



université  
PARIS-SACLAY

# Modeling Reverberant Mixtures for Multichannel Audio-Source Separation

Simon Leglaive

Ph.D. supervisors: Roland Badeau and Gaël Richard

LTCI, Télécom ParisTech, Université Paris Saclay

Seminar at PERCEPTION group, INRIA Grenoble Rhône-Alpes,  
Montbonnot Saint-Martin

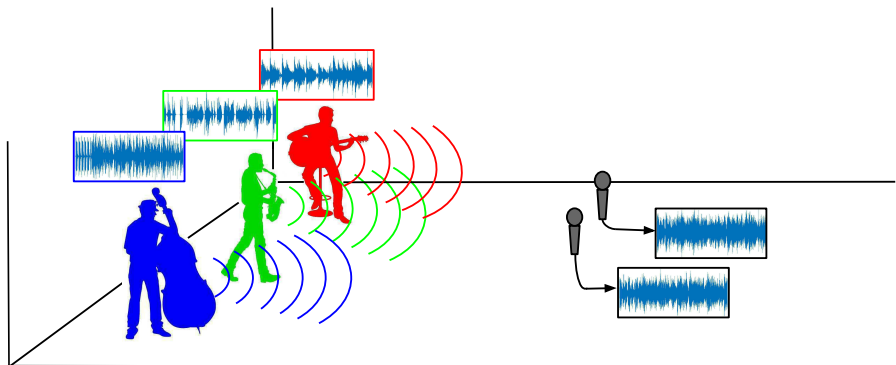
March 22, 2017



# Multichannel audio source separation

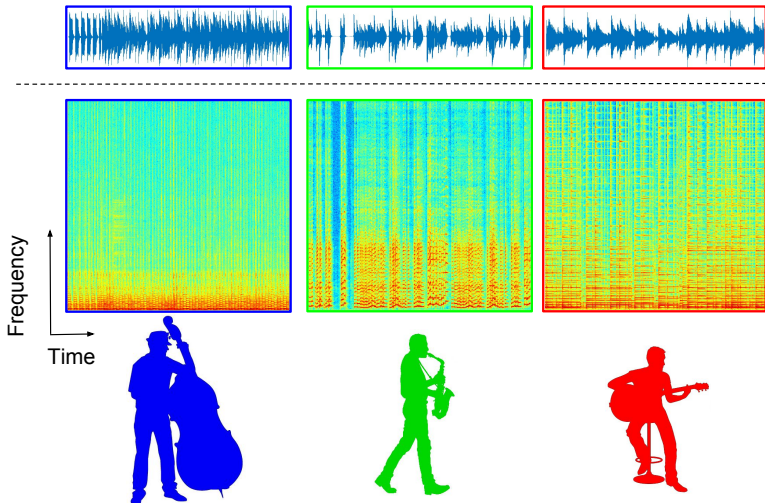
Objective: Recover source signals from the observation of several mixtures.

Context: Under-determined and reverberant.



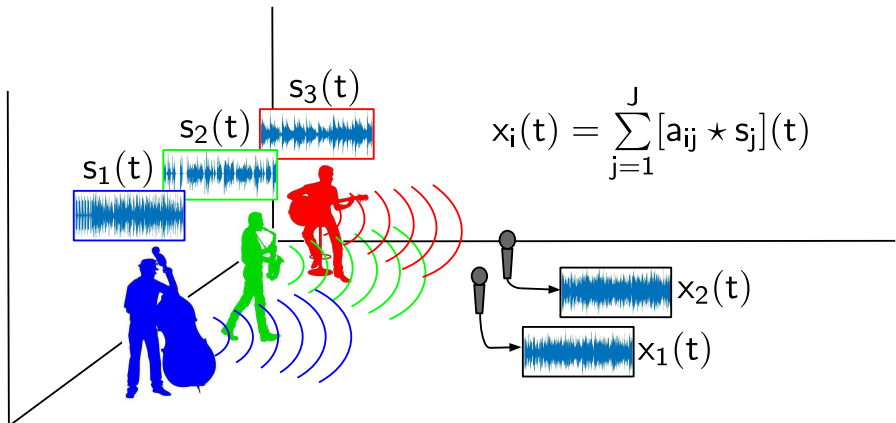
# Time-frequency source representation

Time-frequency (TF) transforms provide meaningful representations.



