



Separake

Source Separation with a Little Help from Echoes

Robin Scheibler Diego Di Carlo Antoine Deleforge Ivan Dokmanić
DECEMBER 3, 2020



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 ?

Full RIRs or a few early reflections ?

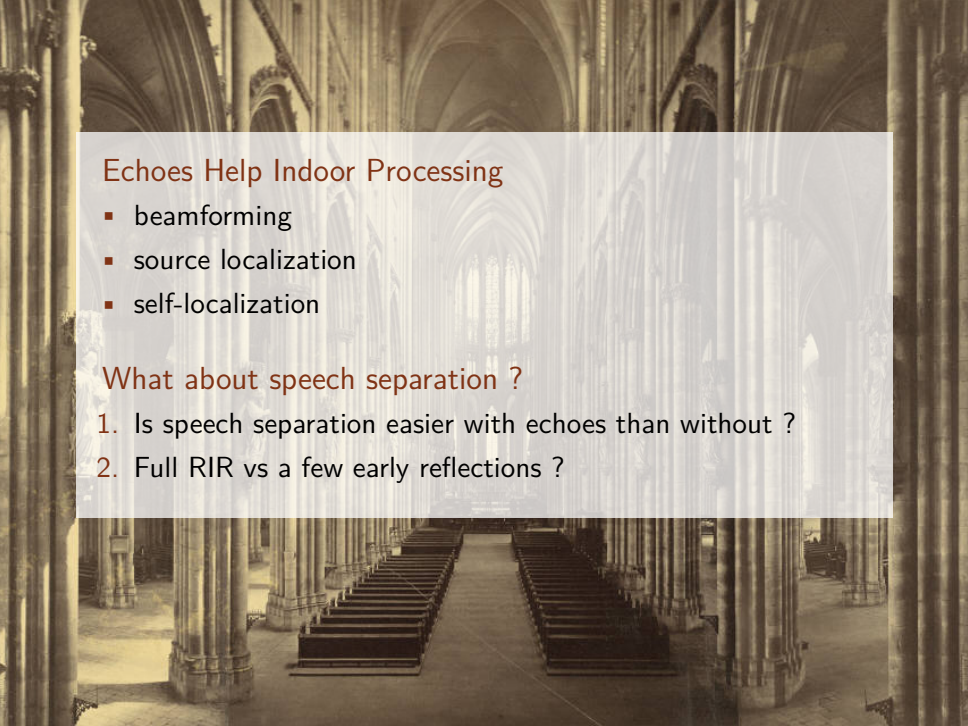
Echoes Help Indoor Processing

- beamforming
- source localization
- self-localization

What about speech separation ?

1. Is speech separation easier with echoes than without ?

Full R Room Acoustics References ?



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 in multichannel NMF ¹
3. Three baseline scenarios
 - *Anechoic* conditions
 - *Learn* transfer functions
 - Ignore reverberation (i.e. consider 0 echoes)
4. Numerical Experiments

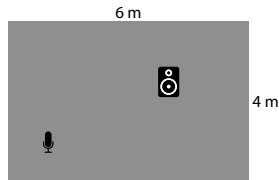
¹Ozerov & Févotte, 2010

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

Full Image Microphone Model

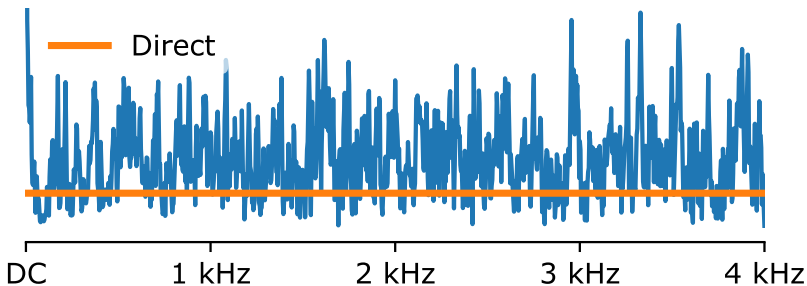
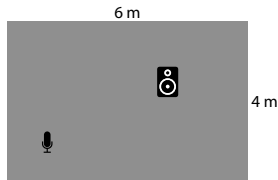
Partial Room Impulse Responses

