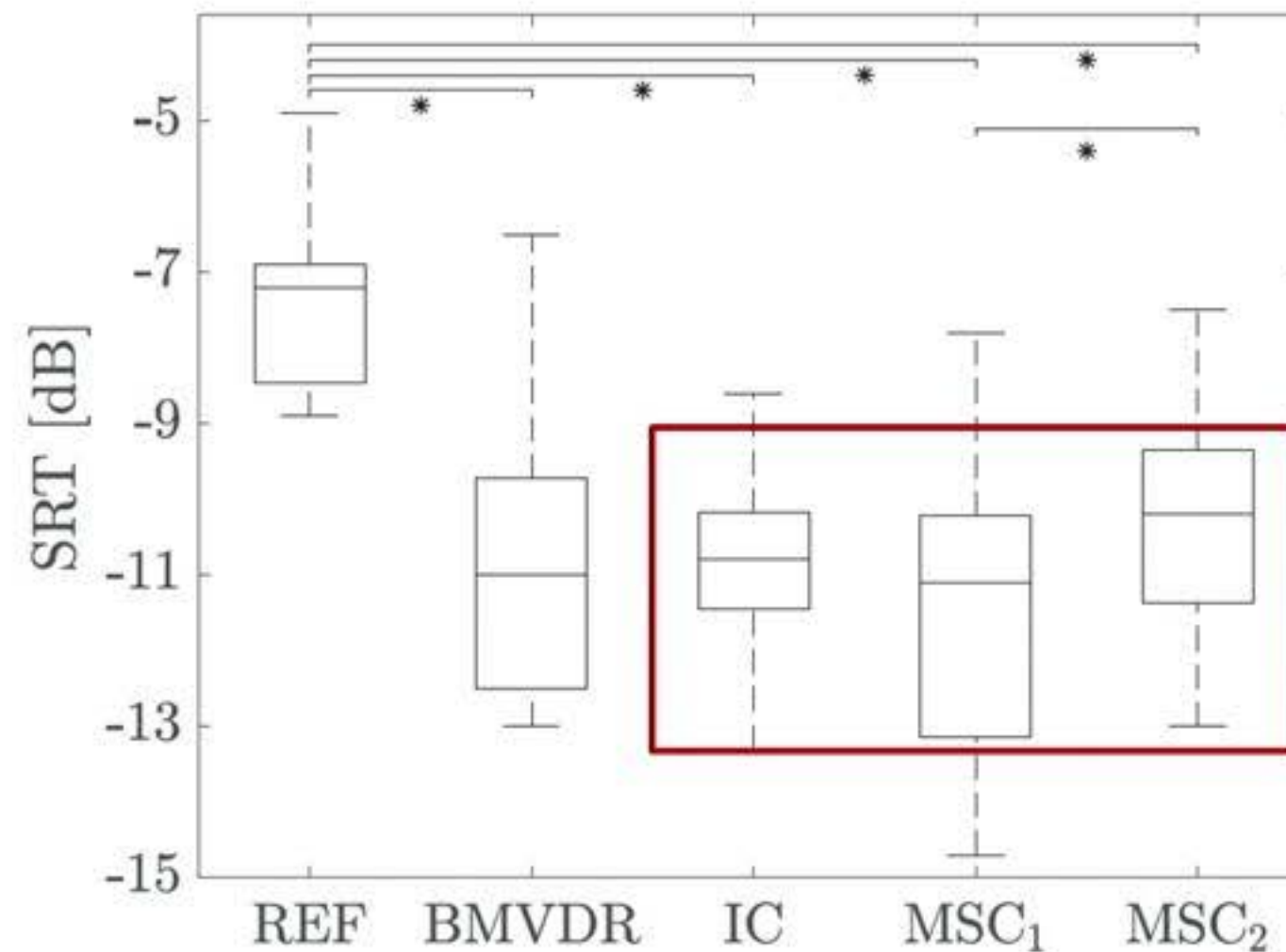
Evaluation: Speech intelligibility (SRT)

- All algorithms show a highly significant speech reception threshold (SRT) improvement



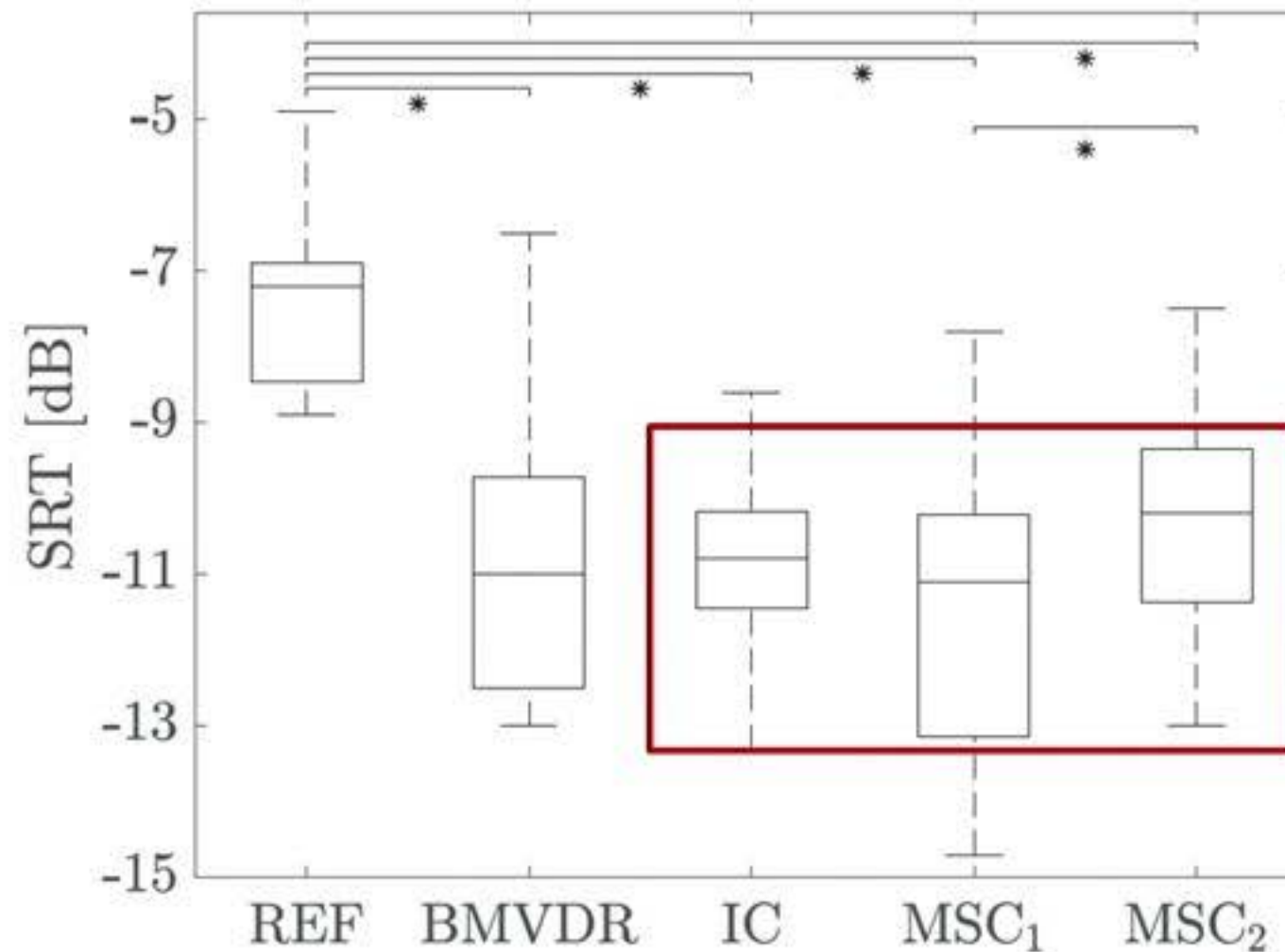
Evaluation: Speech intelligibility (SRT)

- All algorithms show a highly significant speech reception threshold (SRT) improvement
- **No significant SRT difference between BMVDR and MVDR-N**




Evaluation: Speech intelligibility (SRT)

