

# Separake

## Source Separation with a Little Help from Echoes

Robin Scheibler   Diego Di Carlo   Antoine Deleforge   Ivan Dokmanić  
APRIL 17, 2018



ICASSP 2018



## Echoes Help Indoor Processing

- beamforming
- source localization
- self-localization

What about speech separation ?

1. Is speech separation easier with echoes than without ?
2. Full RIR vs a few early reflections ?

## Echoes Help Indoor Processing

- beamforming
- source localization
- self-localization

## What about speech separation ?

Is speech separation easier with echoes than without ?

Q: Full RIR or a few early reflections ?

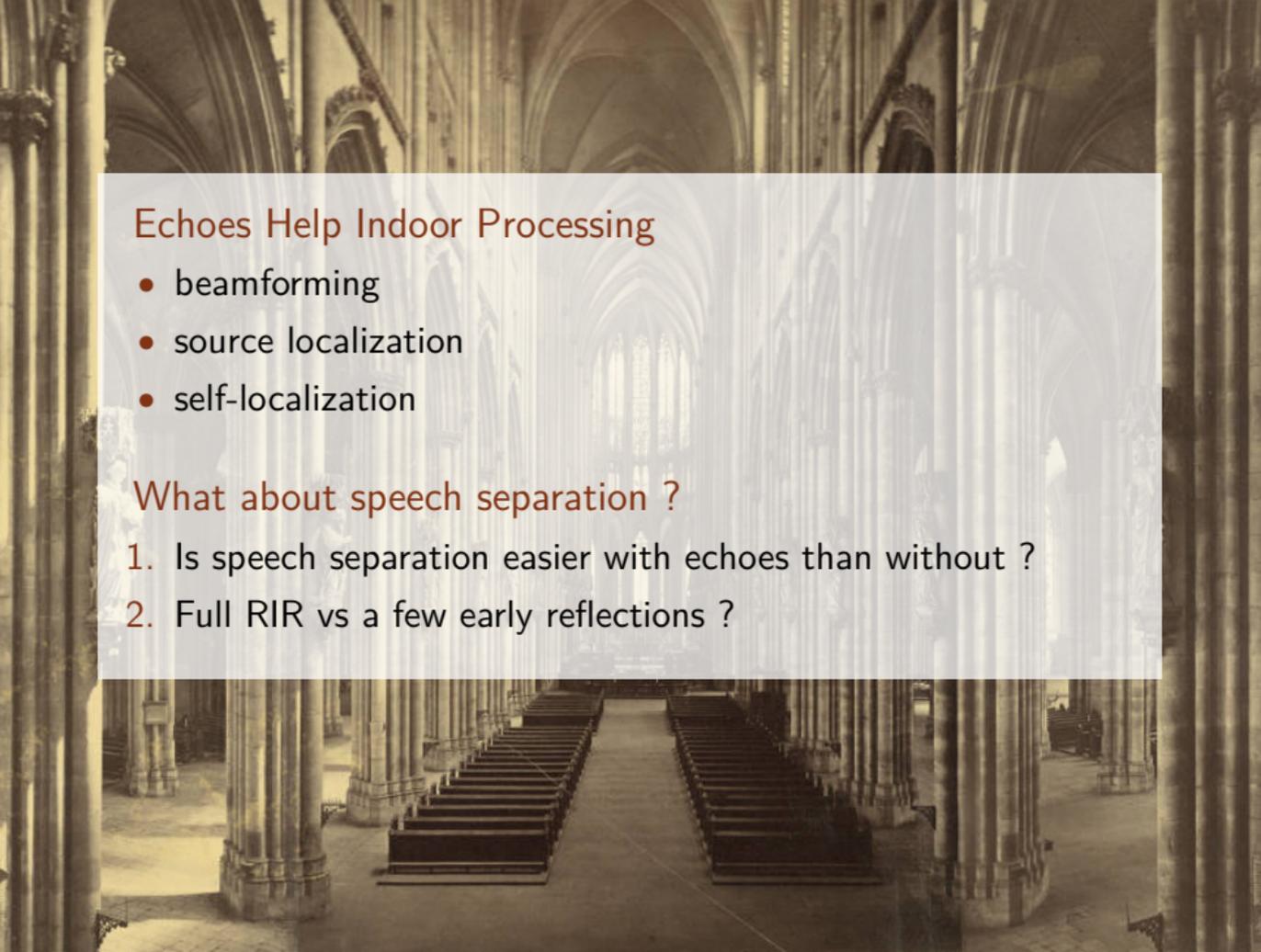
## Echoes Help Indoor Processing

- beamforming
- source localization
- self-localization

## What about speech separation ?

1. Is speech separation easier with echoes than without ?

Full Research Paper: <https://arxiv.org/abs/1508.03591>



## Echoes Help Indoor Processing

- beamforming
- source localization
- self-localization

## What about speech separation ?

1. Is speech separation easier with echoes than without ?
2. Full RIR vs a few early reflections ?

1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

<sup>1</sup>Ozerov & Févotte, 2010

1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

<sup>1</sup>Ozerov & Févotte, 2010

1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

<sup>1</sup>Ozerov & Févotte, 2010

1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

<sup>1</sup>Ozerov & Févotte, 2010

1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

<sup>1</sup>Ozerov & Févotte, 2010

1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

<sup>1</sup>Ozerov & Févotte, 2010

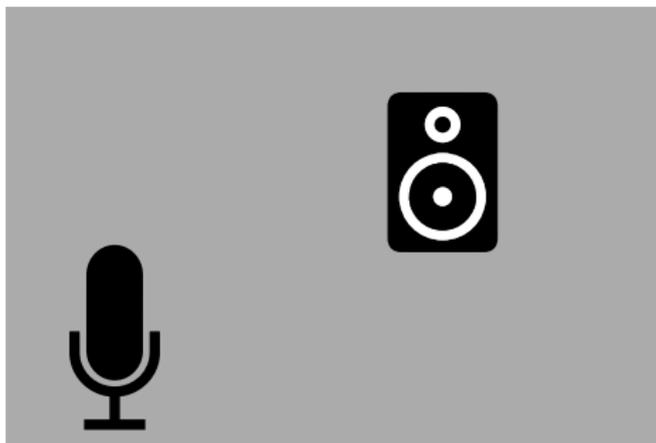
1. Assume knowledge of a few (1-6) early echoes
2. Plug into multichannel NMF <sup>1</sup>
3. Three baseline scenarios
  - *Anechoic* conditions
  - *Learn* transfer functions
  - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