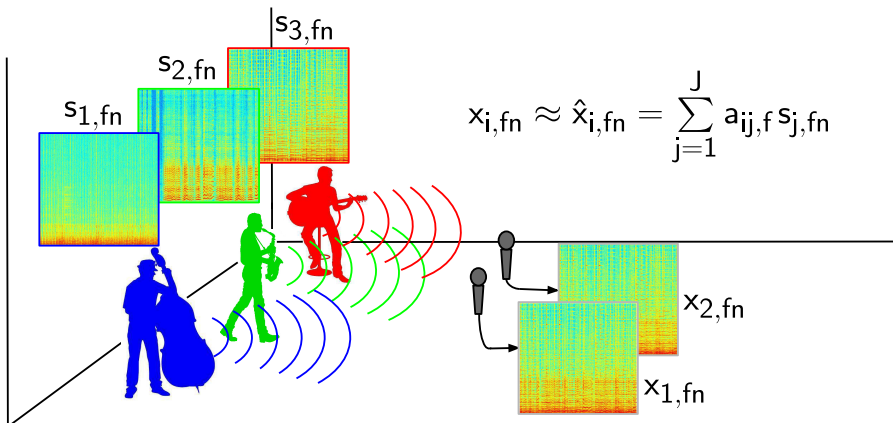
# Modeling reverberant mixtures (1)

Convolutional model in the time domain:



## Modeling reverberant mixtures (2)

Convolutional model in the Short-Term Fourier Transform (STFT) domain:







# Outline

## Convolutional mixture model in the STFT domain

- Baseline source separation framework

- Room frequency response modeling

- Source separation with reverberation priors

- Experiments

- Limitations

## Convolutional mixture model in the time domain

- Model

- Inference

- Experiments

- Ongoing work

## Conclusion



# Outline

## Convolutional mixture model in the STFT domain

- Baseline source separation framework

- Room frequency response modeling

- Source separation with reverberation priors

- Experiments

- Limitations

## Convolutional mixture model in the time domain

- Model

- Inference

- Experiments

- Ongoing work

## Conclusion





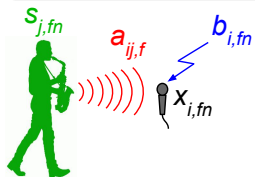
# Mixture model

Convolutional noisy mixture of  $J$  sources on  $I$  channels

STFT domain:

$$\forall (f, n) \in \{0, \dots, F - 1\} \times \{0, \dots, N - 1\}$$

$$x_{i,fn} = \sum_{j=1}^J a_{ij,f} s_{j,fn} + b_{i,fn}$$



- ▶ Frequency response of the mixing filters:  $a_{ij,f}$ ;
- ▶ Source STFT coefficients:  $s_{j,fn}$ ;
- ▶ Additive Gaussian noise:  $b_{i,fn} \sim \mathcal{N}_c(0, \sigma_{b,f}^2)$ .

In matrix form

$$\mathbf{x}_{fn} = \mathbf{A}_f \mathbf{s}_{fn} + \mathbf{b}_{fn}, \quad (1)$$

where  $\mathbf{x}_{fn} = [x_{i,fn}]_i \in \mathbb{C}^I$ ,  $\mathbf{s}_{fn} = [s_{j,fn}]_j \in \mathbb{C}^J$ ,  $\mathbf{A}_f = [a_{ij,f}]_{ij} \in \mathbb{C}^{I \times J}$  and  $\mathbf{b}_{fn} = [b_{i,fn}]_i \in \mathbb{C}^I$ .











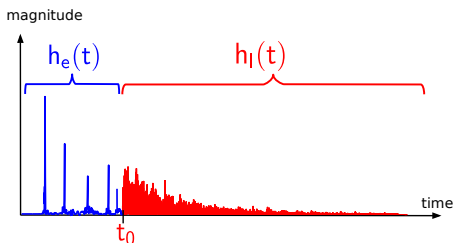
# Room impulse and frequency responses

## Room impulse and frequency responses

For  $t, f \in \{0, \dots, T - 1\}$ :

$$\underbrace{h(t) = h_e(t) + h_l(t)}_{\text{Room impulse response (RIR)}} \begin{matrix} \xrightarrow{\mathcal{F}_T} \\ \xleftarrow{\mathcal{F}_T^{-1}} \end{matrix} \underbrace{H(f) = H_e(f) + H_l(f)}_{\text{Room frequency response (RFR)}}$$

$\mathcal{F}_T$ : Discrete Fourier Transform (DFT)



