

# Modeling Reverberant Mixtures for Multichannel Audio Source Separation

Simon Leglaive

Ph.D. supervisors: Roland Badeau and Gaël Richard

LTCI, Télécom ParisTech, Université Paris Saclay

Ph.D. defense, Télécom ParisTech, Paris

December 12, 2017

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

#### Introduction

Audio source separation Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

### Time-domain mixture model

Model Inference

### Conclusion

December 12, 2017 Modeling Reverberant Mixtures for Multichannel Audio Source Sepa	aration
------------------------------------------------------------------------------------	---------

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

#### Introduction

### Audio source separation

Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

### Time-domain mixture model

- Model Inference
- Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Source separation

Objective: Recover source signals from one or several mixtures.



Many applications:

- Biomedical signal processing (ECG, EEG, MEG, MRI, etc.);
- Astrophysics;
- Underwater acoustics;
- Audio signal processing;
- etc.

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Audio source separation in everyday life



2/55

3/55

STFT-domain mixture model

Time-domain mixture model

Conclusion

### Audio source separation for karaoke



Introduction STET-domain mixture model Time-domain mixture model 

Conclusion

## Audio source separation for music upmixing



STFT-domain mixture model

Time-domain mixture model

Conclusion



#### Under-determined and reverberant multichannel mixture.



5/55

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

### Introduction

Audio source separation

### Source and mixture modeling

Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

### Time-domain mixture model

- Model Inference
- Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Model-based approach

A model aims to explain how are the observed data generated.



STFT-domain mixture model

Time-domain mixture model

Conclusion

## **Time-frequency source representation**

Time-frequency (TF) transforms provide meaningful representations.



Spectrograms computed from the short-term Fourier transform (STFT).

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Reverberant mixtures (1)

Convolutive mixing process in the time domain:  $x_i(t) = \sum_{j=1}^{J} [a_{ij} \star s_j](t)$ 



STFT-domain mixture model

Time-domain mixture model

Conclusion

## **Reverberant mixtures (2)**

Convolutive mixing process in the STFT domain:  $x_{i,fn} \approx \sum_{j=1}^{J} a_{ij,f} s_{j,fn}$ 



STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

### Introduction

Audio source separation Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

### Time-domain mixture model

- Model Inference
- Lxpenner

### Conclusion

December 12, 2017	Modeling Reverberant Mixtures for Multichannel Audio Source Separation
•	

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Probabilistic model with latent variables (1)

- Observed random variables: x
- Latent random variables: z

(e.g. mixture coefficients) (e.g. source coefficients)

### Defining the probabilistic model

$$p(\mathbf{x}, \mathbf{z}; \boldsymbol{\theta}) = p(\mathbf{x}|\mathbf{z}; \boldsymbol{\theta})p(\mathbf{z}; \boldsymbol{\theta}),$$

where  $\theta$  is the set of model parameters.

- What prior knowledge do we have on the latent variables?
- How are the data generated from the latent variables?

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Probabilistic model with latent variables (2)

- Observed random variables: x
- Latent random variables: z

(e.g. mixture coefficients) (e.g. source coefficients)

### Defining the probabilistic model

 $p(\mathbf{x}, \mathbf{z}, \boldsymbol{\theta}) = p(\mathbf{x}|\mathbf{z}, \boldsymbol{\theta})p(\mathbf{z}|\boldsymbol{\theta})p(\boldsymbol{\theta}),$ 

where  $\theta$  is the set of model parameters.

- What prior knowledge do we have on the latent variables?
- How are the data generated from the latent variables?
- What prior knowledge do we have on the model parameters?

STFT-domain mixture model

Time-domain mixture model

Conclusion

## **Statistical inference**

### Posterior inference

We are interested in the posterior distribution of the latent variables:

 $p(\mathbf{z}|\mathbf{x}; \boldsymbol{\theta}^{\star}),$ 

with  $\theta^{\star}$  an estimate of the model parameters.