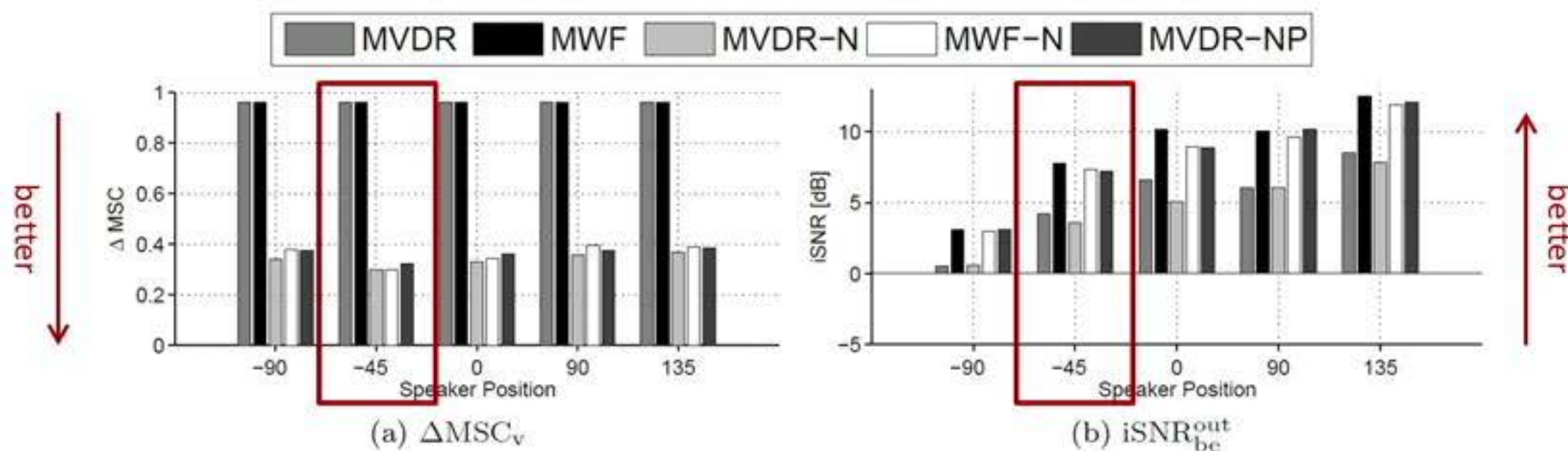
- All algorithms show a highly significant speech reception threshold (SRT) improvement
- **No significant SRT difference between BMVDR and MVDR-N**



**Binaural cue preservation for diffuse noise
does not affect speech intelligibility**

Binaural MVDR/MWF: Sound samples

| Input | MVDR | MWF | MVDR-N | MWF-N | MVDR-NP |
|---|---|---|---|---|---|
|  |  |  |  |  |  |



Cafeteria with recorded ambient noise, speaker at -45° , 0 dB input iSNR (left hearing aid)

MVDR: anechoic ATF, DOA known, spatial coherence matrix calculated from anechoic ATFs / MWF = MVDR + postfilter (SPP-based)

3. Acoustic sensor networks

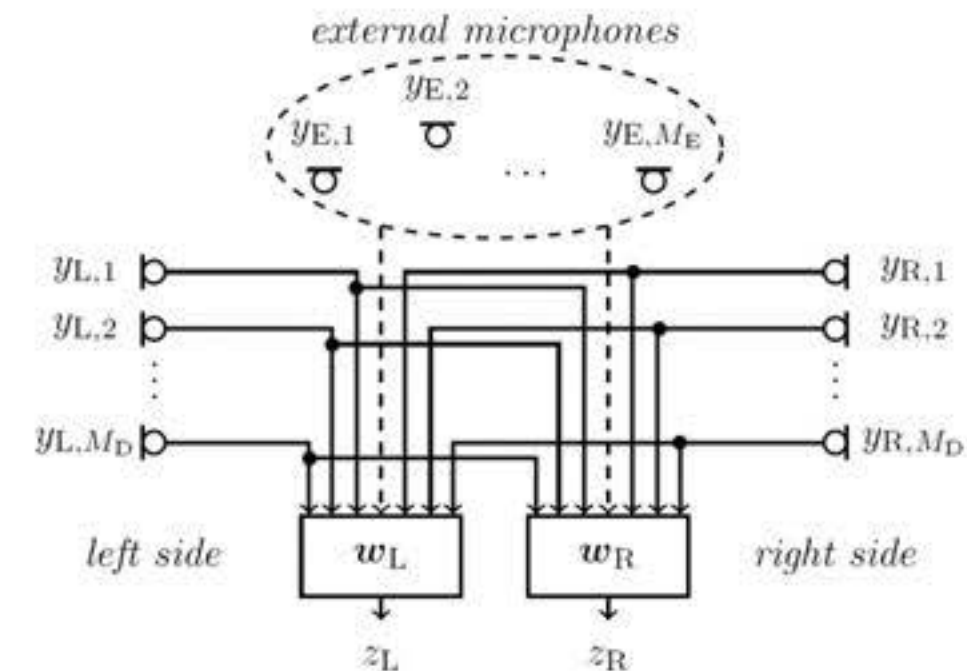
External microphones

- Exploit the availability of one or more external microphones (**acoustic sensor network**) with hearing aids

[Bertrand 2009, Szurley 2016, Yee 2018, Farmani 2018, Kates 2018, Ali 2019, Gößling 2019]

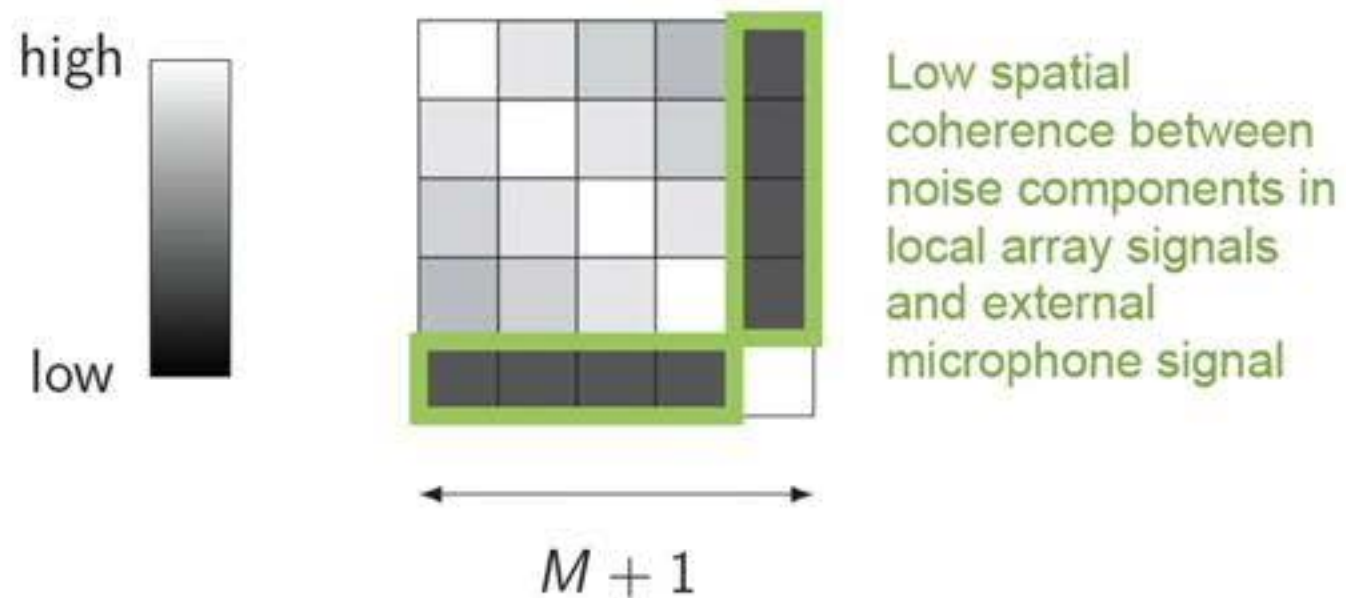
- Integrating external microphone(s) with hearing aid microphones may lead to:
 - Low-complexity method to **estimate relative transfer function (RTF)** vector of target speaker
 - Improved **noise reduction** and **binaural cue preservation** performance

$$\mathbf{W}_L = \frac{\mathbf{R}_v^{-1} \mathbf{a}_L}{\mathbf{a}_L^H \mathbf{R}_v^{-1} \mathbf{a}_L}, \quad \mathbf{W}_R = \frac{\mathbf{R}_v^{-1} \mathbf{a}_R}{\mathbf{a}_R^H \mathbf{R}_v^{-1} \mathbf{a}_R}$$