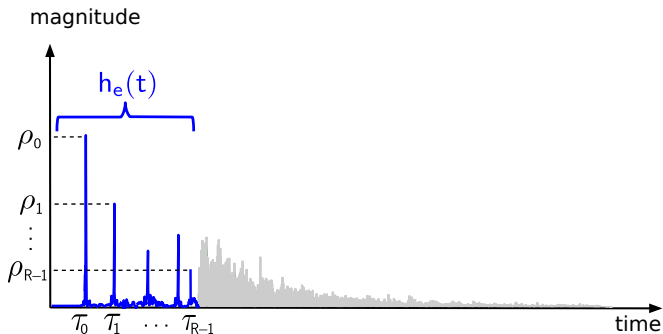
Mixing time:

$$t_0 = \lfloor C_0 \sqrt{V} f_s \rfloor \text{ samples}$$

- ▶  $C_0 = 2 \times 10^{-3}$
- ▶  $V$ : volume of the room
- ▶  $f_s$ : sampling rate



## Early contributions model (1)



$k$ -th early contribution: amplitude  $\rho_k$  and delay  $\tau_k$

$$H_e(f) = \sum_{k=0}^{R-1} \rho_k \delta_k^f \quad \text{with} \quad \delta_k = e^{-j2\pi\tau_k/T}. \quad (2)$$

## Early contributions model (2)

It follows (see, e.g., [Kumaresan, 1983])

$$\{H_e(f)\}_{f=R,\dots,T-1} \quad \text{satisfies} \quad \sum_{r=0}^R \varphi_r^e H_e(f-r) = 0, \quad (3)$$

where  $\{\varphi_r^e\}_{r=0}^R$  and  $\{\delta_k\}_{k=0}^{R-1}$  are the coefficients and roots of the same polynomial of order  $R$ .

Adding an error term  $\rightarrow$  autoregressive model

$$\sum_{r=0}^R \varphi_r^e H_e(f-r) = \kappa(f) \quad \text{with} \quad \kappa(f) \sim \mathcal{N}_c(0, \sigma_\kappa^2). \quad (4)$$

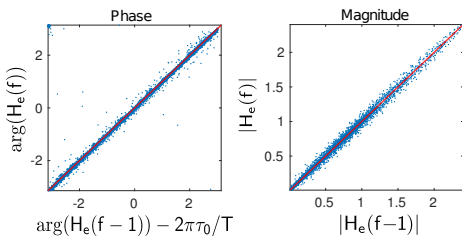
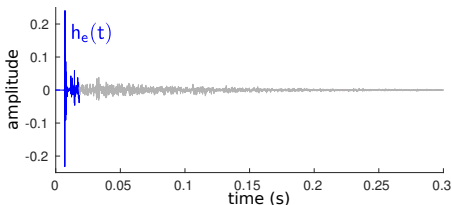
[Kumaresan, 1983] "On the zeros of the linear prediction-error filter for deterministic signals". *IEEE Transactions on Acoustics, Speech, and Signal Processing*.

**For more details:** [Leglaive et al., 2015] "Multichannel audio source separation with probabilistic reverberation modeling". *IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*.

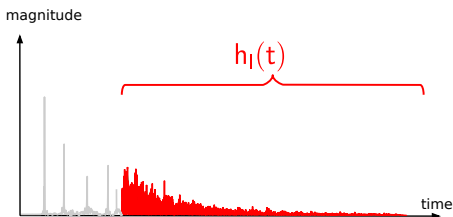
# Autoregressive model of order 1

Assuming that the direct path dominates the early echoes:

$$|H_e(f)| \approx |H_e(f - 1)| \quad \text{and} \quad \arg(H_e(f)) \approx \arg(H_e(f - 1)) - 2\pi \frac{\tau_0}{T}$$



## Late reverberation model (time domain)



Power Temporal Profil (PTP)  $\rightarrow$  exponential decay

$$\bar{h}_l(t) = \mathbb{E} [h_l(t)^2] \propto e^{-2t/\tau} \mathbb{1}_{t \geq t_0}(t),$$

- ▶  $\tau = \frac{T_{60} f_s}{3 \ln(10)}$  samples, with  $T_{60}$  the reverberation time in seconds;
- ▶  $\mathbb{E}[\cdot]$ : spatial averaging.

## Late reverberation model (frequency domain)

**Statistical room acoustics:**  $\{H_l(f)\}_f$  is a proper centered and WSS complex Gaussian random process.

## Late reverberation model (frequency domain)

**Statistical room acoustics:**  $\{H_l(f)\}_f$  is a proper centered and WSS complex Gaussian random process.

### Theoretical Power Spectral Density (PSD)

We can show that the PSD is related to the PTP by:

$$\phi(t) = T\bar{h}_l(T - t).$$

## Late reverberation model (frequency domain)

**Statistical room acoustics:**  $\{H_l(f)\}_f$  is a proper centered and WSS complex Gaussian random process.

### Theoretical Power Spectral Density (PSD)

We can show that the PSD is related to the PTP by:

$$\phi(t) = T\bar{h}_l(T - t).$$

### Theoretical Autocovariance function (ACVF)

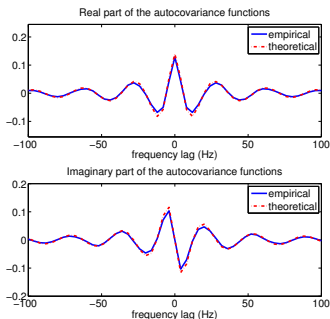
Applying the Wiener-Khinchin theorem we can obtain a theoretical expression of the ACVF:

$$\gamma(m) = \mathcal{F}_T^{-1}\{\phi(t)\}.$$

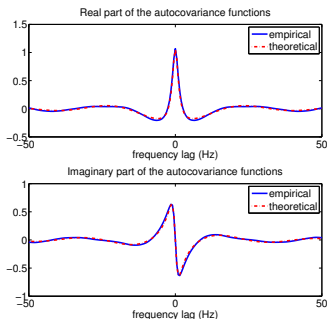
These quantities are theoretically defined according to some room parameters (reverberation time, dimensions).

## Experimental validation

Empirical autocovariance functions computed from a Monte-Carlo simulation on synthesized and real room responses.



(a) 196 simulated RIRs  
 $T_{60} = 0.25$  s  
 $10 \times 6.6 \times 3.3$  m



(b) 130 real RIRs  
 $T_{60} = 1.8$  s  
 $7.5 \times 9 \times 3.5$  m

Figure: Theoretical and empirical autocovariance functions



## ARMA parametrization

### ARMA representation of late reverberation in the frequency domain

We assume that  $\{H_l(f)\}_f$  follows an ARMA( $P, Q$ ) model:

$$\Phi(L)H_l(f) = \Theta(L)\epsilon(f),$$

- ▶  $\Phi(L) = \sum_{p=0}^P \varphi_p^l L^p$  and  $\Theta(L) = \sum_{q=0}^Q \theta_q L^q$  with  $\varphi_0^l = \theta_0 = 1$ ;
- ▶  $L$  is the lag operator, i.e.  $LH_l(f) = H_l(f - 1)$ ;
- ▶  $\epsilon(f) \sim \mathcal{N}_c(0, \sigma_\epsilon^2)$ .

We can compute the ARMA parameters from the theoretical ACVF.

# Experimental validation

Same room parameters as used before for simulated RIRs

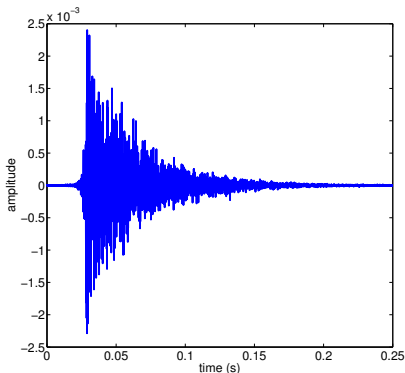
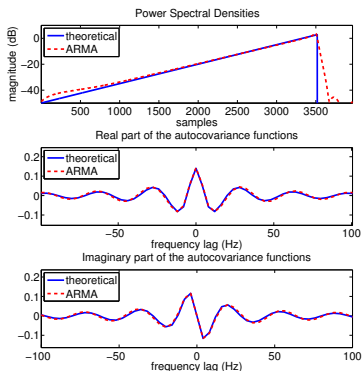


Figure: ARMA(7,2) parametrization

Figure: Synthesized late RIR

For more details: [Leglaive et al., 2016] "Autoregressive moving average modeling of late reverberation in the frequency domain". *European Signal Processing Conference (EUSIPCO)*.



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

Experiments

Ongoing work

## Conclusion

# EM algorithm

- Mixing matrix:  $\mathbf{A}_f = \underbrace{\mathbf{A}_{e,f}}_{\text{early reverb.}} + \underbrace{\mathbf{A}_{l,f}}_{\text{late reverb.}}$

## E-step

$$Q(\boldsymbol{\eta}|\boldsymbol{\eta}_{\text{old}}) = \mathbb{E}_{\mathbf{s}|\mathbf{x},\boldsymbol{\eta}_{\text{old}}} [\ln p(\mathbf{x}, \mathbf{s}|\boldsymbol{\eta})]$$

## M-step

ML estimation:

$$\boldsymbol{\eta}^* = \arg \max_{\boldsymbol{\eta}} Q(\boldsymbol{\eta}|\boldsymbol{\eta}_{\text{old}})$$

MAP estimation:

$$\boldsymbol{\eta}^* = \arg \max_{\boldsymbol{\eta}} Q(\boldsymbol{\eta}|\boldsymbol{\eta}_{\text{old}}) + \ln p(\{\mathbf{A}_{e,f}\}) + \ln p(\{\mathbf{A}_{l,f}\})$$

⇒ ML and MAP estimations only differ in the mixing filters update at the M-step.