### Parameters estimation

Maximum likelihood (ML) - parameters are treated as deterministic variables:

$$\boldsymbol{ heta}^{\star}_{\mathsf{ML}} = rg\max_{\boldsymbol{ heta}} p(\mathbf{x}; \boldsymbol{ heta}).$$

Maximum a posteriori (MAP) - parameters are treated as random variables:

$$oldsymbol{ heta}_{\mathsf{MAP}}^{\star} = rg\max_{oldsymbol{ heta}} p(oldsymbol{ heta} | \mathbf{x}) = rg\max_{oldsymbol{ heta}} p(\mathbf{x} | oldsymbol{ heta}) p(oldsymbol{ heta}).$$

STFT-domain mixture model

Time-domain mixture model

Conclusion

## **Research problem**

- Mixing filters are usually treated as deterministic parameters only estimated from the observed data.
- We know that they correspond to room responses.



How can we guide the estimation of the mixing filters?

Two approaches:

- 1. STFT-domain convolutive mixture model;
- 2. Time-domain convolutive mixture model.

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

### Time-domain mixture model

Model Inference Experiment

### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

### Baseline source separation framework

Room frequency response modeling Source separation with reverberation priors Limitations

### Time-domain mixture model

- Model Inference
- Conclusion

S<sub>j,fn</sub>

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Mixture model

### Convolutive noisy mixture of J sources on I channels

STFT domain:  

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

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

- Source STFT coefficients: s<sub>j,fn</sub>
- Additive Gaussian noise:  $b_{i,fn} \sim \mathcal{N}_c(0, \sigma_{b,f}^2)$

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Mixture model

### Convolutive noisy mixture of J sources on I channels

STFT domain:  

$$\forall (f, n) \in \{0, ..., F - 1\} \times \{0, ..., N - 1\}$$
 $x_{i,fn} = \sum_{j=1}^{J} a_{ij,f} s_{j,fn} + b_{i,fn}$ 

- ► Source STFT coefficients: *s<sub>j,fn</sub>*
- Additive Gaussian noise:  $b_{i,fn} \sim \mathcal{N}_c(0, \sigma_{b,f}^2)$

### In matrix form

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

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  $\mathbf{b}_{fn} = [b_{i,fn}]_i \in \mathbb{C}^I$ .

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Source model

Gaussian source model based on non-negative matrix factorization (NMF) [1]:

$$s_{j,fn} \sim \mathcal{N}_c \Big( 0, (\mathbf{W}_j \mathbf{H}_j)_{f,n} \Big).$$

All sources and TF points are further assumed to be independent.



[1] C. Févotte, N. Bertin and J. L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis", *Neural computation*, 2009.

Time-domain mixture model

Conclusion

## Statistical inference (1)

- ► Latent source random variables: s = {s<sub>j,fn</sub>}<sub>j,f,n</sub>;
- Observed random variables:  $\mathbf{x} = \{x_{i,fn}\}_{i,f,n}$ ;

• Model parameters: 
$$\boldsymbol{\theta} = \left\{ \left\{ \mathbf{W}_{j}, \mathbf{H}_{j} \right\}_{j}, \left\{ \mathbf{A}_{f}, \sigma_{b, f}^{2} \right\}_{f} \right\}.$$

#### Source posterior mean

$$\hat{\mathbf{s}} = \mathbb{E}_{\mathbf{s}|\mathbf{x}; \boldsymbol{ heta}_{\mathsf{ML}}^{\star}}[\mathbf{s}]$$

Maximum likelihood estimation of the parameters

$$oldsymbol{ heta}_{\mathsf{ML}}^{\star} = rg\max_{oldsymbol{ heta}} p(\mathbf{x};oldsymbol{ heta})$$

### $\rightarrow$ Expectation-maximization (EM) algorithm [2].

[2] A. Ozerov and C. Févotte, "Multichannel nonnegative matrix factorization in convolutive mixtures for audio source separation", *IEEE Transactions on Audio, Speech and Language Processing*, 2010.

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Statistical inference (2)

- ► Latent source random variables: s = {s<sub>j,fn</sub>}<sub>j,f,n</sub>;
- Observed random variables:  $\mathbf{x} = \{x_{i,fn}\}_{i,f,n}$ ;
- Model parameters:  $\boldsymbol{\theta} = \left\{ \left\{ \mathbf{W}_{j}, \mathbf{H}_{j} \right\}_{j}, \left\{ \mathbf{A}_{f}, \sigma_{b, f}^{2} \right\}_{f} \right\}.$

#### Source posterior mean

$$\hat{\mathbf{s}} = \mathbb{E}_{\mathbf{s} | \mathbf{x}; \boldsymbol{\theta}_{\mathsf{MAP}}^{\star}}[\mathbf{s}]$$

Maximum a posteriori estimation of the parameters

$$m{ heta}^{\star}_{\mathsf{MAP}} = rg\max_{m{ heta}} p(\mathbf{x}|m{ heta}) p\left(\left\{\mathbf{A}_{f}
ight\}_{f}
ight)$$

 $\rightarrow$  Expectation-maximization (EM) algorithm.