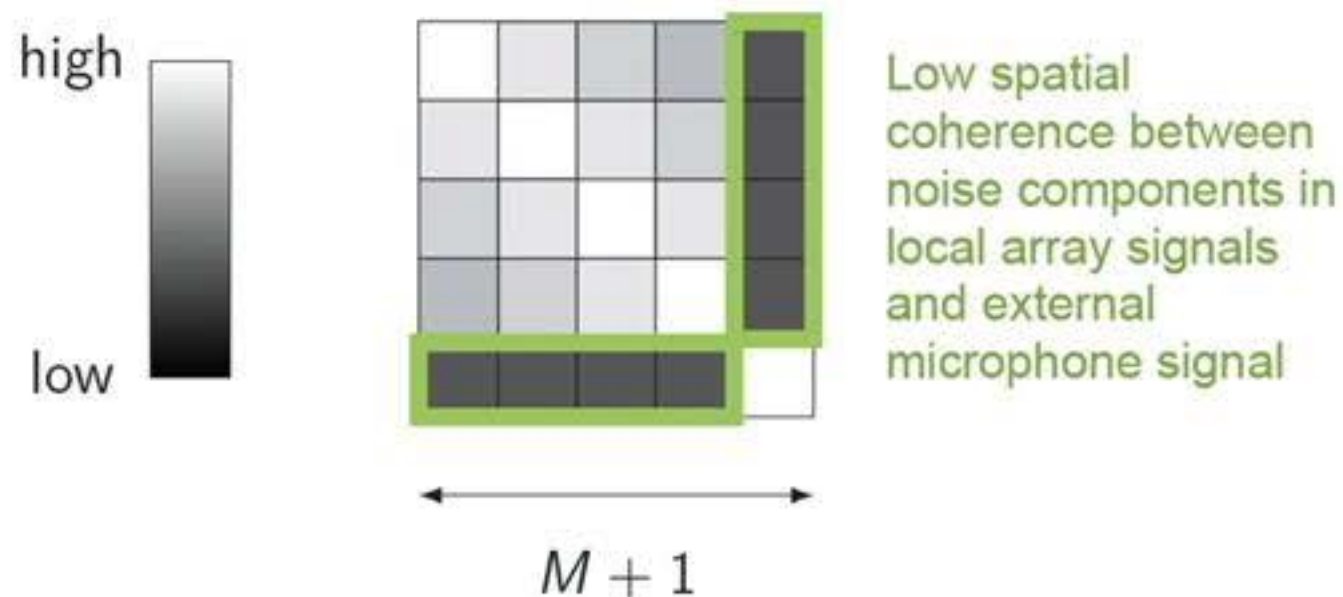
One external microphone: RTF estimation

- **Estimate RTF vector of target speaker** to steer binaural MVDR beamformer
- **Spatial coherence (SC) method:** assume that noise components in external microphone and HA microphones are uncorrelated, e.g., when external microphone is spatially separated from HA microphones + diffuse noise field



One external microphone: RTF estimation

- **Estimate RTF vector of target speaker** to steer binaural MVDR beamformer
- **Spatial coherence (SC) method:** assume that noise components in external microphone and HA microphones are uncorrelated, e.g., when external microphone is spatially separated from HA microphones + diffuse noise field
→ correlate HA microphone signals with external microphone signals and normalize by reference element

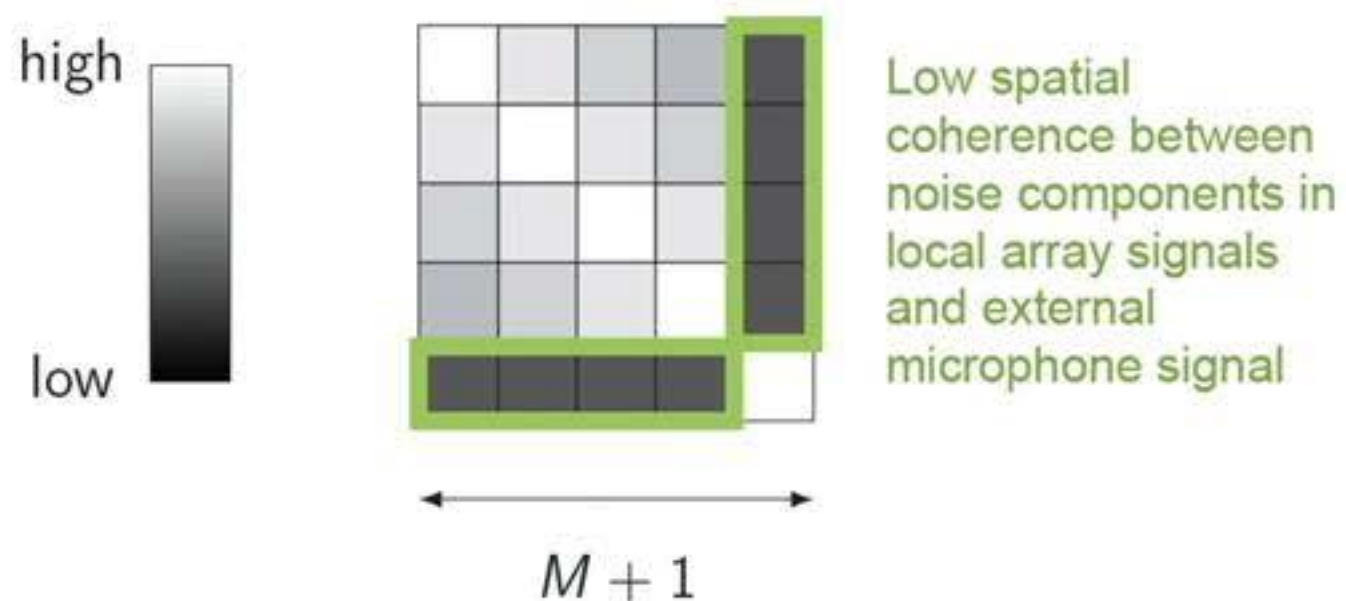


$$\bar{\mathbf{a}}_L^{\text{SCE}} = \frac{\bar{\mathbf{R}}_y \mathbf{e}_E}{\mathbf{e}_L^T \bar{\mathbf{R}}_y \mathbf{e}_E}, \quad \bar{\mathbf{a}}_R^{\text{SCE}} = \frac{\bar{\mathbf{R}}_y \mathbf{e}_E}{\mathbf{e}_R^T \bar{\mathbf{R}}_y \mathbf{e}_E}$$

$$\bar{\mathbf{w}}_L^{\text{SCE}} = \begin{bmatrix} \alpha \cdot [\mathbf{I}_{2M}, \mathbf{0}_{2M \times 1}] \bar{\mathbf{w}}_L \\ \alpha(1 + \beta) \cdot \mathbf{e}_E^T \bar{\mathbf{w}}_L \end{bmatrix}$$

One external microphone: RTF estimation

- **Estimate RTF vector of target speaker** to steer binaural MVDR beamformer
- **Spatial coherence (SC) method:** assume that noise components in external microphone and HA microphones are uncorrelated, e.g., when external microphone is spatially separated from HA microphones + diffuse noise field
→ correlate HA microphone signals with external microphone signals and normalize by reference element



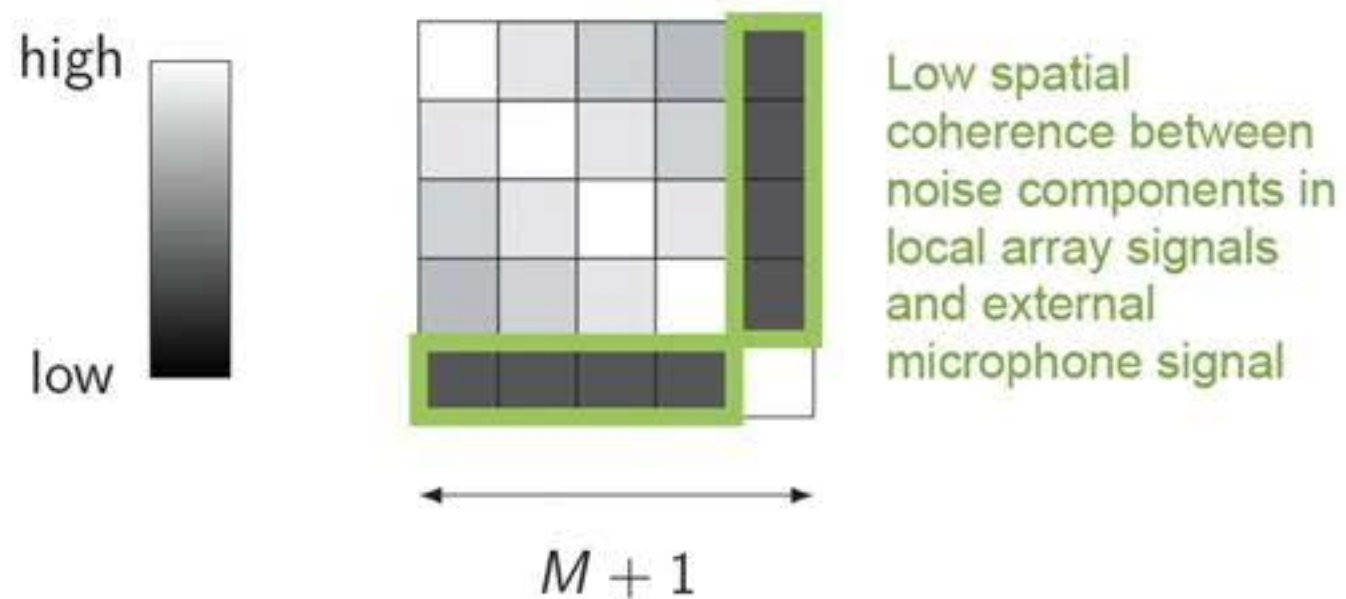
$$\bar{\mathbf{a}}_L^{\text{SCE}} = \frac{\bar{\mathbf{R}}_y \mathbf{e}_E}{\mathbf{e}_L^T \bar{\mathbf{R}}_y \mathbf{e}_E}, \quad \bar{\mathbf{a}}_R^{\text{SCE}} = \frac{\bar{\mathbf{R}}_y \mathbf{e}_E}{\mathbf{e}_R^T \bar{\mathbf{R}}_y \mathbf{e}_E}$$