## Early reverberation prior

### Early reverberation prior

We consider an AR(1) model for the early part of the mixing filters:

$$\ln p(\{\mathbf{A}_{e,f}\}_f) \stackrel{c}{=} -\frac{1}{\sigma_\kappa^2} \sum_{f=1}^{F-1} \left\| \mathbf{A}_{e,f} - \mathbf{\Delta} \circ \mathbf{A}_{e,f-1} \right\|_F^2, \quad (5)$$

where  $\mathbf{\Delta} = [\delta_{ij}]_{ij} \in \mathbb{C}^{I \times J}$ ,  $\|\cdot\|_F^2$  is the Frobenius norm and  $\circ$  is the element-wise matrix product.

### Hyperparameters

- ▶ AR coefficients  $\{\delta_{ij}\}_{ij}$ : Estimated within the M-step.
- ▶ Noise variance  $\sigma_\kappa^2$ : Expresses how confident we are about the prior (assumed to be fixed).

# Late reverberation prior

## Late reverberation prior

We consider an ARMA(7,2) model for the late part of the mixing filters:

$$\ln p(\{\mathbf{A}_f^l\}_f) \stackrel{c}{=} - \sum_{f=0}^{F-1} \text{Trace} \left[ \left( \frac{\Phi(L)}{\Theta(L)} \mathbf{A}_{l,f} \right)^H \boldsymbol{\Sigma}_{\epsilon,f}^{-1} \left( \frac{\Phi(L)}{\Theta(L)} \mathbf{A}_{l,f} \right) \right], \quad (6)$$

where  $\boldsymbol{\Sigma}_{\epsilon,f} = \sigma_{\epsilon}^2 \mathbf{I}_l$ .

## Hyperparameters

- ▶ ARMA coefficients: **Learned and fixed** from the theoretical ACVF, knowing some room parameters.
- ▶ Noise variance  $\sigma_{\epsilon}^2$ : Expresses how confident we are about the prior (assumed to be fixed).



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

Experiments

Ongoing work

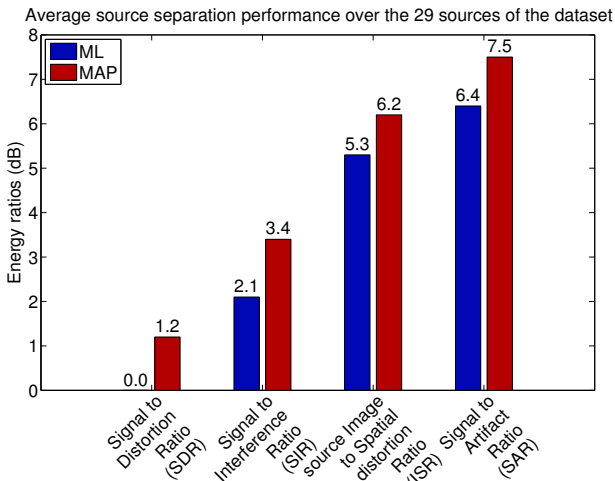
## Conclusion

# Experiments

- ▶ Dataset composed of 8 stereo mixtures:
  - ▶ Created using synthetic room impulse responses;
  - ▶ Reverberation time: 128 ms;
  - ▶ Duration: 12 to 28 seconds;
  - ▶ Number of musical sources: 3 to 5.
- ▶ Source separation results: ML (w/o priors) vs. MAP (w/ priors).
- ▶ Both algorithms are run from the same blind initialization.



# Source separation results















For more details: [Leglaive et al., 2016] "Multichannel audio source separation with probabilistic reverberation priors". *IEEE Transactions on Audio, Speech and Language Processing*.

## Audio example

- ▶ Separation of 4 sources from a stereo mixture.
- ▶ Both algorithms are run from the same blind initialization.

Stereo mixture: 

| source    | drums   |  | guitar 1  |  | guitar 2  |  | voice   |  |
|-----------|---|--|---|--|---|--|---|--|
| original  |        |  |        |  |        |  |        |  |
| estimated | ML<br> | MAP<br> | ML<br> | MAP<br> | ML<br> | MAP<br> | ML<br> | MAP<br> |

Song: "TV On" by Kismet. MTG MASS database.