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework

### Room frequency response modeling

Source separation with reverberation priors Limitations

### Time-domain mixture model

Model Inference

### Conclusion

Time-domain mixture model

Conclusion

## Room impulse response

Mixing filters are room responses. They exhibit a simple specific structure in the time domain.



STFT-domain mixture model

Time-domain mixture model

Conclusion

## Room impulse and frequency responses

Room impulse and frequency responses

 $\forall t, f \in \{0, ..., T - 1\}$ :

$$\underbrace{a(t) = a_e(t) + a_I(t)}_{\text{Room impulse response (RIR)}} \stackrel{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\text{Room frequency response (RFR)}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\overset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^{-1}}}{\underset{\mathcal{F}_T^{-1}}{\underset{\mathcal{F}_T^$$

 $\mathcal{F}_T$ : Discrete Fourier transform (DFT) over T points.



Time-domain mixture model

Conclusion

## Early contributions model (1)



### Geometrical room acoustics

The *k*-th early contribution is characterized by an amplitude  $\rho_k$  and a delay  $\tau_k$ :

$$A_e(f) = \sum_{k=0}^{R-1} \rho_k \delta_k^f \quad \text{with} \quad \delta_k = e^{-i2\pi\tau_k/T}.$$
 (2)

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Early contributions model (2)

It follows (see, e.g., [3])

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

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

#### Adding an error term $\rightarrow$ autoregressive model

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

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

Time-domain mixture model

Conclusion

### Autoregressive model of order 1

Assuming that the direct path dominates the early echoes:

 $|A_e(f)| \approx |A_e(f-1)|$  and  $\arg(A_e(f)) \approx \arg(A_e(f-1)) - 2\pi \frac{\tau_0}{T}$ .



## Late reverberation model (time domain)

Conclusion



#### Statistical room acoustics

Exponentially decaying energy temporal profile (ETP):

$$ar{a}_l(t) = \mathbb{E}\left[a_l^2(t)
ight] \propto e^{-2t/ au}\mathbbm{1}_{t\geq t_0}(t).$$

- $\tau$  is related to the reverberation time  $T_{60}$  in seconds.
- ▶ E[·]: spatial averaging.

## Late reverberation model (frequency domain)

**Statistical room acoustics**:  $\{A_l(f)\}_f$  is a proper, centered and wide-sense stationary complex Gaussian random process.

### Theoretical power spectral density (PSD)

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

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

### 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).

24/55 De

December 12, 2017

Modeling Reverberant Mixtures for Multichannel Audio Source Separation

STFT-domain mixture model

Time-domain mixture model

Conclusion

## **Experimental validation**

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



Figure: Theoretical and empirical autocovariance functions

STFT-domain mixture model

Time-domain mixture model

Conclusion

## **ARMA** parametrization

ARMA representation of late reverberation in the frequency domain

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

$$\Phi(L)A_{I}(f) = \Theta(L)\epsilon(f),$$

• 
$$\Phi(L) = \sum_{p=0}^{P} \varphi_p^l L^p$$
 and  $\Theta(L) = \sum_{q=0}^{Q} \vartheta_q L^q$  with  $\varphi_0^l = \vartheta_0 = 1$ ;

• L is the lag operator, i.e.  $LA_{I}(f) = A_{I}(f-1)$ ;

•  $\epsilon(f) \stackrel{i.i.d}{\sim} \mathcal{N}_c(0, \sigma_{\epsilon}^2).$ 

We can compute the ARMA parameters from the theoretical ACVF.