$$\bar{\mathbf{w}}_L^{\text{SCE}} = \begin{bmatrix} \alpha \cdot [\mathbf{I}_{2M}, \mathbf{0}_{2M \times 1}] \bar{\mathbf{w}}_L \\ \alpha(1 + \beta) \cdot \mathbf{e}_E^T \bar{\mathbf{w}}_L \end{bmatrix}$$

real-valued bias

One external microphone: RTF estimation

- **Estimate RTF vector of target speaker** to steer binaural MVDR beamformer
- **Spatial coherence (SC) method:** assume that noise components in external microphone and HA microphones are uncorrelated, e.g., when external microphone is spatially separated from HA microphones + diffuse noise field
→ correlate HA microphone signals with external microphone signals and normalize by reference element
- **Low computational complexity** with similar (even better in practice) performance than state-of-the-art covariance whitening (CW) approach



$$\bar{\mathbf{a}}_L^{\text{SCE}} = \frac{\bar{\mathbf{R}}_y \mathbf{e}_E}{\mathbf{e}_L^T \bar{\mathbf{R}}_y \mathbf{e}_E}, \quad \bar{\mathbf{a}}_R^{\text{SCE}} = \frac{\bar{\mathbf{R}}_y \mathbf{e}_E}{\mathbf{e}_R^T \bar{\mathbf{R}}_y \mathbf{e}_E}$$

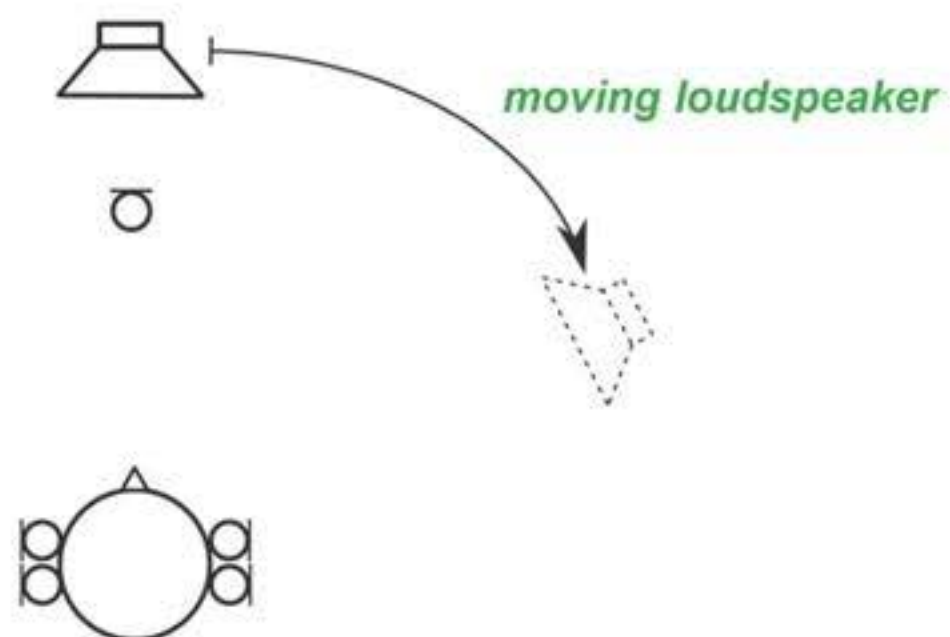
$$\bar{\mathbf{w}}_L^{\text{SCE}} = \begin{bmatrix} \alpha \cdot [\mathbf{I}_{2M}, \mathbf{0}_{2M \times 1}] \bar{\mathbf{w}}_L \\ \alpha(1 + \beta) \cdot \mathbf{e}_E^T \bar{\mathbf{w}}_L \end{bmatrix}$$

One external microphone: Simulation results



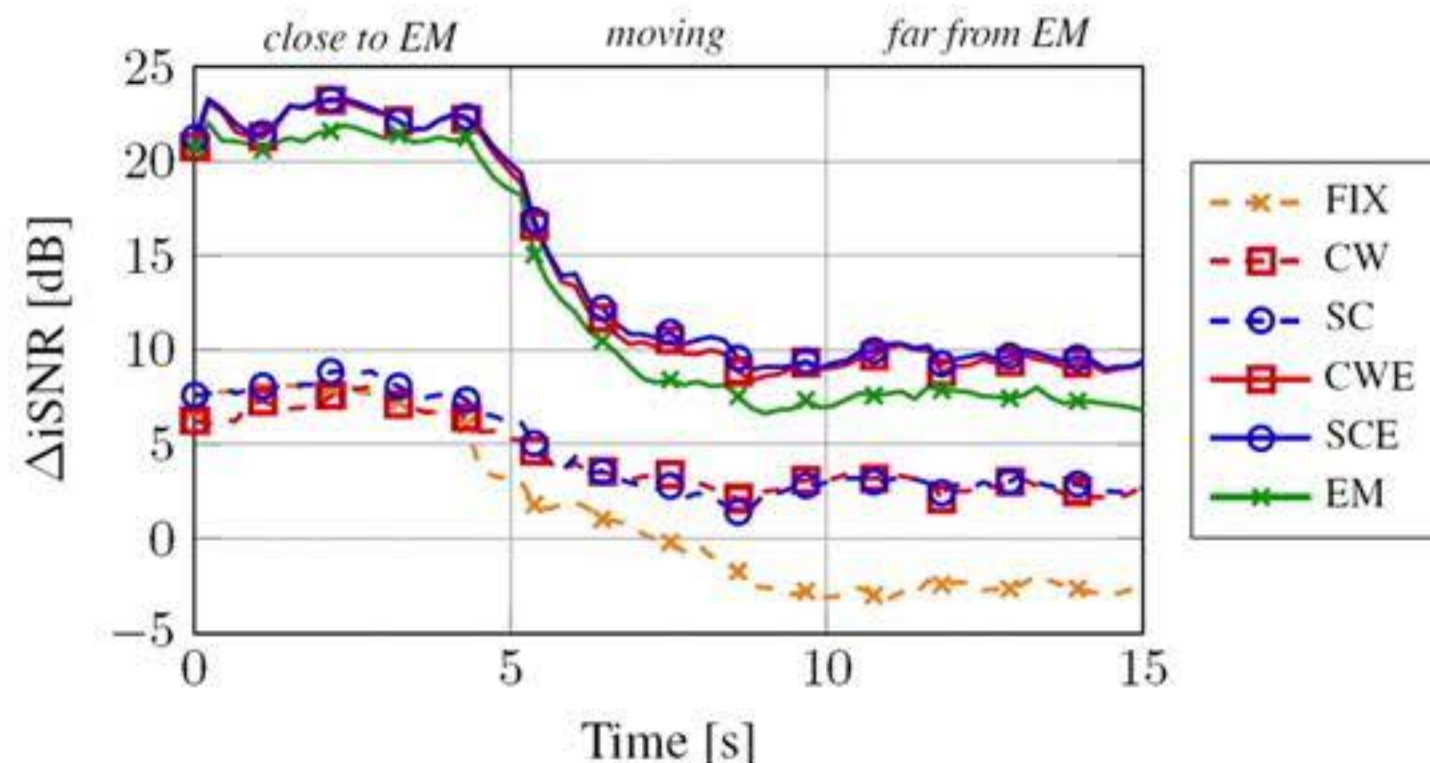
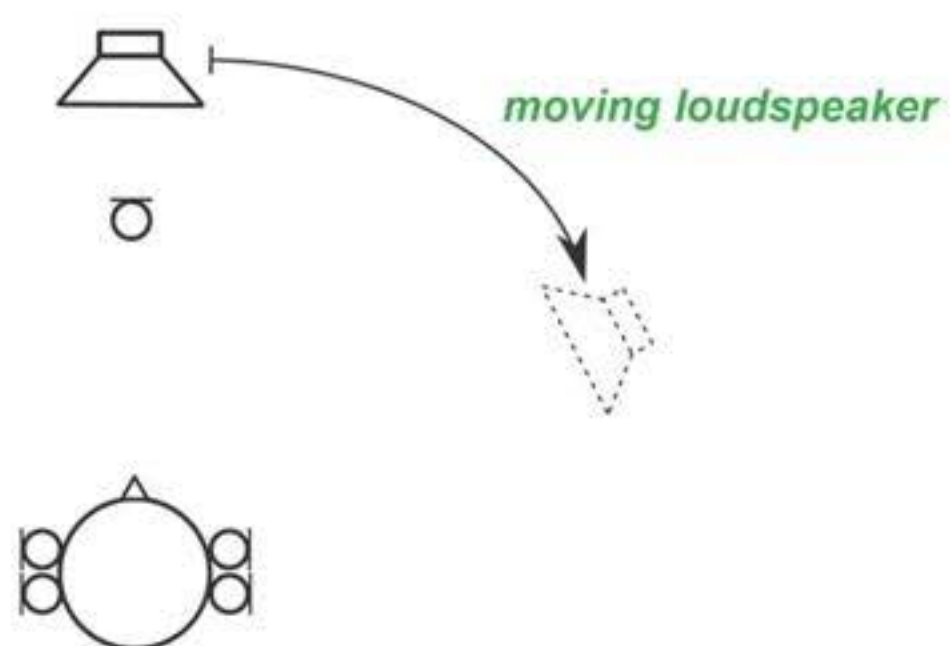
Oldenburg VARECHOIC Lab ($T_{60} \approx 350\text{ms}$), $M=4 + 1$ external mic (1.5m/0.5m), moving speaker, pseudo-diffuse babble noise, $i\text{SNR}=0\text{dB}$ (right HA)
STFT: 32 ms, 50% overlap, sqrt-Hann; SPP in external microphone; smoothing: 100 ms (speech), 1 s (noise)