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

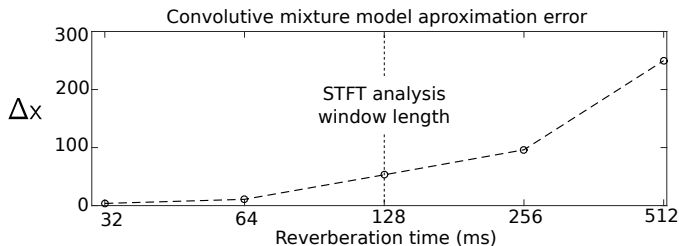
Experiments

Ongoing work

## Conclusion

## Main limitation

Error due to the STFT approximation of the convolutive mixing process.



Average relative squared error: 
$$\Delta x = \frac{1}{IFN} \sum_{i,f,n} \frac{|x_{i,fn} - \hat{x}_{i,fn}|^2}{|x_{i,fn}|^2},$$

▶  $x_{i,fn} = \text{STFT}\{x_i(t)\};$

▶  $\hat{x}_{i,fn} = \sum_{j=1}^J a_{ij,f} s_{j,fn}.$

## Overcoming this limitation

- ▶ More accurate time-frequency convolutional mixture models:
  - ▶ 2D filtering [Badeau and Plumbley, 2014];
  - ▶ Subband filtering (convolutional transfer function) [Li et al., 2017];
- ▶ **Time-domain convolutional mixture model** [Kowalski et al., 2010].

---

[Badeau and Plumbley, 2014] "Multichannel high-resolution NMF for modeling convolutional mixtures of non-stationary signals in the time-frequency domain". *IEEE Transactions on Audio, Speech and Language Processing*.

[Li et al., 2017] "Audio source separation based on convolutional transfer function and frequency-domain Lasso optimization". *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*.

[Kowalski et al., 2010] "Beyond the narrowband approximation: Wideband convex methods for under-determined reverberant audio source separation". *IEEE Transactions on Audio, Speech, and Language Processing*.



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

Experiments

Ongoing work

## Conclusion



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

Experiments

Ongoing work

## Conclusion





## Mixture model

### Time-domain convolutional mixture model

$$x_i(t) = \sum_{j=1}^J [a_{ij} \star s_j](t) + b_i(t), \quad (7)$$

with  $b_i(t) \stackrel{i.i.d}{\sim} \mathcal{N}(0, \sigma_i^2)$ .

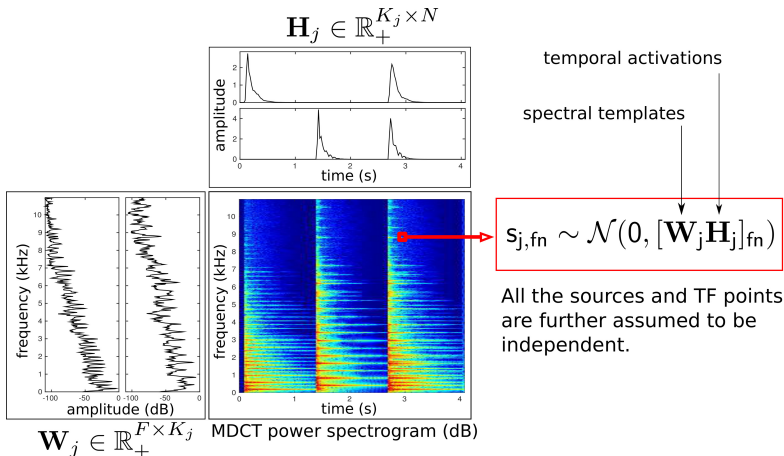
### Time-frequency source representation

$$s_j(t) = \sum_{f=0}^{F-1} \sum_{n=0}^{N-1} s_{j,fn} \psi_{fn}(t), \quad (8)$$

with  $\psi_{fn}(t)$  a Modified Discrete Cosine Transform (MDCT) atom.

# Source model

Gaussian source model based on Non-negative Matrix Factorization.



**Remark:** Source time-frequency coefficients are real-valued.



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

**Inference**

Experiments

Ongoing work

## Conclusion

## Statistical inference

- ▶ Latent **time-frequency** random variables:  $\mathbf{s} = \{s_{j,fn}\}_{j,f,n}$ ;
- ▶ Observed **time-domain** random variables:  $\mathbf{x} = \{x_i(t)\}_{i,t}$ ;
- ▶ Model parameters:  $\eta = \{\{\mathbf{W}_j, \mathbf{H}_j\}_j, \{a_{ij}(t)\}_{i,j,t}, \{\sigma_i^2\}_i\}$ .