Time-domain mixture model

Conclusion

## **Experimental validation**

#### Same room parameters as used before for simulated RIRs.





Figure: ARMA(7,2) parametrization

#### Figure: Synthesized late RIR
STFT-domain mixture model

Time-domain mixture model

Conclusion

### Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

#### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling

#### Source separation with reverberation priors

Limitations

#### Time-domain mixture model

Model Inference

Experiments

#### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

# **EM** algorithm

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

E-step

$$Q(oldsymbol{ heta};oldsymbol{ heta}_{\mathsf{old}}) = \mathbb{E}_{\mathbf{s}|\mathbf{x},oldsymbol{ heta}_{\mathsf{old}}}ig[\ln p(\mathbf{x},\mathbf{s}|oldsymbol{ heta})ig]$$

#### M-step

$$\frac{\text{ML estimation:}}{\theta^{\star} = \arg \max_{\theta} Q(\theta; \theta_{\text{old}})} \begin{vmatrix} \frac{\text{MAP estimation:}}{\theta} \\ \theta^{\star} = \arg \max_{\theta} Q(\theta; \theta_{\text{old}}) \\ + \ln p(\{\mathbf{A}_{e,f}\}_{f}) + \ln p(\{\mathbf{A}_{l,f}\}_{f}) \end{vmatrix}$$

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

28/55	December 12, 2017	Modeling Reverberant Mixtures for Multichannel Audio Source Separation
-------	-------------------	------------------------------------------------------------------------

STFT-domain mixture model  Time-domain mixture model

Conclusion

### **Experiments**

#### Dataset:

- 8 stereo mixtures created with synthetic room impulse responses;
- Reverberation time: 128 ms:
- Number of sources per mixture: 3 to 5; ►
- Musical sources: drums, piano, bass, guitar, voice;
- Duration: 12 to 28 seconds.
- Source separation results: ML (w/o priors) vs. MAP (w/ priors).
- Both algorithms are run from the same initialization.

STFT-domain mixture model

Time-domain mixture model

Conclusion

### Source separation results



**Remark**: Some hyperparameters of the priors (noise variances) are optimized knowing the true source signals.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Outline

#### Introduction

Audio source separation Source and mixture modeling Probabilistic modeling

#### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

#### Time-domain mixture model

Model Inference

#### Conclusion

 Time-domain mixture model

Conclusion

### **Main limitation**

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



# **Overcoming this limitation**

- More accurate time-frequency convolutive mixture representations:
  - 2-dimensional filtering [4];
  - Subband filtering (convolutive transfer function) [5, 6].

#### ► Time-domain convolutive mixture representation [7]:

- Exact representation of the convolutive mixing process;
- Suitable for incorporating simple priors on the mixing filters.

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

[5] X. Li, L. Girin and R. Horaud, "Audio source separation based on convolutive transfer function and frequency-domain Lasso optimization", *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2017.

[6] X. Li, L. Girin and R. Horaud, "An EM algorithm for audio source separation based on the convolutive transfer function", IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA), 2017.

[7] M. Kowalski, E. Vincent and R. Gribonval, "Beyond the narrowband approximation: Wideband convex methods for under-determined reverberant audio source separation". IEEE Transactions on Audio, Speech, and Language Processing, 2010.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

#### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

#### Time-domain mixture model

Model Inference Experiments

#### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

#### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

#### Time-domain mixture model

Model

Inference

Experiments

#### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

### **Proposed approach**

Time-domain mixture representation and time-frequency source representation.



STFT-domain mixture model

Time-domain mixture model

Conclusion

(5)

(6)

# Mixture model

#### Time-domain convolutive mixture model

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

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

Time-frequency source representation

$$egin{aligned} & s_j(t) = \sum_{(f,n)\in\mathcal{B}_j} s_{j,fn} \psi_{j,fn}(t), \end{aligned}$$

with  $\psi_{j,fn}(t) \in \mathbb{R}$  a source-dependent modified discrete cosine transform (MDCT) atom and  $\mathcal{B}_j = \{0, ..., F_j - 1\} \times \{0, ..., N_j - 1\}$ .

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

STFT-domain mixture model

Time-domain mixture model

Conclusion

### Student's t distribution

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

- Shape:  $\alpha$ ;
- Location: μ;
- Scale: σ.