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



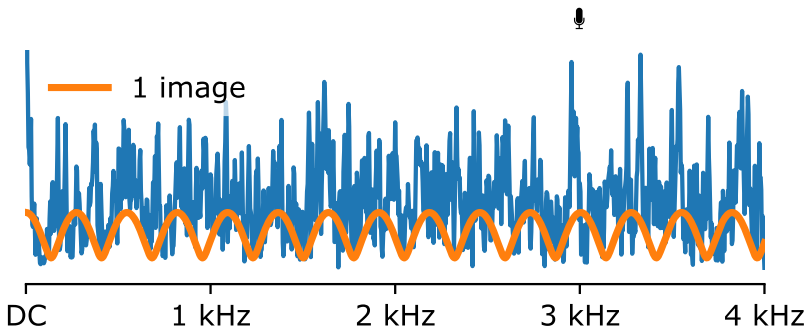
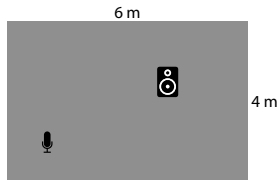
Partial Room Impulse Responses

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



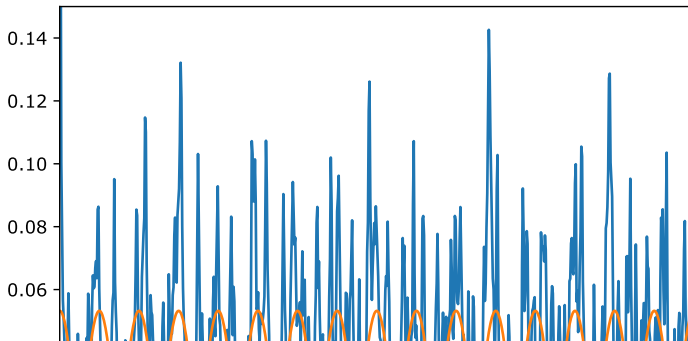
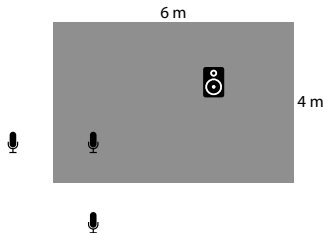
Partial Room Impulse Responses

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



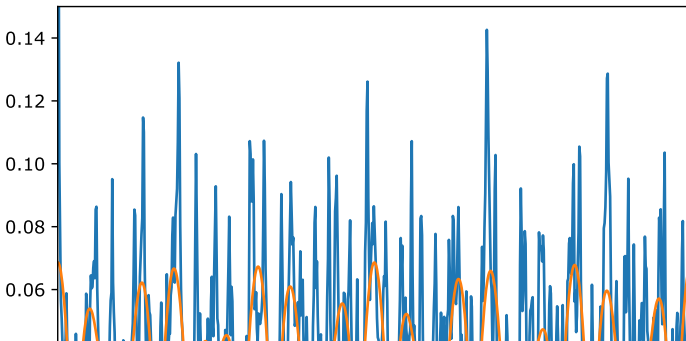
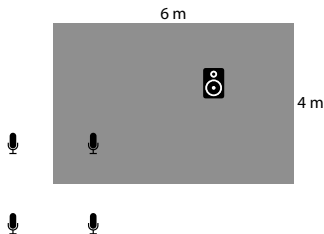
Partial Room Impulse Responses

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

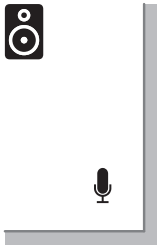


Partial Room Impulse Responses

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



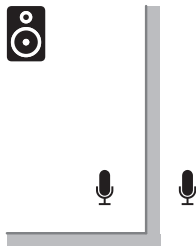
Why should that help ?



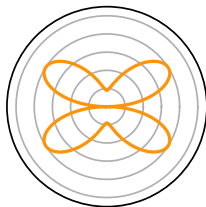
Why should that help ?



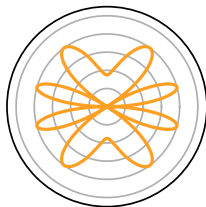
Why should that help ?