---

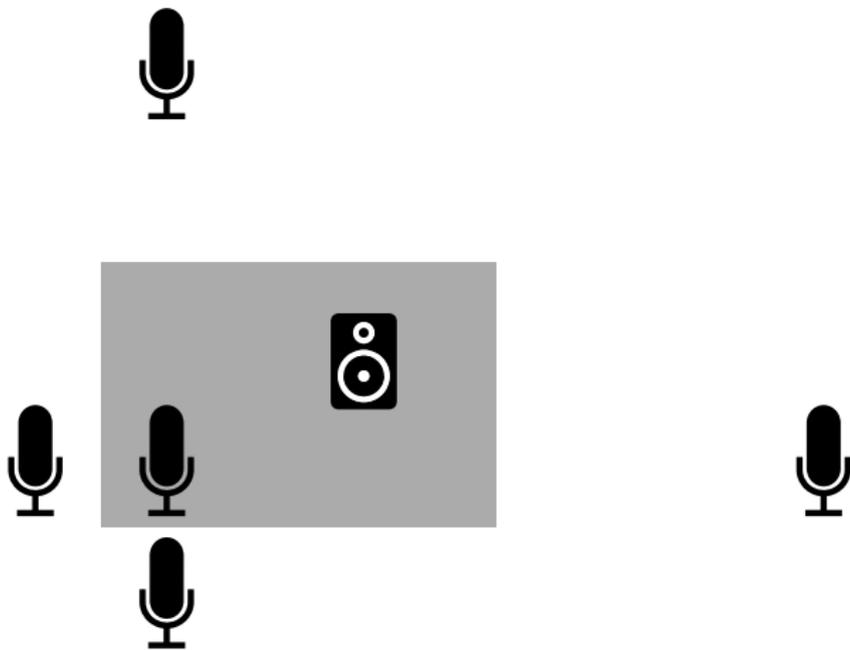
<sup>1</sup>Ozerov & Févotte, 2010

1. Approximate Propagation Model
2. NMF Algorithms
3. Results from Numerical Experiments

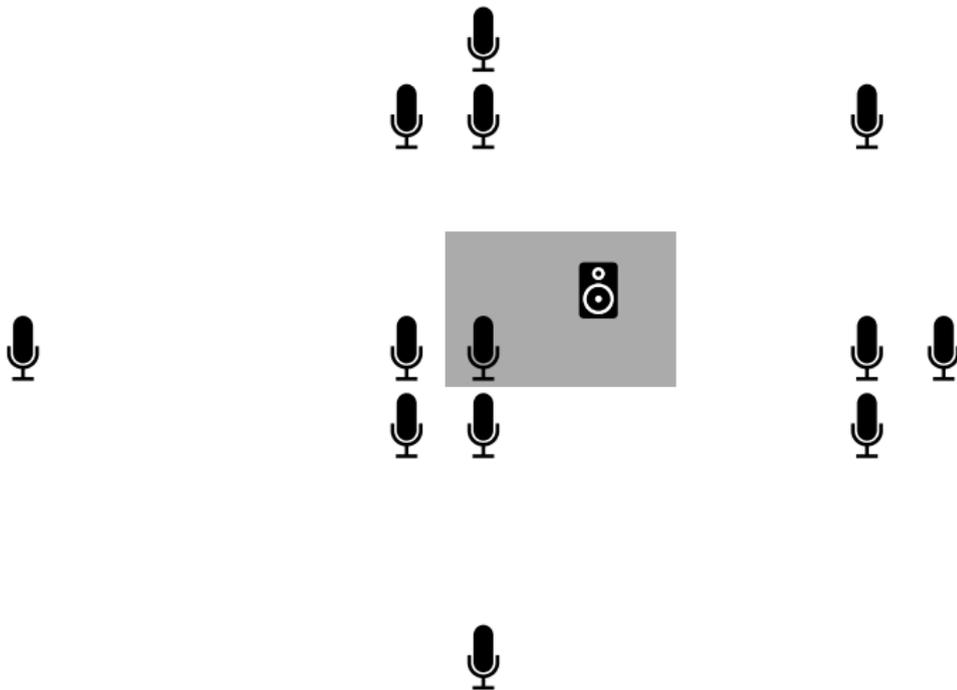
# Approximate Propagation Model



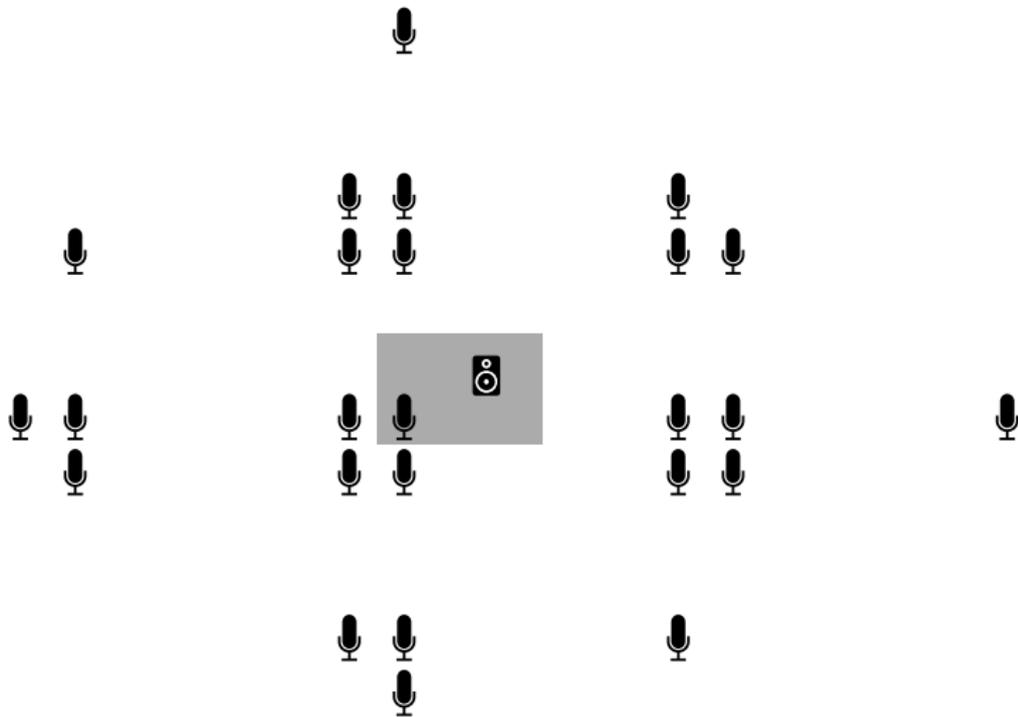
# Image Microphone Model



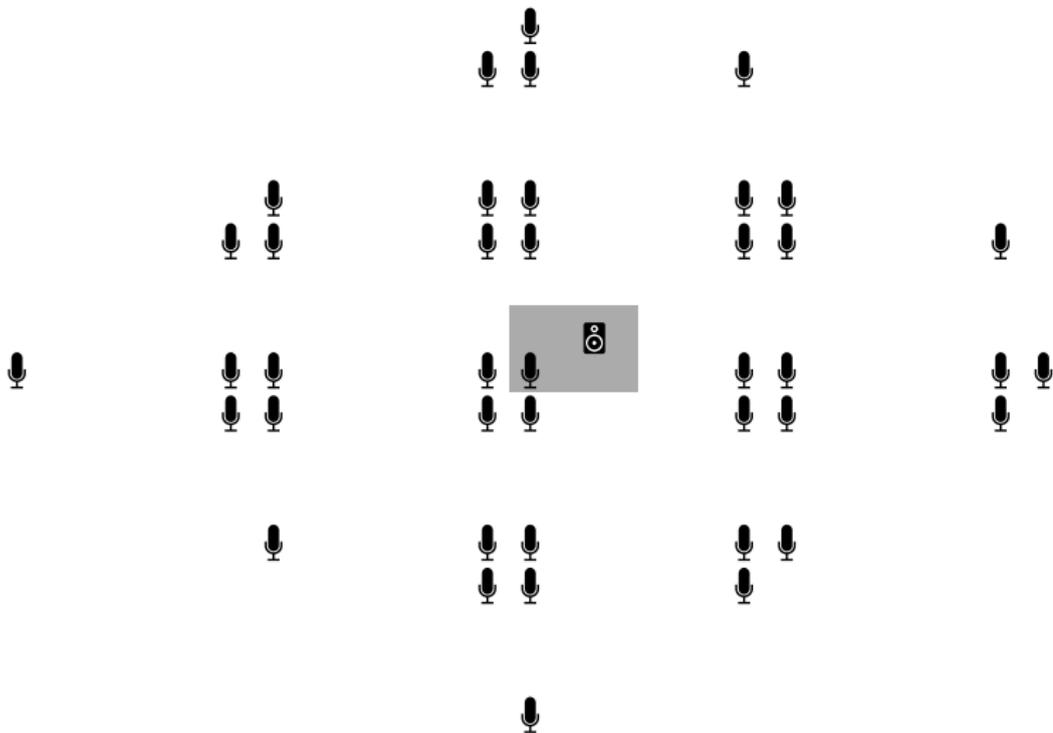
# Image Microphone Model



# Image Microphone Model

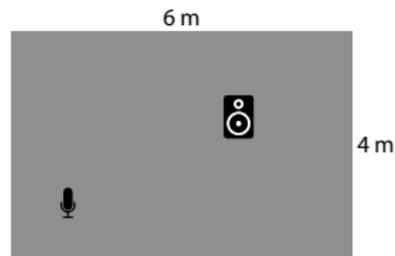


# Image Microphone Model



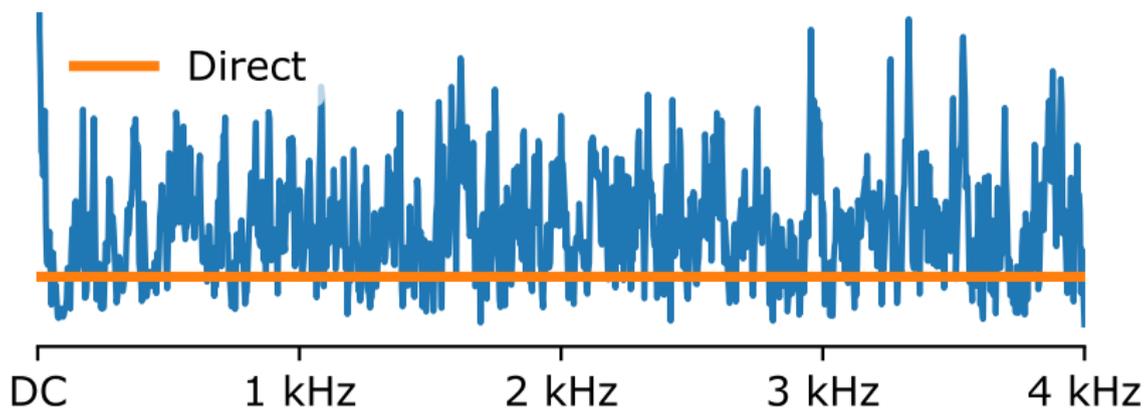