#### Scale mixture of Gaussians

$$z \sim \mathcal{T}_{\alpha}(\mu, \sigma) \quad \Leftrightarrow \begin{cases} z | v \ \sim \mathcal{N}\left(\mu, v \sigma^{2}\right) \\ v \ \sim \mathcal{IG}\left(rac{lpha}{2}, rac{lpha}{2}
ight) \end{cases}$$

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Source model

#### Student's t source model based on NMF [8]:

$$s_{j,fn} \sim \mathcal{T}_{\alpha_v} \Big( \mathbf{0}, \lambda_{j,fn} \Big),$$

with 
$$\lambda_{j,fn}^2 = (\mathbf{W}_j \mathbf{H}_j)_{f,n}$$
.



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

[8] K. Yoshii, K. Itoyama and M. Goto, "Student's t nonnegative matrix factorization and positive semidefinite tensor factorization for single-channel audio source separation", *IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, 2016.

36/55	December 12, 2017	Modeling Reverberant Mixtures for Multichannel Audio Source Separation
-------	-------------------	------------------------------------------------------------------------

STFT-domain mixture model

Time-domain mixture model

Conclusion

### Gaussian RIR model (1)



STFT-domain mixture model

Time-domain mixture model

Conclusion

# Gaussian RIR model (1)



#### Gaussian model with exponential decay [9]

Independently for all  $t \in \{0, ..., L_a - 1\}$ :

$$a(t) \sim \exp(-t/\tau) \mathcal{N}(0, \sigma_r^2).$$

Theoretically valid only for late reverberation (diffuse sound field).

[9] J. D. Polack, "La transmission de l'énergie sonore dans les salles", Ph.D. dissertation, Université du Maine, 1988.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Gaussian RIR model (2)



#### Gaussian model with exponential decay [9]

Equivalently:

$$a(t)/\exp(-t/ au) \stackrel{i.i.d}{\sim} \mathcal{N}(0,\sigma_r^2).$$

Theoretically valid only for late reverberation (diffuse sound field).

[9] J. D. Polack, "La transmission de l'énergie sonore dans les salles", Ph.D. dissertation, Université du Maine, 1988.

STFT-domain mixture model

Time-domain mixture model

Conclusion

## Student's t RIR model



[10] E. Hadad, F. Heese, P. Vary, and S. Gannot, "Multichannel audio database in various acoustic environments", IEEE International Workshop on Acoustic Signal Enhancement (IWAENC), 2014.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Student's t RIR model



[10] E. Hadad, F. Heese, P. Vary, and S. Gannot, "Multichannel audio database in various acoustic environments", IEEE International Workshop on Acoustic Signal Enhancement (IWAENC), 2014.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Student's t RIR model



#### Student's t model with exponential decay

Independently for all i, j and t:

$$a_{ij}(t)\sim \mathcal{T}_{lpha_u}(0,r(t)), \qquad r^2(t)=\sigma_r^2\exp(-2t/ au).$$

[10] E. Hadad, F. Heese, P. Vary, and S. Gannot, "Multichannel audio database in various acoustic environments", IEEE International Workshop on Acoustic Signal Enhancement (IWAENC), 2014.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# **Bayesian network**



- z: set of all latent variables (empty circles);
- x: set of observations (shaded circles);
- $\theta$ : set of model parameters to be estimated (dots).

**Remark**: We recall that  $\lambda_{j,fn}^2 = (\mathbf{W}_j \mathbf{H}_j)_{f,n}$ .

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

#### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

#### Time-domain mixture model

Model Inference

Experiments

#### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion



Exact posterior inference is analytically intractable.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Variational inference

• Exact posterior inference is analytically intractable.

▶ Variational inference: Find  $q \in \mathcal{F}$  which approximates  $p(\mathbf{z}|\mathbf{x}; \theta)$ .

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Variational inference

- Exact posterior inference is analytically intractable.
- ▶ Variational inference: Find  $q \in \mathcal{F}$  which approximates  $p(\mathbf{z}|\mathbf{x}; \theta)$ .
- ▶ We take the Kullback-Leibler divergence as a measure of fit:

$$D_{\mathcal{K}L}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x};\boldsymbol{\theta})) = \underbrace{\ln p(\mathbf{x};\boldsymbol{\theta})}_{\text{log-(marginal) likelihood}} - \underbrace{\mathcal{L}(q;\boldsymbol{\theta})}_{\text{variational free energy}}, \quad (7)$$
  
where  $\mathcal{L}(q;\boldsymbol{\theta}) = \left\langle \ln \left( \frac{p(\mathbf{x},\mathbf{z};\boldsymbol{\theta})}{q(\mathbf{z})} \right) \right\rangle_{q}$  and  $\langle f(\mathbf{z}) \rangle_{q} = \int f(\mathbf{z})q(\mathbf{z})d\mathbf{z}.$ 