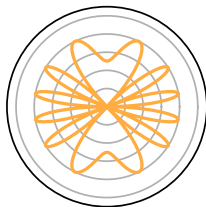
1000 Hz



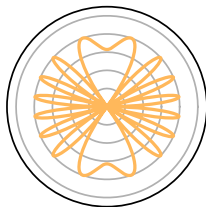
2000 Hz



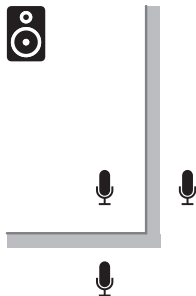
3000 Hz



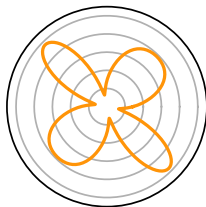
4000 Hz



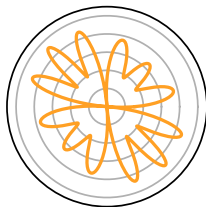
Why should that help ?



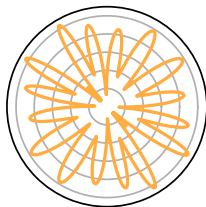
1000 Hz



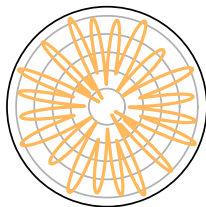
2000 Hz



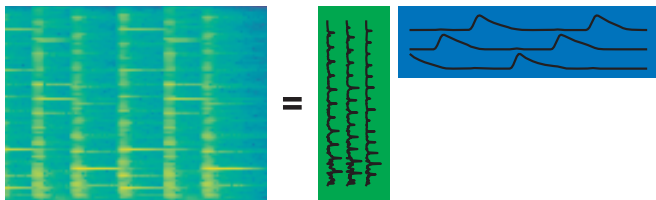
3000 Hz



4000 Hz



Non-negative Spectrogram Source Model



Multiplicative Updates View (Lee & Seung 2001)

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

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

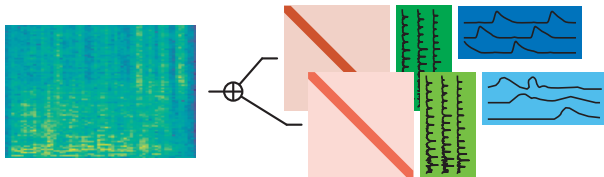
Expectation Minimization View (Ozerov & Févotte 2010)

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

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

Microphone magnitude spectrogram model

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



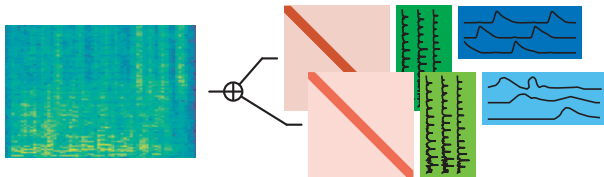
Minimize *Itakura-Saito* divergence

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

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

Microphone magnitude spectrogram model

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



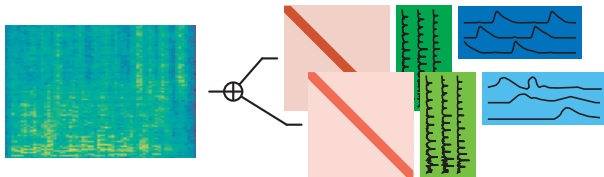
Minimize *Itakura-Saito* divergence

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

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

Microphone magnitude spectrogram model

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



Minimize *Itakura-Saito* divergence

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

- 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

$$\Sigma_{\mathbf{y}}[f, n] = \hat{\mathbf{H}}[f] \Sigma_{\mathbf{x}}[f, n] \hat{\mathbf{H}}^H[f] + \Sigma_{\mathbf{b}}[f, n],$$