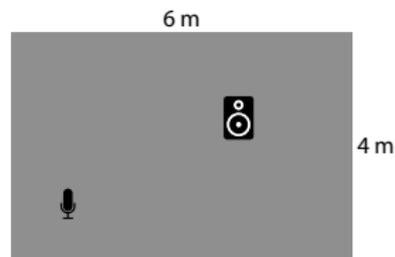
# Partial Room Impulse Responses

$$h_{jm}(t) = \sum_{k=0}^K \alpha_{jm}^k \delta(t - t_{jm}^k) + \epsilon_{jm}(t)$$



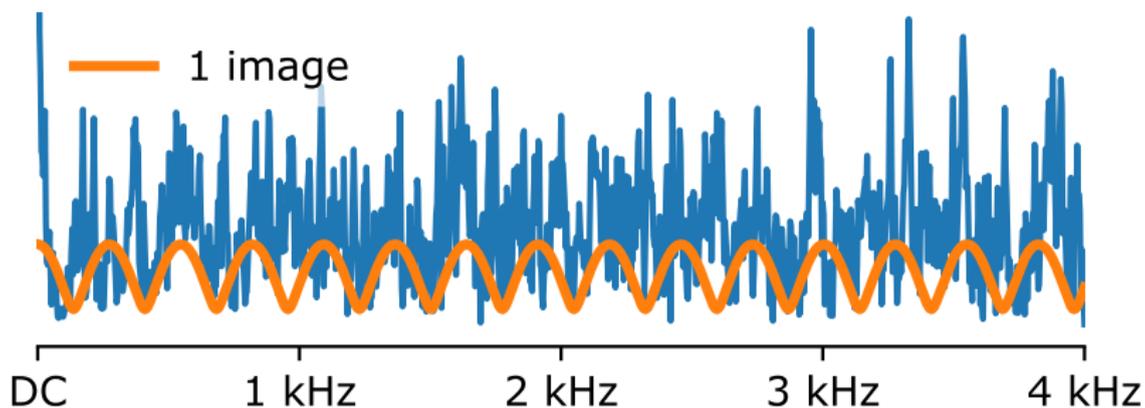
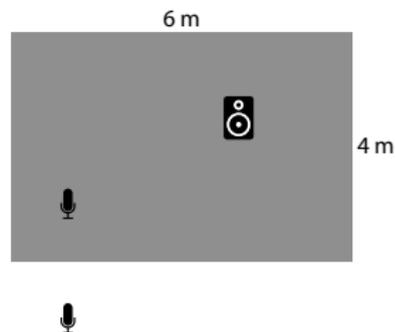
# Partial Room Impulse Responses

$$h_{jm}(t) = \sum_{k=0}^K \alpha_{jm}^k \delta(t - t_{jm}^k) + \epsilon_{jm}(t)$$



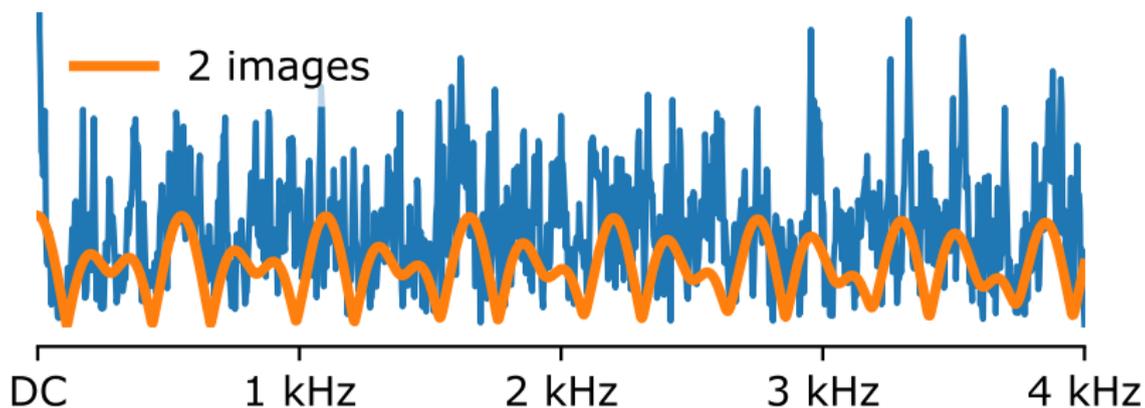
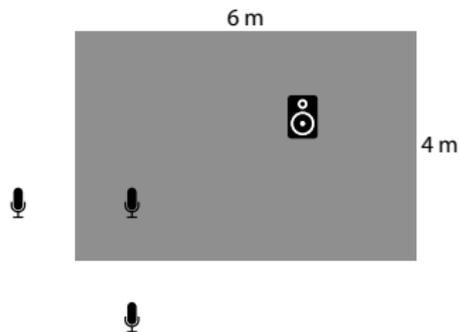
# Partial Room Impulse Responses