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Variational inference

- Exact posterior inference is analytically intractable.
- ▶ Variational inference: Find  $q \in \mathcal{F}$  which approximates  $p(\mathbf{z}|\mathbf{x}; \theta)$ .
- We take the Kullback-Leibler divergence as a measure of fit:

$$D_{KL}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x};\theta)) = \underbrace{\ln p(\mathbf{x};\theta)}_{\text{log-(marginal) likelihood}} - \underbrace{\mathcal{L}(q;\theta)}_{\text{variational free energy}}, \quad (7)$$
  
where  $\mathcal{L}(q;\theta) = \left\langle \ln \left( \frac{p(\mathbf{x},\mathbf{z};\theta)}{q(\mathbf{z})} \right) \right\rangle_q$  and  $\langle f(\mathbf{z}) \rangle_q = \int f(\mathbf{z})q(\mathbf{z})d\mathbf{z}.$ 

Variational EM algorithm:

► **E-step**: 
$$q^* = \underset{q \in \mathcal{F}}{\arg\min} D_{\mathcal{KL}}(q(\mathbf{z})||p(\mathbf{z}|\mathbf{x}; \theta)) = \underset{q \in \mathcal{F}}{\arg\max} \mathcal{L}(q; \theta_{\text{old}});$$

• **M-step**: 
$$\theta_{\text{new}} = \arg \max_{\theta} \mathcal{L}(q^*; \theta).$$

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Mean-field approximation

 $\ensuremath{\mathcal{F}}$  is the set of probability density functions that factorize as:

$$q(\mathbf{z}) = \prod_{z_k \in \mathbf{z}} q_k(z_k).$$



STFT-domain mixture model

Time-domain mixture model

Conclusion

# Mean-field approximation

 ${\cal F}$  is the set of probability density functions that factorize as:

$$q(\mathbf{z}) = \prod_{z_k \in \mathbf{z}} q_k(z_k).$$

- true posterior - mean field approximation

Under the mean-field approximation we can show that:

$$\begin{aligned} q_{jfn}^{s}(s_{j,fn})^{*} &= N(\hat{s}_{j,fn}, \gamma_{j,fn}); \\ q_{ijt}^{a}(a_{ij}(t))^{*} &= N\left(\hat{a}_{ij}(t), \rho_{ij}(t)\right); \\ q_{jfn}^{v}(v_{j,fn})^{*} &= IG(\nu_{v}, \beta_{j,fn}); \\ q_{ijt}^{u}(u_{ij}(t))^{*} &= IG(\nu_{u}, d_{ij}(t)). \end{aligned}$$

 $\textbf{E-Step} \rightarrow \textbf{update}$  all the variational parameters.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# M-Step

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

NMF parameters

$$\min_{\mathbf{W}_{j},\mathbf{H}_{j}\geq 0}\sum_{(f,n)\in\mathcal{B}_{j}}d_{IS}\left(\left\langle v_{j,fn}^{-1}\right\rangle _{q^{\star}}\left\langle s_{j,fn}^{2}\right\rangle _{q^{\star}},\,(\mathbf{W}_{j}\mathbf{H}_{j})_{f,n}\right),$$

where  $d_{IS}(\cdot, \cdot)$  is the Itakura-Saito divergence.

 $\rightarrow$  multiplicative update rules [1].

Noise variance

 $\sigma_i^2$  is manually decreased along the iterations.



STFT-domain mixture model

Time-domain mixture model

Conclusion

# Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

#### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

#### Time-domain mixture model

Model Inference

Experiments

#### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion

# **Experimental setup**

#### Dataset:

- 8 stereo mixtures created with RIRs from the MIRD database [10];
- Reverberation times: 160, 360 and 610 ms;
- Number of sources per mixture: 3 to 5;
- Musical sources: drums, piano, bass, guitar, voice;
- Duration: 12 to 28 seconds.
- Semi-blind scenario:
  - ▶ NMF dictionaries **W**<sub>j</sub> are pre-trained using the true source signals;
  - Reverberation time is assumed to be known;
  - All other parameters are blindly estimated.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Model hyperparameters

How do we choose  $\alpha_v$ ,  $\alpha_u$  and the MDCT window length?