One external microphone: Simulation results



Oldenburg Varechoic Lab ($T_{60} \approx 350\text{ms}$), $M=4 + 1$ external mic (1.5m/0.5m), moving speaker, pseudo-diffuse babble noise, $i\text{SNR}=0\text{dB}$ (right HA)
STFT: 32 ms, 50% overlap, sqrt-Hann; SPP in external microphone; smoothing: 100 ms (speech), 1 s (noise)

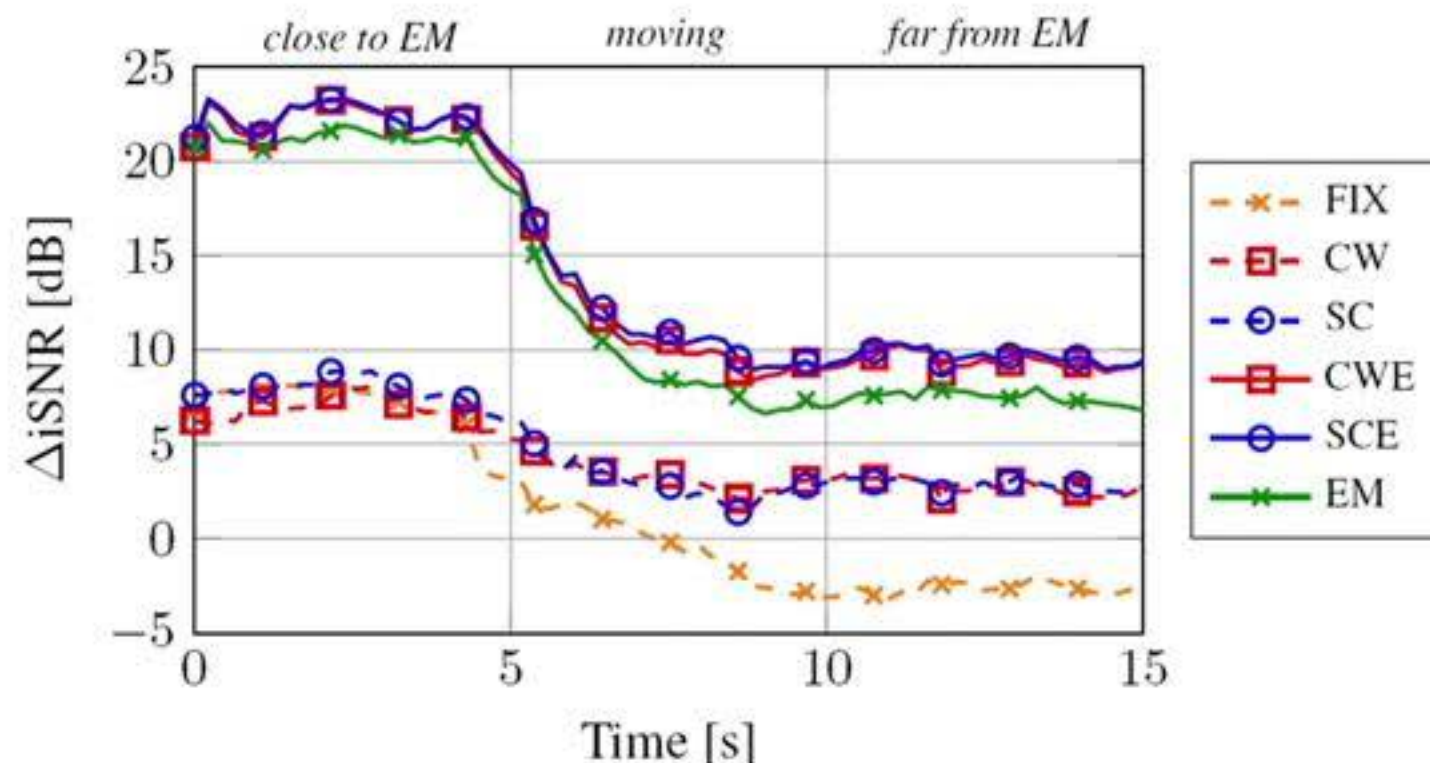
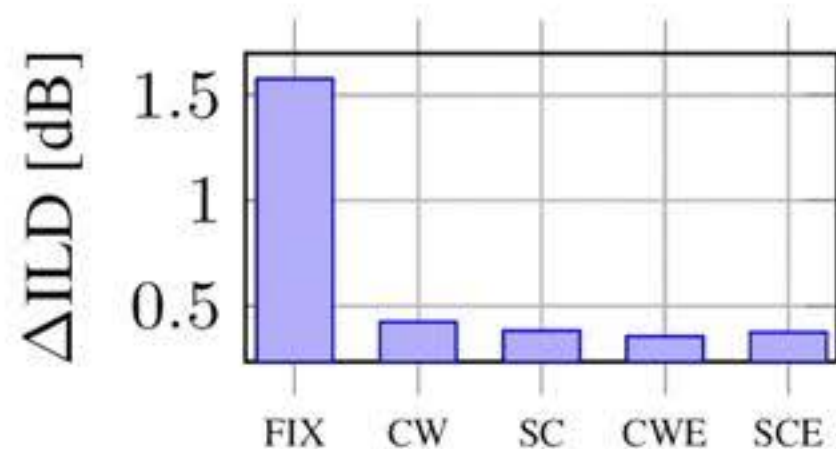
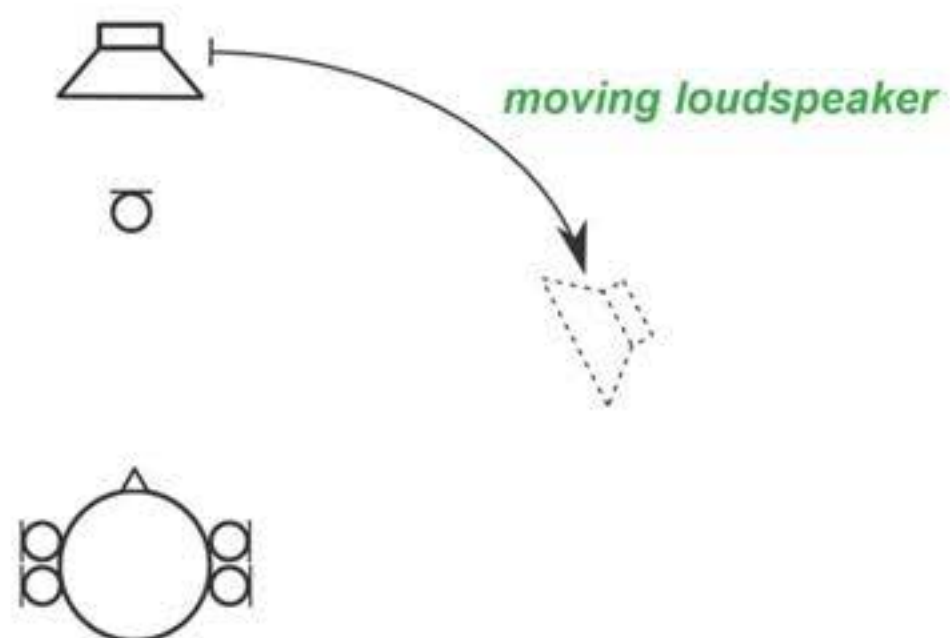
One external microphone: Simulation results