$$h_{jm}(t) = \sum_{k=0}^K \alpha_{jm}^k \delta(t - t_{jm}^k) + \epsilon_{jm}(t)$$



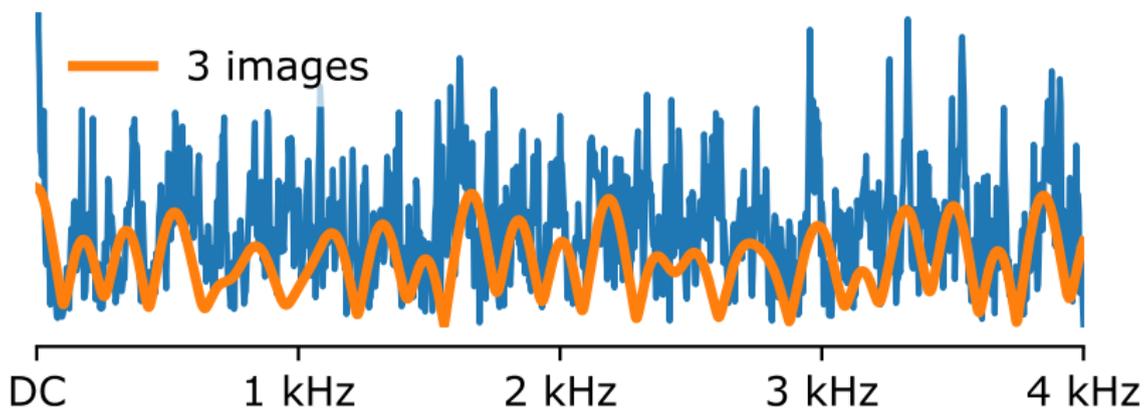
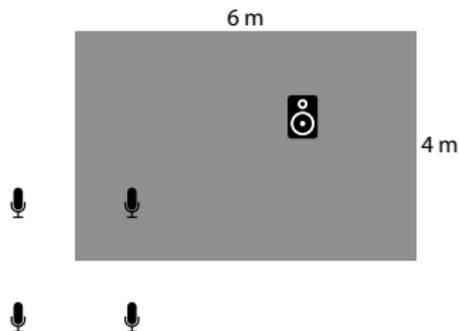
# Partial Room Impulse Responses

$$h_{jm}(t) = \sum_{k=0}^K \alpha_{jm}^k \delta(t - t_{jm}^k) + \epsilon_{jm}(t)$$



# Partial Room Impulse Responses

$$h_{jm}(t) = \sum_{k=0}^K \alpha_{jm}^k \delta(t - t_{jm}^k) + \epsilon_{jm}(t)$$



# Why should that help ?



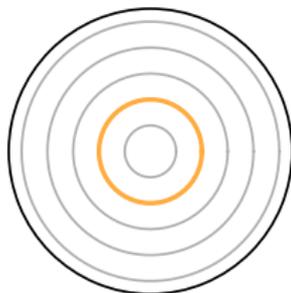
1000 Hz



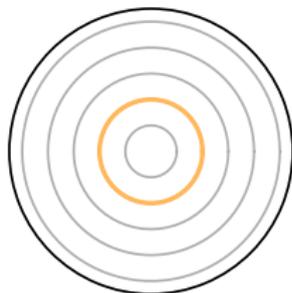
2000 Hz



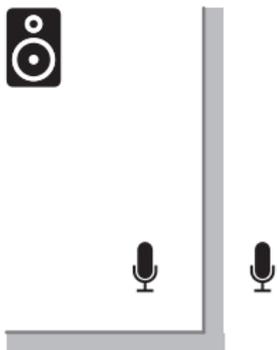
3000 Hz



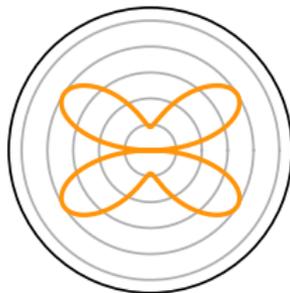
4000 Hz



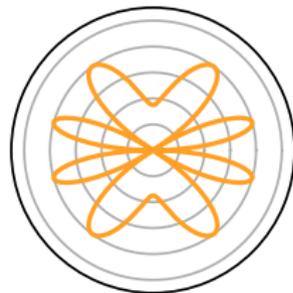
# Why should that help ?



1000 Hz



2000 Hz



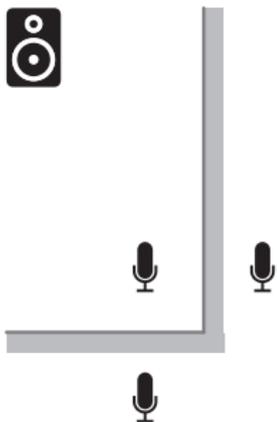
3000 Hz



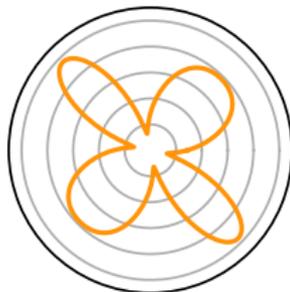
4000 Hz



# Why should that help ?



1000 Hz



2000 Hz



3000 Hz

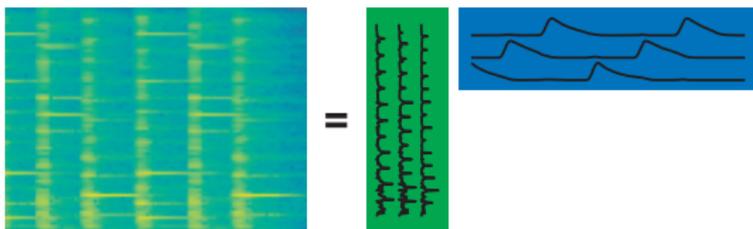


4000 Hz



# NMF Algorithms

# Non-negative Spectrogram Source Model



Multiplicative Updates View (Lee & Seung 2001)