### Source and parameter estimation

- ▶ Source estimation according to the posterior mean:

$$\hat{\mathbf{s}} = \mathbb{E}_{\mathbf{s}|\mathbf{x};\eta^*} [\mathbf{s}].$$

- ▶ Maximum likelihood estimation of the parameters:

$$\eta^* = \arg \max_{\eta} p(\mathbf{x}; \eta).$$

The posterior distribution is Gaussian but with a high-dimensional full covariance matrix → **Variational inference**.

## Variational inference

- ▶ We want to find  $q \in \mathcal{F}$  which approximates  $p(\mathbf{s}|\mathbf{x}; \boldsymbol{\eta})$ .
- ▶ Taking the KL divergence as a measure of fit, we can show that:

$$KL(q||p(\mathbf{s}|\mathbf{x}; \boldsymbol{\eta})) = \underbrace{\ln p(\mathbf{x}; \boldsymbol{\eta})}_{\text{Log-likelihood}} - \underbrace{\mathcal{L}(q; \boldsymbol{\eta})}_{\text{Variational Free Energy}}, \quad (9)$$

where  $\mathcal{L}(q; \boldsymbol{\eta}) = \left\langle \ln \left( \frac{p(\mathbf{x}, \mathbf{s}; \boldsymbol{\eta})}{q(\mathbf{s})} \right) \right\rangle_q$  and  $\langle f(\mathbf{s}) \rangle_q = \int f(\mathbf{s})q(\mathbf{s})ds$ .

- ▶ Variational Expectation-Maximization algorithm:
  - ▶ **E-step:**  $q^* = \arg \min_{q \in \mathcal{F}} KL(q||p(\mathbf{s}|\mathbf{x}; \boldsymbol{\eta}_{\text{old}})) = \arg \max_{q \in \mathcal{F}} \mathcal{L}(q; \boldsymbol{\eta}_{\text{old}});$
  - ▶ **M-step:**  $\boldsymbol{\eta}_{\text{new}} = \arg \max_{\boldsymbol{\eta}} \mathcal{L}(q^*; \boldsymbol{\eta}).$



## M-Step

Maximize (or only increase) the variational free energy w.r.t  $\eta$ .

### NMF parameters

Compute an NMF with the Itakura-Saito divergence on:

$$\langle s_{j,fn}^2 \rangle_{q^*} = m_{j,fn}^2 + \gamma_{j,fn},$$

→ standard multiplicative update rules.

### Mixing filters

Solve a Toeplitz system of equations for  $\mathbf{a}_{ij} = [a_{ij}(0), \dots, a_{ij}(L_a - 1)]^T$ .

### Noise variance

$$\sigma_i^2 = \frac{1}{T} \sum_{t=0}^{T-1} \left\langle \left( x_i(t) - \sum_{j=1}^J [a_{ij} \star s_j](t) \right)^2 \right\rangle_{q^*}.$$



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

Experiments

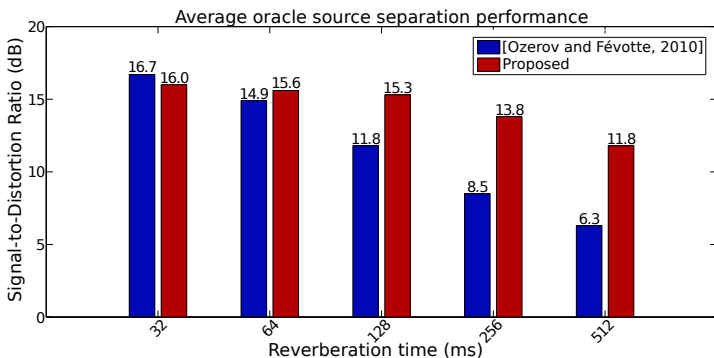
Ongoing work

## Conclusion



## Oracle experiment

- ▶ **Dataset:** Same as before with different reverberation times;
- ▶ **Oracle initialization** of the parameters;
- ▶ STFT and MDCT analysis/synthesis window length: 128 ms.



For more details: [Leglaive et al., 2017] "Multichannel audio source separation: variational inference of time-frequency sources from time-domain observations". *IEEE International Conference on Acoustics Speech and Signal Processing (ICASSP)*.