- MVDR with external microphone (**SCE**) leads to **better SNR** compared to MVDR using only HA microphones (**SC**, **FIX**) and external microphone (**EM**)

Oldenburg Varechoic Lab ($T_{60} \approx 350\text{ms}$), $M=4 + 1$ external mic (1.5m/0.5m), moving speaker, pseudo-diffuse babble noise, $iSNR=0\text{dB}$ (right HA)
STFT: 32 ms, 50% overlap, sqrt-Hann; SPP in external microphone; smoothing: 100 ms (speech), 1 s (noise)

One external microphone: Simulation results



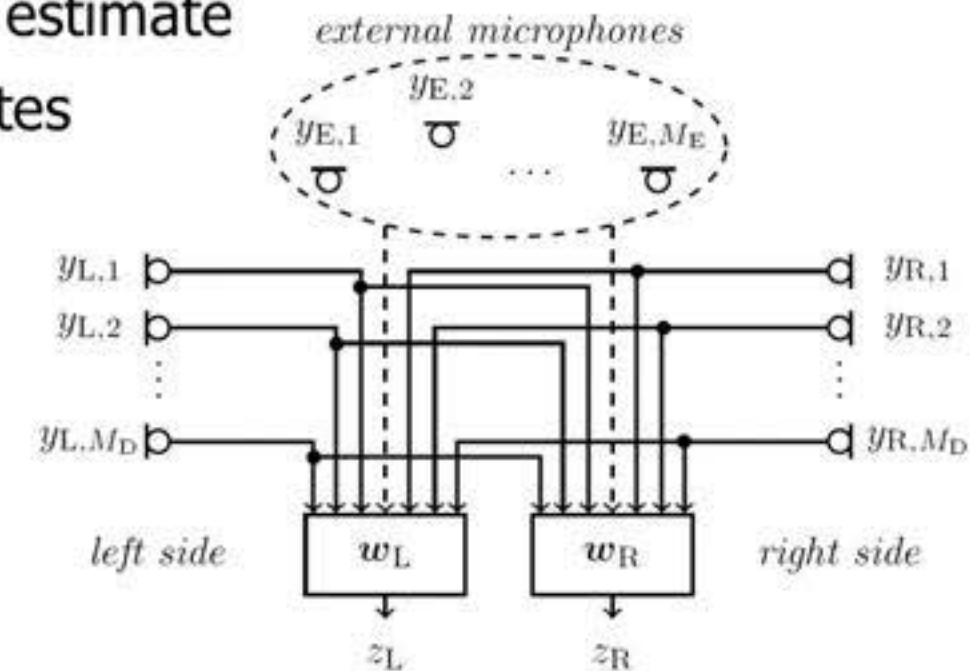
- MVDR with external microphone (**SCE**) leads to **better SNR** compared to MVDR using only HA microphones (**SC**, **FIX**) and external microphone (**EM**)
- MVDR using estimated RTFs (**SCE**, **SC**) **preserves binaural cues of target speaker** compared to fixed MVDR (**FIX**) and external microphone (**EM**)

Oldenburg Varechoic Lab ($T_{60} \approx 350\text{ms}$), $M=4 + 1$ external mic (1.5m/0.5m), moving speaker, pseudo-diffuse babble noise, $i\text{SNR}=0\text{dB}$ (right HA)
STFT: 32 ms, 50% overlap, sqrt-Hann; SPP in external microphone; smoothing: 100 ms (speech), 1 s (noise)

Multiple external microphones

- Each external microphone yields (different) RTF estimate
- Linear combination/selection** of RTF estimates (per frequency)

$$\mathbf{a}_L^{\text{SC-C}} = \frac{\mathbf{A}_L^{\text{SC}} \mathbf{c}}{\mathbf{e}_L^T \mathbf{A}_L^{\text{SC}} \mathbf{c}}$$



Multiple external microphones

- Each external microphone yields (different) RTF estimate
- Linear combination/selection** of RTF estimates (per frequency)

$$\mathbf{a}_L^{\text{SC-C}} = \frac{\mathbf{A}_L^{\text{SC}} \mathbf{c}}{\mathbf{e}_L^T \mathbf{A}_L^{\text{SC}} \mathbf{c}}$$

1. Input SNR-based selection

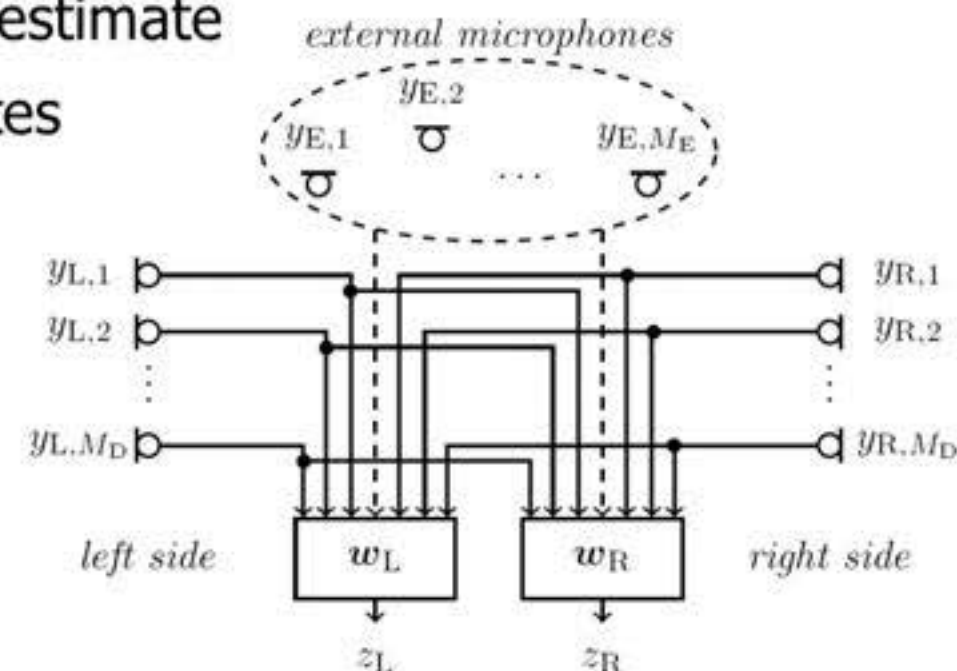
$$\mathbf{c}^{\text{iSNR}} = \mathbf{e}_{E,\hat{i}}, \quad \hat{i} = \arg \max_i \frac{\mathbf{e}_{E,i}^T \mathbf{R}_y \mathbf{e}_{E,i}}{\mathbf{e}_{E,i}^T \mathbf{R}_n \mathbf{e}_{E,i}}$$

2. Simple averaging

$$\mathbf{c}^{\text{AV}} = \left[\frac{1}{M_E}, \dots, \frac{1}{M_E} \right]^T$$

3. Output SNR-maximizing combination

$$\mathbf{c}^{\text{mSNR}} = \arg \max_{\mathbf{c}} \text{SNR}_{\text{BMVDR,L}}^{\text{out}} = \mathcal{P}\{\mathbf{\Lambda}_2^{-1} \mathbf{\Lambda}_1\}$$



Multiple external microphones

- Each external microphone yields (different) RTF estimate
- Linear combination/selection** of RTF estimates (per frequency)

$$\mathbf{a}_L^{\text{SC-C}} = \frac{\mathbf{A}_L^{\text{SC}} \mathbf{c}}{\mathbf{e}_L^T \mathbf{A}_L^{\text{SC}} \mathbf{c}}$$