Source signal's **magnitude spectrogram** decomposes non-negatively

$$|\mathbf{X}_j| = \mathbf{D}_j \mathbf{Z}_j$$

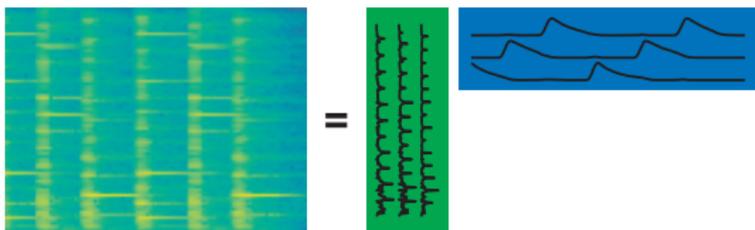
Expectation Maximization View (Ozerov & Févotte 2010)

Source signal's **variance spectrogram** decomposes non-negatively

$$X_j[f, n] \sim \mathcal{CN}(0, (\mathbf{D}_j \mathbf{Z}_j)_{fn})$$

In this work:  $\mathbf{D}_j$  is pre-trained, known dictionary

# Non-negative Spectrogram Source Model



## Multiplicative Updates View (Lee & Seung 2001)

Source signal's **magnitude spectrogram** decomposes non-negatively

$$|\mathbf{X}_j| = \mathbf{D}_j \mathbf{Z}_j$$

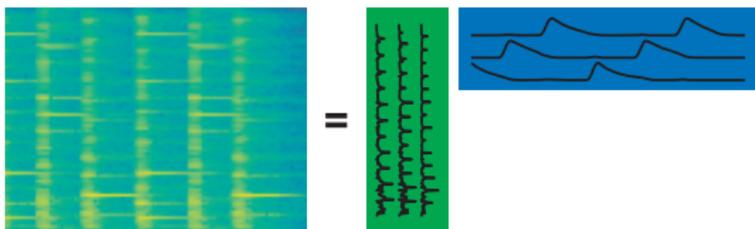
## Expectation Maximization View (Ozerov & Févotte 2010)

Source signal's **variance spectrogram** decomposes non-negatively

$$X_j[f, n] \sim \mathcal{CN}(0, (\mathbf{D}_j \mathbf{Z}_j)_{fn})$$

In this work:  $\mathbf{D}_j$  is pre-trained, known dictionary

# Non-negative Spectrogram Source Model



## Multiplicative Updates View (Lee & Seung 2001)

Source signal's **magnitude spectrogram** decomposes non-negatively

$$|X_j| = D_j Z_j$$

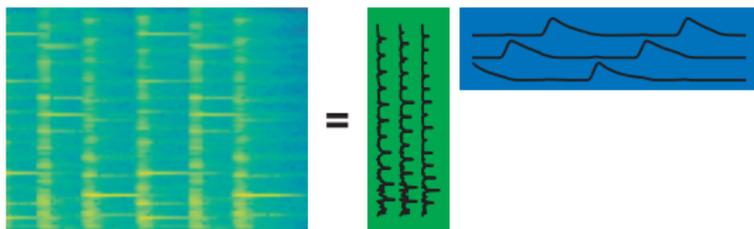
## Expectation Maximization View (Ozerov & Févotte 2010)

Source signal's **variance spectrogram** decomposes non-negatively

$$X_j[f, n] \sim \mathcal{CN}(0, (D_j Z_j)_{fn})$$

In this work:  $D_j$  is pre-trained, known dictionary

# Non-negative Spectrogram Source Model



## Multiplicative Updates View (Lee & Seung 2001)

Source signal's **magnitude spectrogram** decomposes non-negatively

$$|\mathbf{X}_j| = \mathbf{D}_j \mathbf{Z}_j$$

## Expectation Maximization View (Ozerov & Févotte 2010)

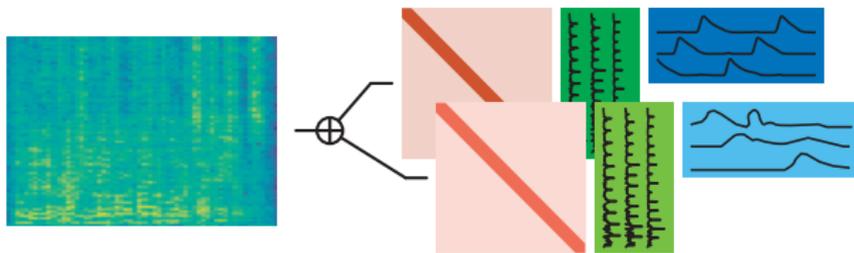
Source signal's **variance spectrogram** decomposes non-negatively

$$X_j[f, n] \sim \mathcal{CN}(0, (\mathbf{D}_j \mathbf{Z}_j)_{fn})$$

In this work:  $\mathbf{D}_j$  is pre-trained, known dictionary

## Microphone magnitude spectrogram model

$$\hat{\mathbf{V}}_m = \sum_j \text{diag}(|\hat{\mathbf{H}}_{mj}|) \mathbf{D}_j \mathbf{Z}_j$$



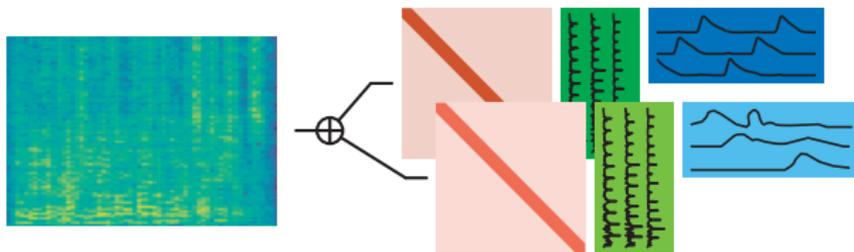
Minimize *Itakura-Saito* divergence

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mf n} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

- Efficient multiplicative update rules (Ozerov & Févotte 2010)
- Regularization needed for large number of latent variables

## Microphone magnitude spectrogram model

$$\hat{\mathbf{V}}_m = \sum_j \text{diag}(|\hat{\mathbf{H}}_{mj}|) \mathbf{D}_j \mathbf{Z}_j$$



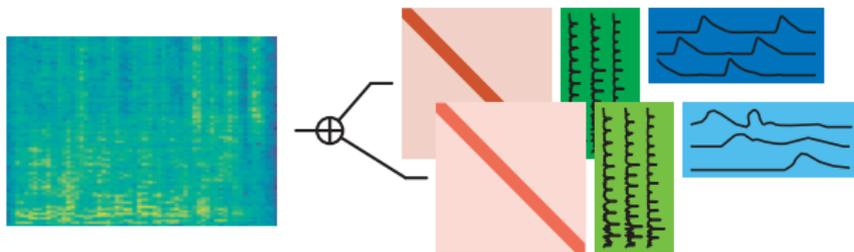
## Minimize *Itakura-Saito* divergence

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mf n} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

- Efficient **multiplicative update** rules (Ozerov & Févotte 2010)
- **Regularization** needed for large number of latent variables

## Microphone magnitude spectrogram model

$$\hat{\mathbf{V}}_m = \sum_j \text{diag}(|\hat{\mathbf{H}}_{mj}|) \mathbf{D}_j \mathbf{Z}_j$$