## Semi-blind experiment

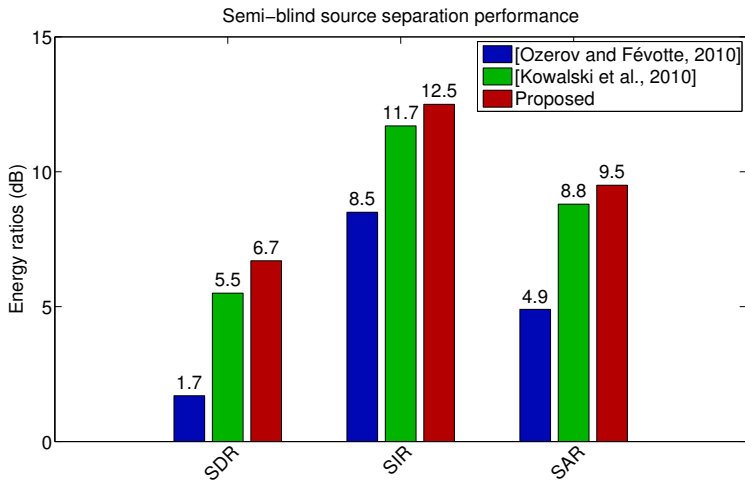
- ▶ Mixing filters are known and fixed;
- ▶ All other parameters are blindly estimated;
- ▶ Compared methods:

|                   |                     | Convolutional mixture model |                            |
|-------------------|---------------------|-----------------------------|----------------------------|
|                   |                     | Exact (time)                | Approximate (TF)           |
| Source model (TF) | Sparse ( $\ell_1$ ) | [Kowalski et al., 2010]     | -                          |
|                   | Gaussian NMF-based  | Proposed                    | [Ozerov and Févotte, 2010] |

- ▶ Dataset: Same as before with a reverberation time of 256 ms.
















[Kowalski et al., 2010] "Beyond the narrowband approximation: Wideband convex methods for under-determined reverberant audio source separation". *IEEE Transactions on Audio, Speech, and Language Processing*.

## Semi-blind experiment results



## Semi-blind audio example

Stereo mixture: 

|          | Original  | [Ozerov and Févotte, 2010] | [Kowalski et al., 2010]   | Proposed  |
|----------|---|----------------------------|---|---|
| Drums    |  |                            |  |  |
| Guitar 1 |  |                            |  |  |
| Guitar 2 |  |                            |  |  |
| Voice    |  |                            |  |  |
| Bass     |  |                            |  |  |

Musical excerpt from "Ana" by Vieux Farka Toure. MTG MASS database.



# Outline

## Convolutional mixture model in the STFT domain

Baseline source separation framework

Room frequency response modeling

Source separation with reverberation priors

Experiments

Limitations

## Convolutional mixture model in the time domain

Model

Inference

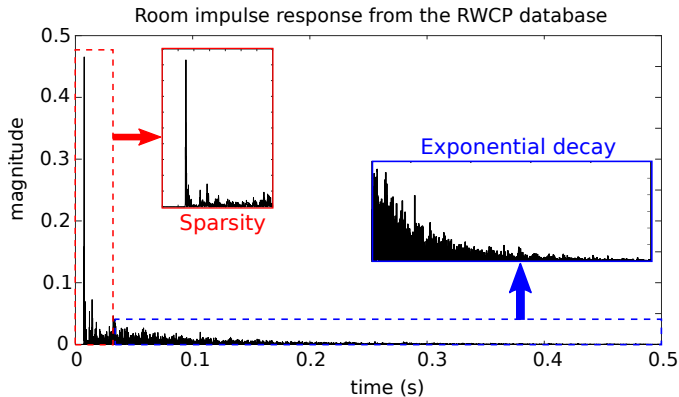
Experiments

Ongoing work

## Conclusion

## Ongoing work

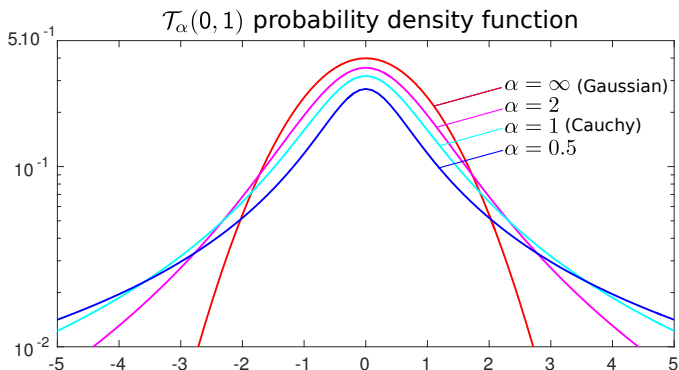
Probabilistic priors on the mixing filters **in the time domain**.



# Student's $t$ distribution

Student's  $t$  distribution:  $\mathcal{T}_\alpha(\mu, \sigma)$

- ▶ Shape:  $\alpha$ ;
- ▶ Location:  $\mu$ ;
- ▶ Scale:  $\sigma$ .









# Thank you

More audio examples and Matlab code available at:

<https://perso.telecom-paristech.fr/leglaive/>