1. Input SNR-based selection

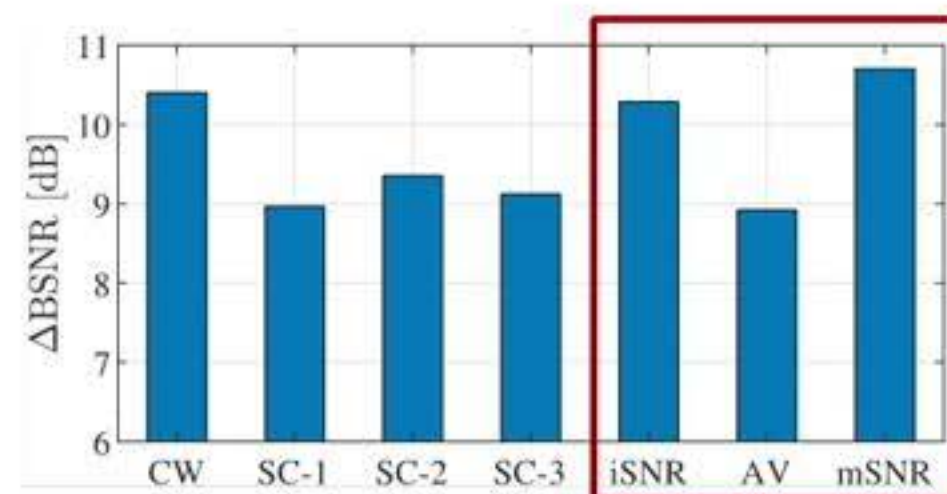
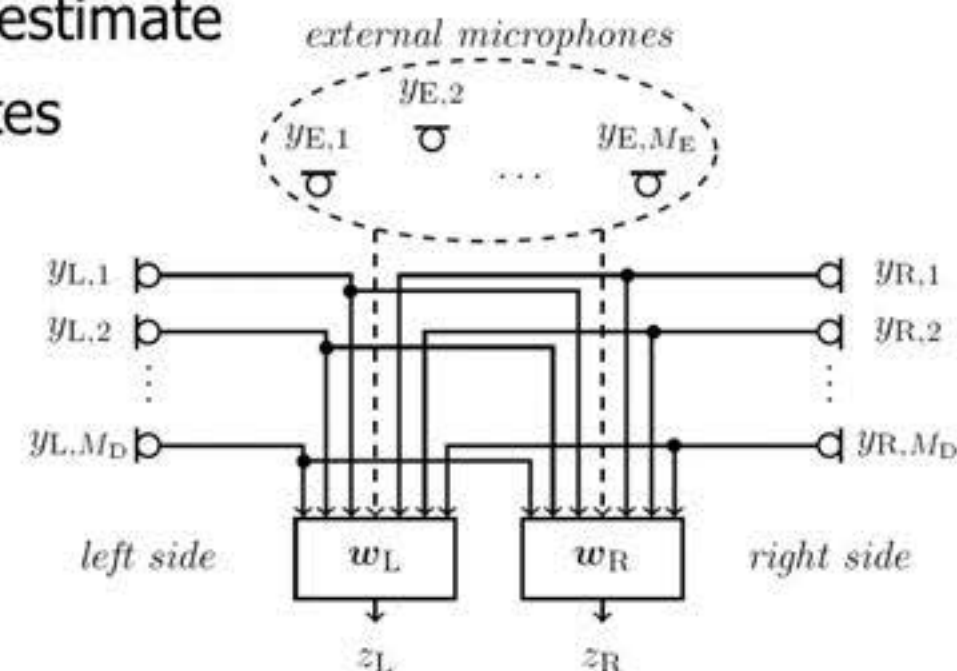
$$\mathbf{c}^{\text{iSNR}} = \mathbf{e}_{E,\hat{i}}, \quad \hat{i} = \arg \max_i \frac{\mathbf{e}_{E,i}^T \mathbf{R}_y \mathbf{e}_{E,i}}{\mathbf{e}_{E,i}^T \mathbf{R}_n \mathbf{e}_{E,i}}$$

2. Simple averaging

$$\mathbf{c}^{\text{AV}} = \left[\frac{1}{M_E}, \dots, \frac{1}{M_E} \right]^T$$

3. Output SNR-maximizing combination

$$\mathbf{c}^{\text{mSNR}} = \arg \max_{\mathbf{c}} \text{SNR}_{\text{BMVDR,L}}^{\text{out}} = \mathcal{P}\{\mathbf{\Lambda}_2^{-1} \mathbf{\Lambda}_1\}$$



Audio Demo

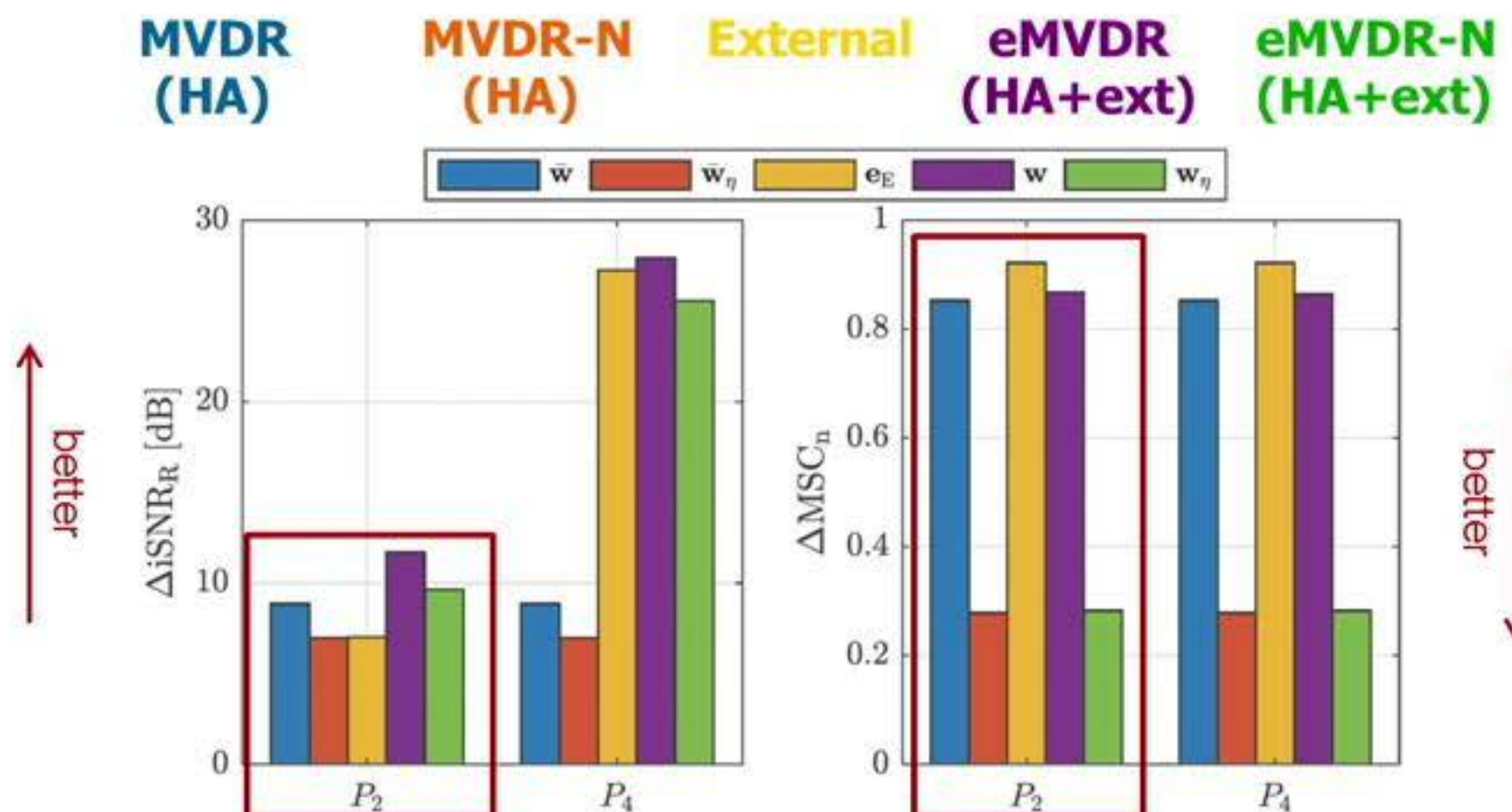
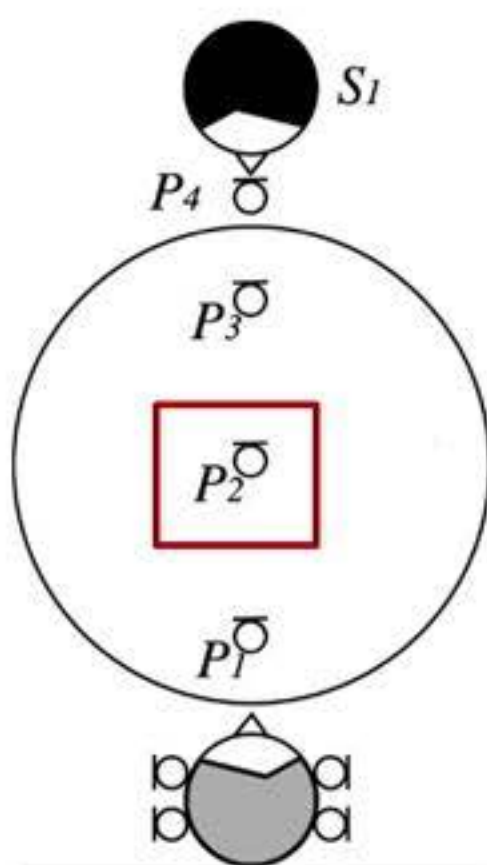
- Real-world recordings ($T_{60} \approx 300$ ms), **moving speaker**
- KEMAR with **two BTE hearing aids** (2 mics each) and **one external mic**
- Pseudo-diffuse babble noise