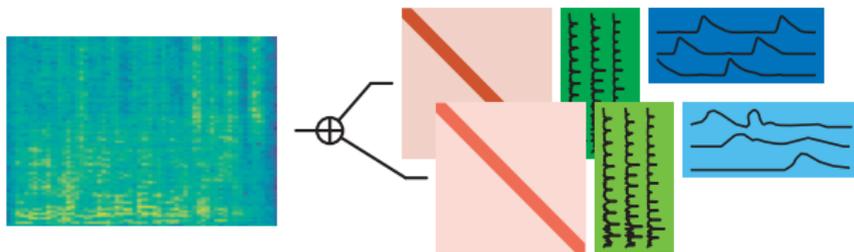
## Minimize *Itakura-Saito* divergence

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mf n} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

- Efficient **multiplicative update** rules (Ozerov & Févotte 2010)
- **Regularization** needed for large number of latent variables

## Microphone magnitude spectrogram model

$$\hat{\mathbf{V}}_m = \sum_j \text{diag}(|\hat{\mathbf{H}}_{mj}|) \mathbf{D}_j \mathbf{Z}_j$$



## Minimize *Itakura-Saito* divergence

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mf n} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

- Efficient **multiplicative update** rules (Ozerov & Févotte 2010)
- **Regularization** needed for large number of latent variables

## Probabilistic Model

Source are complex Gaussian with low-rank spectrogram

$$X_j[f, n] \sim \mathcal{CN}(0, (\mathbf{D}_j \mathbf{Z}_j)_{fn})$$

Microphone signals have variance

$$\Sigma_y[f, n] = \hat{\mathbf{H}}[f] \Sigma_x[f, n] \hat{\mathbf{H}}^H[f] + \Sigma_b[f, n],$$

Minimize Negative Log-likelihood

$$C_{EM}(\mathbf{Z}_j) = \sum_{fn} \text{trace} \left( \mathbf{y}[f, n] \mathbf{y}[f, n]^H \Sigma_y^{-1}[f, n] \right) + \log \det \Sigma_y[f, n]$$

Efficiently minimized by **Expectation-Maximization** algorithm  
(Ozerov & Févotte 2010)

## Probabilistic Model

Source are complex Gaussian with low-rank spectrogram

$$X_j[f, n] \sim \mathcal{CN}(0, (\mathbf{D}_j \mathbf{Z}_j)_{fn})$$

Microphone signals have variance

$$\boldsymbol{\Sigma}_y[f, n] = \hat{\mathbf{H}}[f] \boldsymbol{\Sigma}_x[f, n] \hat{\mathbf{H}}^H[f] + \boldsymbol{\Sigma}_b[f, n],$$

## Minimize Negative Log-likelihood

$$C_{EM}(\mathbf{Z}_j) = \sum_{fn} \text{trace} \left( \mathbf{y}[f, n] \mathbf{y}[f, n]^H \boldsymbol{\Sigma}_y^{-1}[f, n] \right) + \log \det \boldsymbol{\Sigma}_y[f, n]$$

Efficiently minimized by Expectation-Maximization algorithm  
(Ozerov & Févotte 2010)

## Probabilistic Model

Source are complex Gaussian with low-rank spectrogram

$$X_j[f, n] \sim \mathcal{CN}(0, (\mathbf{D}_j \mathbf{Z}_j)_{fn})$$

Microphone signals have variance

$$\boldsymbol{\Sigma}_y[f, n] = \hat{\mathbf{H}}[f] \boldsymbol{\Sigma}_x[f, n] \hat{\mathbf{H}}^H[f] + \boldsymbol{\Sigma}_b[f, n],$$

## Minimize Negative Log-likelihood

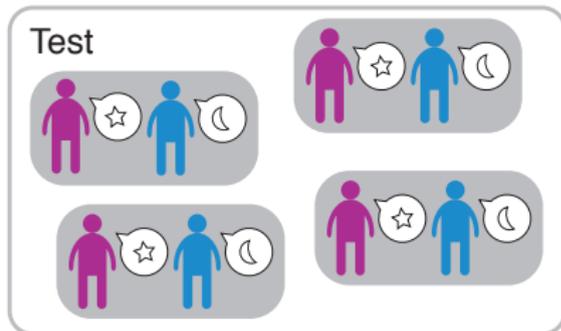
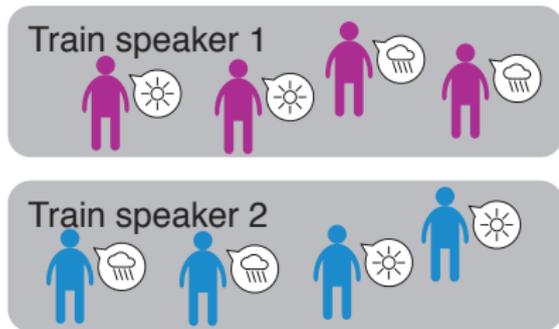
$$C_{EM}(\mathbf{Z}_j) = \sum_{fn} \text{trace} \left( \mathbf{y}[f, n] \mathbf{y}[f, n]^H \boldsymbol{\Sigma}_y^{-1}[f, n] \right) + \log \det \boldsymbol{\Sigma}_y[f, n]$$

Efficiently minimized by **Expectation-Maximization** algorithm  
(Ozerov & Févotte 2010)

Speaker Dependent

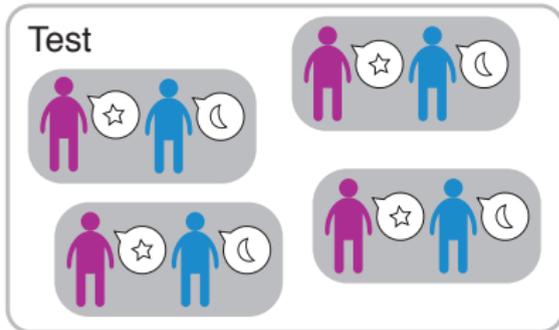
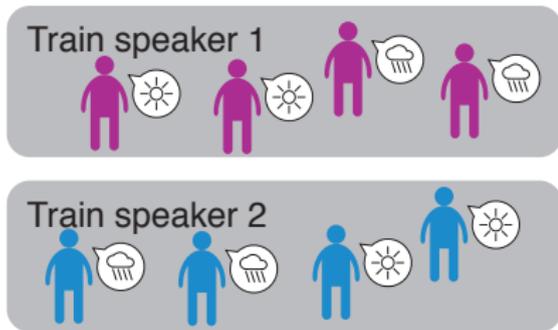
Universal

## Speaker Dependent

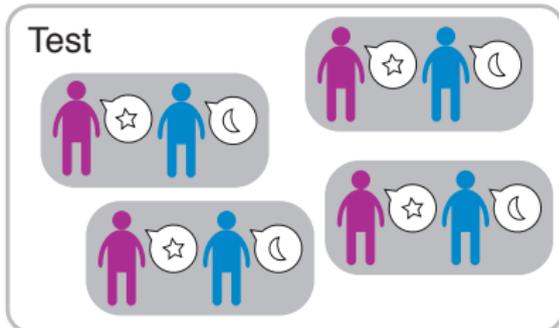


Universal

## Speaker Dependent



## Universal



Remark 1: Anechoic separation *cannot* work!

$$\hat{\mathbf{V}}_m = \sum_j \mathbf{D}_j \mathbf{Z}_j \quad \rightarrow \quad \hat{\mathbf{V}}_m = \sum_j \mathbf{D} \mathbf{Z}_j = \mathbf{D} \sum_j \mathbf{Z}_j$$