Minimize Negative Log-likelihood

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

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

Speaker Dependent

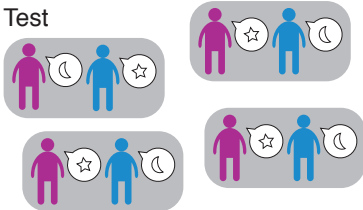
Train speaker 1



Train speaker 2



Test

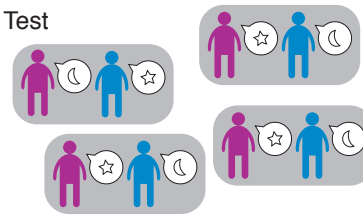


Universal

Train



Test



Remark 1: Anechoic separation with MU-NMF

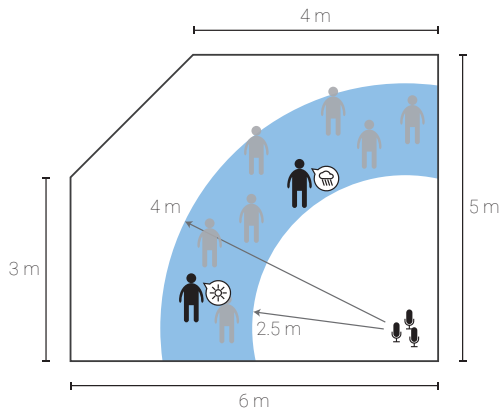
Anechoic separation **cannot work!**

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

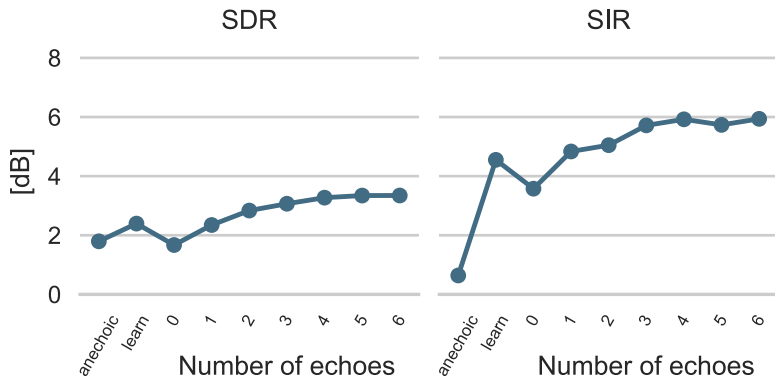
Remark 2: EM-NMF with Universal Dictionary

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

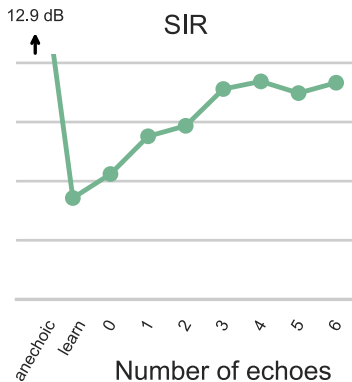
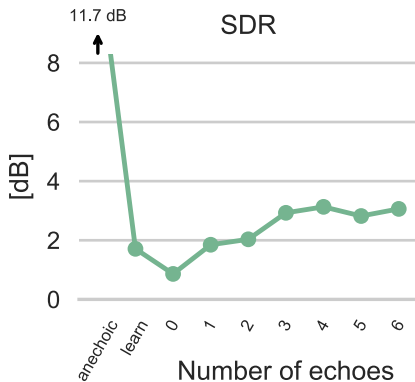
Experimental Setup

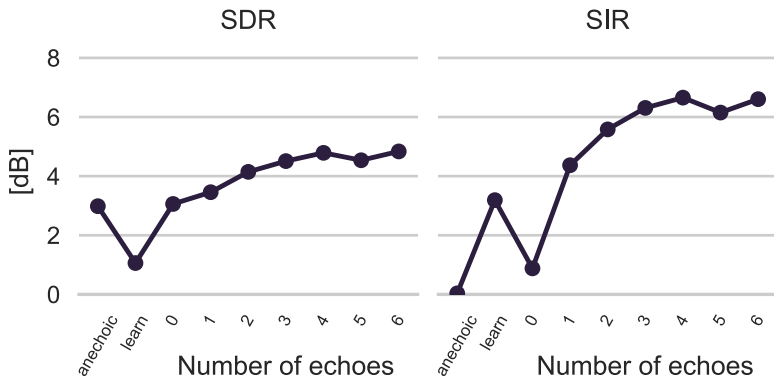


MU-NMF – Speaker Dependent



EM-NMF – Speaker Dependent





Conclusion

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

Future Work

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

Conclusion

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

Future Work

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

Conclusion

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

Future Work

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

Conclusion

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

Future Work

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

Conclusion

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

Future Work

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

Conclusion

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

Future Work

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

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 ?

Numerical Experiments Results