Audio Demo



Binaural MVDR-N beamformer

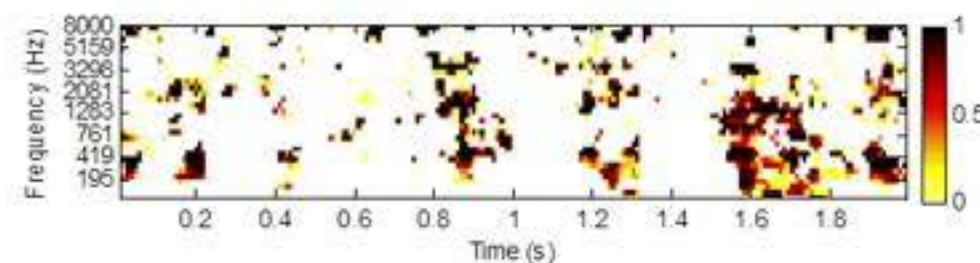
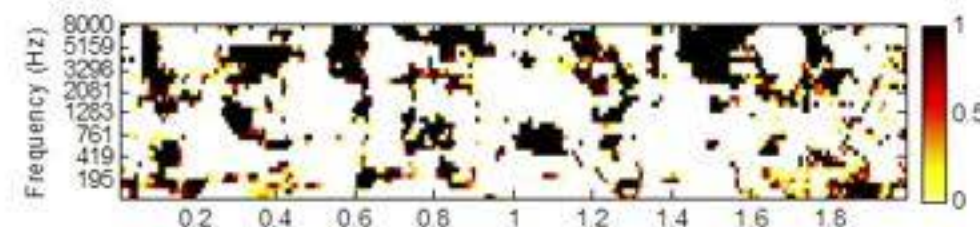
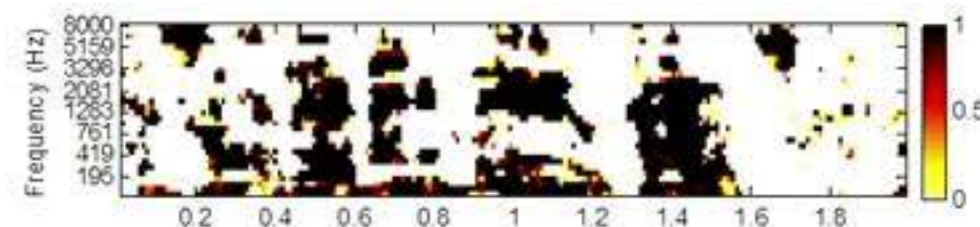
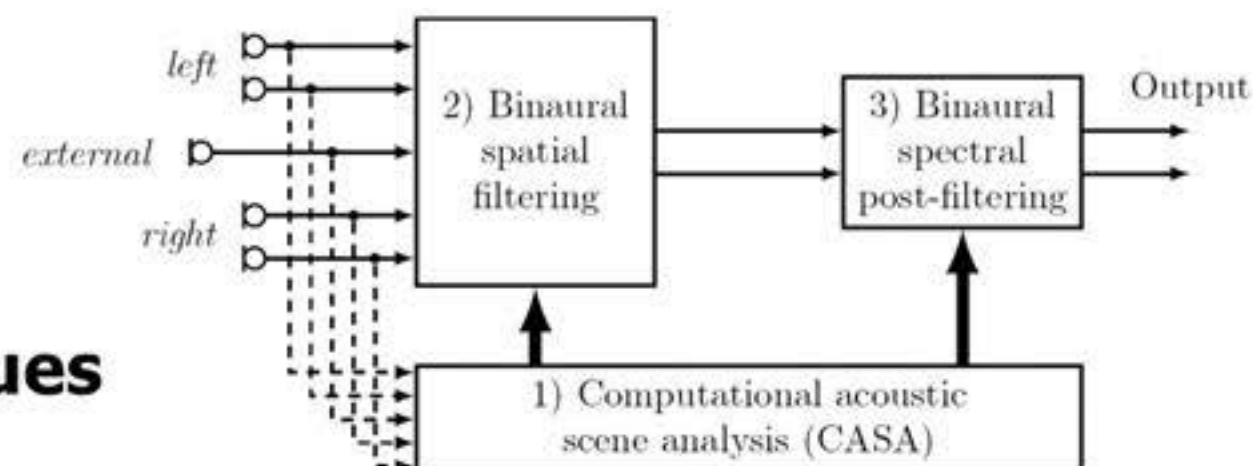
- Including external microphone in **binaural MVDR-N beamformer** leads to:
 - Larger output SNR** for same trade-off parameter η
 - Same output SNR with larger trade-off parameter $\eta \rightarrow$ **better cue preservation**



Starkey database with real-world recordings ($T_{60} \approx 620\text{ms}$), $M=4$, target speaker S_1 , multi-talker babble noise, 0 dB input iSNR (right hearing aid)
MVDR: perfectly estimated noise correlation matrix, RTF of target speaker estimated using covariance whitening method

Current/future work

- **Performance analysis** for different acoustic scenarios (interfering speakers)
- **Synchronization/latency issues**
- **Complex and time-varying scenarios:** incorporate computational acoustic scene analysis (CASA) into control path of developed algorithms
- **Subjective evaluation** of binaural speech enhancement algorithms with **HA/CI users** ongoing



Conclusions

- **Speech communication applications:** on-line speech enhancement algorithms for dynamic acoustic scenarios required

Conclusions

- **Speech communication applications:** on-line speech enhancement algorithms for dynamic acoustic scenarios required
- **Joint noise reduction and dereverberation** using multiple microphones:
 - MVDR beamformer + spectral postfiltering: estimates of time-varying spatial and spectral variables (RETF vector, PSDs)
 - Reverberation suppression: multi-channel linear prediction

Conclusions

- **Speech communication applications:** on-line speech enhancement algorithms for dynamic acoustic scenarios required
- **Joint noise reduction and dereverberation** using multiple microphones:
 - MVDR beamformer + spectral postfiltering: estimates of time-varying spatial and spectral variables (RETF vector, PSDs)
 - Reverberation suppression: multi-channel linear prediction
- **Binaural hearing devices** with binaural output signals:
 - Extensions of binaural MVDR/MWF enable to improve speech intelligibility while preserving spatial awareness (binaural cues)
 - Improved performance when integrating external microphones (acoustic sensor networks)

Acknowledgments



Dr. Ina
Kodrasi



Dr. Ante
Jukić



Dr. Daniel
Marquardt



Marvin
Tammen



Jonas
Klug



Nico
Gößling



Wiebke
Middelberg



Prof. Timo
Gerkmann



Prof. Sharon
Gannot

□ Funding:

- Cluster of Excellence Hearing4all (DFG), Research Unit Individualized Hearing Acoustics (DFG)
- Marie-Curie Initial Training Network "Dereverberation and Reverberation of Audio, Music, and Speech" (EU)
- Joint Lower-Saxony Israel Project "Acoustic scene aware speech enhancement for binaural hearing aids" (Partner: Bar-Ilan University, Israel)
- German-Israeli Foundation Project "Signal Dereverberation Algorithms for Next-Generation Binaural Hearing Aids" (Partners: International Audiolabs Erlangen; Bar-Ilan University, Israel)



Questions ?



House of Hearing, Oldenburg