Remark 2: TF makes universal dict. speaker specific

$$\hat{\mathbf{V}}_m = \sum_j (\mathbf{H}_{mj} \mathbf{D}) \mathbf{Z}_j$$

Remark 3: EM-NMF with Universal Dictionary

- Unclear how to enforce sparsity in EM (to us)
- Left for future work

Remark 1: Anechoic separation *cannot* work!

$$\hat{\mathbf{V}}_m = \sum_j \mathbf{D}_j \mathbf{Z}_j \quad \rightarrow \quad \hat{\mathbf{V}}_m = \sum_j \mathbf{D} \mathbf{Z}_j = \mathbf{D} \sum_j \mathbf{Z}_j$$

Remark 2: TF makes universal dict. speaker specific

$$\hat{\mathbf{V}}_m = \sum_j (\mathbf{H}_{mj} \mathbf{D}) \mathbf{Z}_j$$

Remark 3: EM-NMF with Universal Dictionary

- Unclear how to enforce sparsity in EM (to us)
- Left for future work

Remark 1: Anechoic separation *cannot* work!

$$\hat{\mathbf{V}}_m = \sum_j \mathbf{D}_j \mathbf{Z}_j \quad \rightarrow \quad \hat{\mathbf{V}}_m = \sum_j \mathbf{D} \mathbf{Z}_j = \mathbf{D} \sum_j \mathbf{Z}_j$$

Remark 2: TF makes universal dict. speaker specific

$$\hat{\mathbf{V}}_m = \sum_j (\mathbf{H}_{mj} \mathbf{D}) \mathbf{Z}_j$$

Remark 3: EM-NMF with Universal Dictionary

- Unclear how to enforce sparsity in EM (to us)
- Left for future work

# Results from Numerical Experiments

## Conditions

# sources 2

# mics 3

STFT 2048

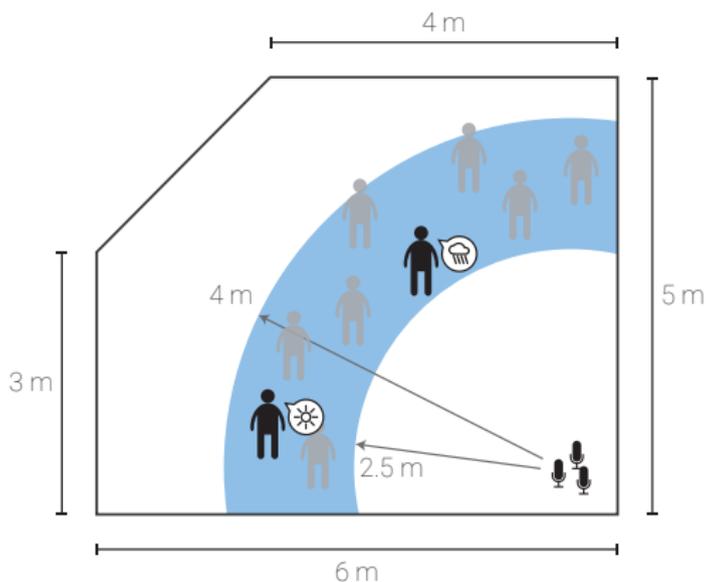
half-overlap, Hann win

Simulation with  
*pyroomacoustics*

T60  $\sim 100$  ms

## Baselines

- *Anechoic*
- *Learn TF*
- *Ignore reverb*



## Conditions

# sources 2

# mics 3

STFT 2048

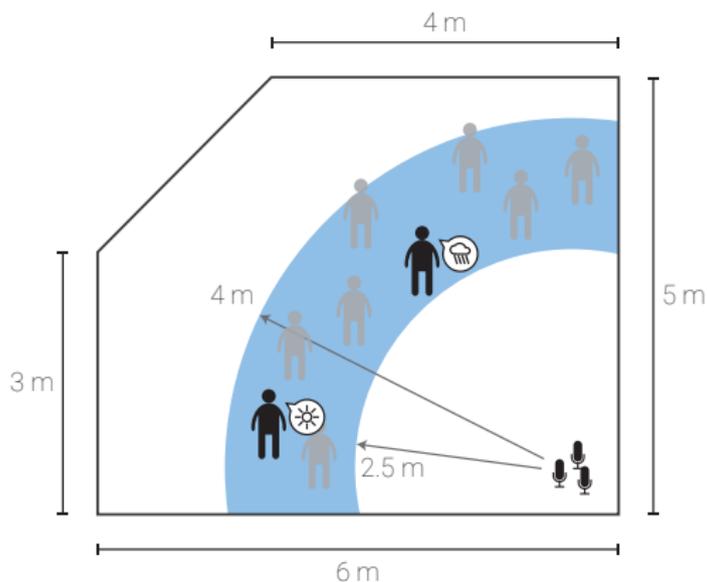
half-overlap, Hann win

Simulation with  
*pyroomacoustics*

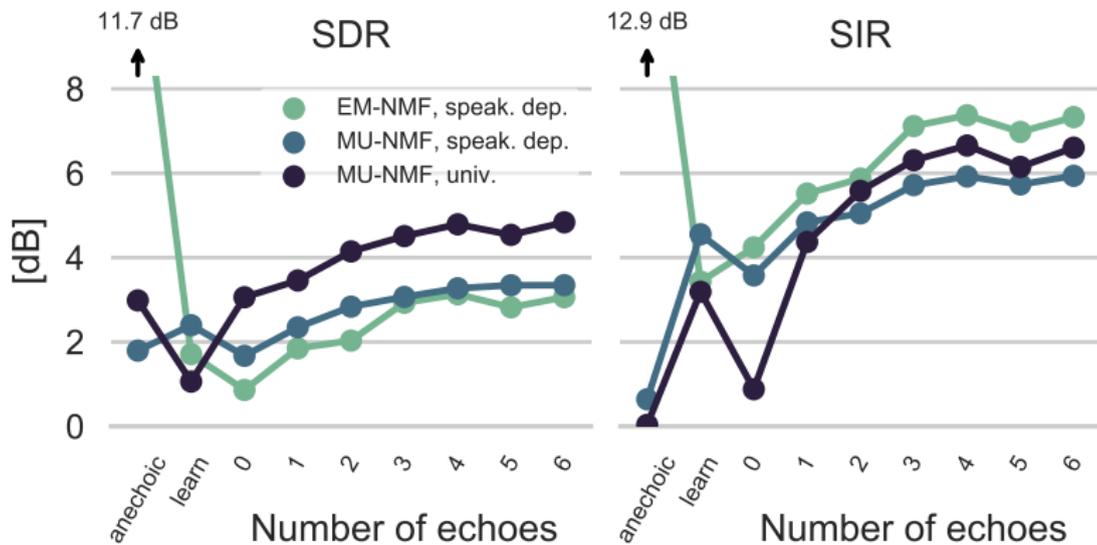
T60  $\sim 100$  ms

## Baselines

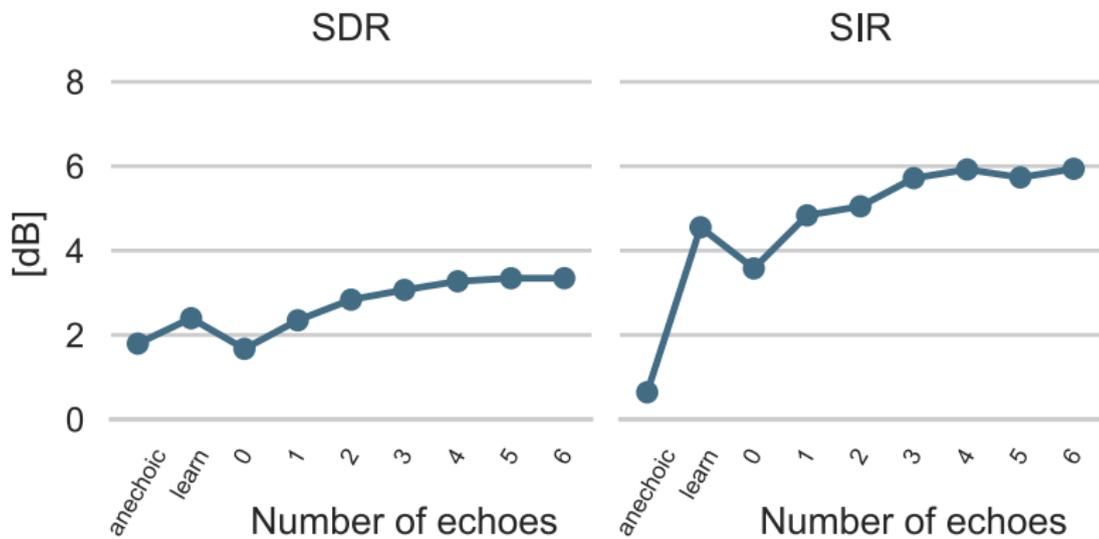
- *Anechoic*
- *Learn TF*
- *Ignore reverb*



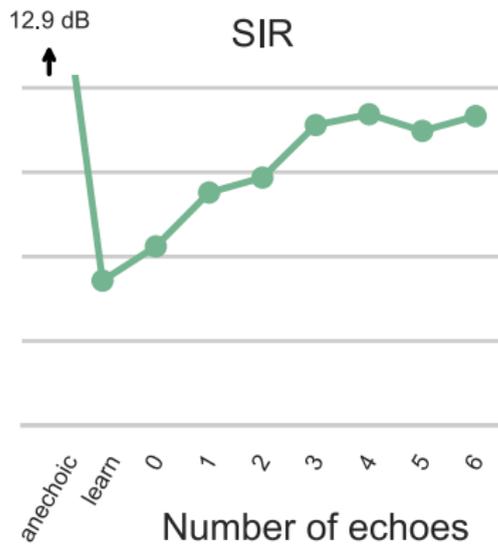
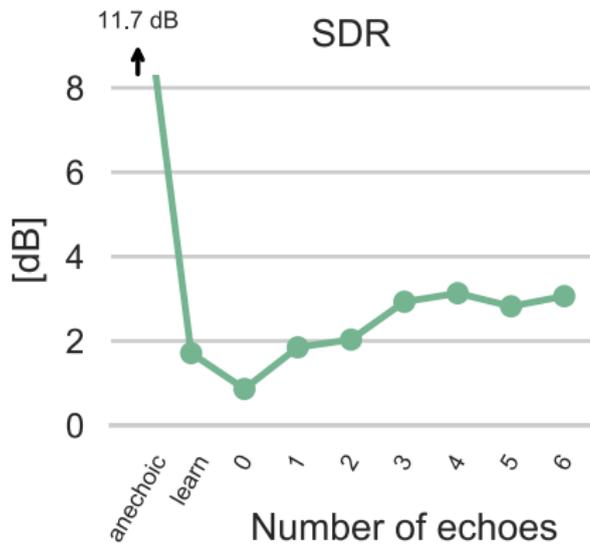
# Numerical Experiments Results

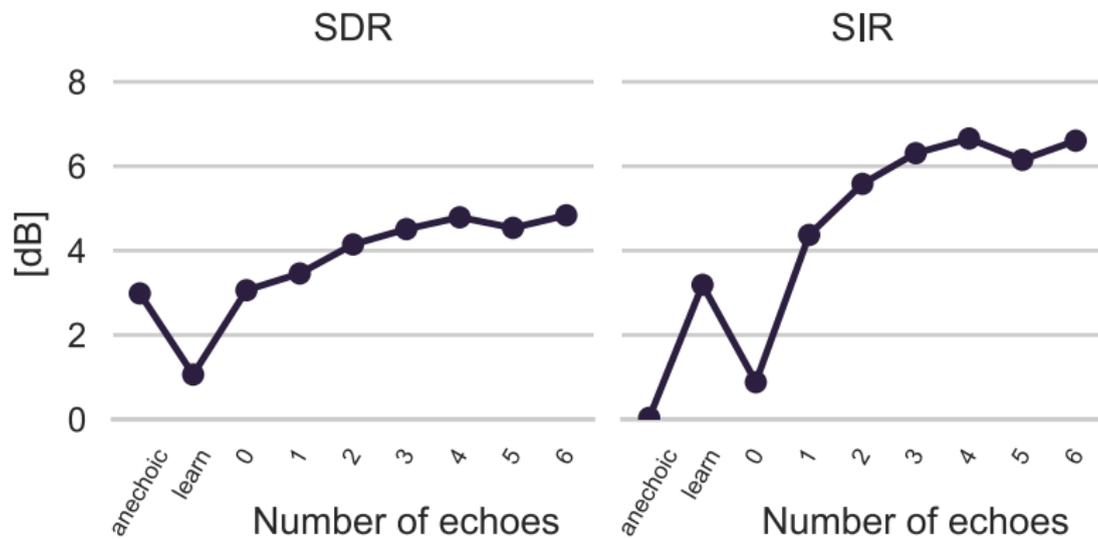


# MU-NMF – Speaker Dependent



# EM-NMF – Speaker Dependent





## Recall

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mfn} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

		Number of echoes $K$								
		anechoic	learn	0	1	2	3	4	5	6
$\gamma =$		10	$10^{-1}$	10	$10^{-4}$	0	0	0	0	0

Table : Value of regularization parameter.

Partial RIR regularizes universal dictionary!

## Recall

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mfn} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

		Number of echoes $K$								
		anechoic	learn	0	1	2	3	4	5	6
$\gamma =$		10	$10^{-1}$	10	$10^{-4}$	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>

Table : Value of regularization parameter.

Partial RIR regularizes universal dictionary!

## Recall

$$C_{\text{MU}}(\mathbf{Z}_j) = \sum_{mfn} d_{\text{IS}}(V_m[f, n] | \hat{V}_m[f, n]) + \gamma \sum_j \|\mathbf{Z}_j\|_1$$

		Number of echoes $K$								
		anechoic	learn	0	1	2	3	4	5	6
$\gamma =$		10	$10^{-1}$	10	$10^{-4}$	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>

Table : Value of regularization parameter.

Partial RIR regularizes universal dictionary!

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?

## Conclusion

- Single echo improves performance
- Enables universal dictionary
- First few echoes most important

## Future Work

- Compare to BSS
- Include (deeply) learnt models
- Underdetermined case

Thank you! Questions ?