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Model hyperparameters

How do we choose  $\alpha_v$ ,  $\alpha_u$  and the MDCT window length?

Average SDR in dB using the mixtures with a reverberation time of 360 ms



# **Remark**: MDCT window length is fixed to 64 ms.

STFT-domain mixture model

Time-domain mixture model

Conclusion

# Model hyperparameters

How do we choose  $\alpha_v$ ,  $\alpha_u$  and the MDCT window length?

Average SDR in dB using the mixtures with a reverberation time of 360 ms



# **Remark**: MDCT window length is fixed to 64 ms.

**Remark**:  $(\alpha_v, \alpha_u)$  is fixed to (100, 1).

STFT-domain mixture model

Time-domain mixture model

Conclusion

# **Baseline methods**

#### First method: "Gaussian - SCM rank 1-2" [11]

- Source model: STFT Gaussian NMF.
- Convolutive mixture model:
  - Approximate STFT-domain convolutive mixture representation.
  - Spatial covariance matrix (SCM) model:
    - Rank 1: punctual source;
    - Rank 2: diffuse source.

Second method: "Gaussian - unconst. TD filters" [12]

- Source model: MDCT Gaussian NMF.
- Convolutive mixture model:
  - Exact time-domain (TD) convolutive mixture representation.
  - Unconstrained mixing filters (w/o prior).

[11] A. Ozerov, E. Vincent and F. Bimbot, "A general flexible framework for the handling of prior information in audio source separation", *IEEE Transactions on Audio, Speech, and Language Processing*, 2012.

[12] S. L., R. Badeau and G. Richard, "Multichannel audio source separation: Variational inference of time-frequency sources from time-domain observations", IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP), 2017.

STFT-domain mixture model

Time-domain mixture model

Conclusion

### Source separation results



STFT-domain mixture model

Time-domain mixture model

Conclusion

### Source separation results


STFT-domain mixture model

Time-domain mixture model

Conclusion

### Source separation results



STFT-domain mixture model

Time-domain mixture model

Conclusion

### Source separation results



STFT-domain mixture model

Time-domain mixture model

Conclusion

### Mixing filters: Influence of the prior

- Convolutive mixture equation:  $x_i(t) = \sum_{i=1}^{J} [a_{ij} \star s_j](t)$ .
- Multiple solutions can explain the same observed data.





STFT-domain mixture model

Time-domain mixture model

Conclusion

### Audio example (1)

- Separation of 3 sources from a stereo mixture ( $T_{60} = 610$  ms).
- ► All algorithms are run using oracle NMF dictionaries.

Stereo mixture: 🧕

source image	drums	guitar	bass
original	0	0	0
Gaussian - SCM rank-1	0	0	٥
Gaussian - unconstrained TD filters	0	0	٥
Prop. w/o adapted TF window	0	0	٥

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

50/55	December 12, 2017	Modeling Reverberant Mixtures for Multichannel Audio Source Separation
-------	-------------------	------------------------------------------------------------------------

STFT-domain mixture model

Time-domain mixture model

Conclusion

### Audio example (2)

- Stereo mixture provided by Radio France (Edison 3D ANR project).
- Blind separation of voice and instrumental.



0

Song: "C'est magnifique" by Ella Fitzgerald (Nice Jazz Festival 1972 - Recording: ORTF).

STFT-domain mixture model

Time-domain mixture model

Conclusion

0000

### Outline

#### Introductior

Audio source separation Source and mixture modeling Probabilistic modeling

### STFT-domain mixture model

Baseline source separation framework Room frequency response modeling Source separation with reverberation priors Limitations

#### Time-domain mixture model

Model Inference Experiment

### Conclusion

STFT-domain mixture model

Time-domain mixture model

Conclusion ●○○○

### Conclusion

Probabilistic modeling of reverberation for audio source separation

Guiding the estimation of the mixing filters can help improving source separation results.

- Convolutive mixture model in the STFT domain:
  - ► Time-domain dynamics of the mixing filters → frequency-domain correlations;
  - Limited to weakly reverberant mixtures.
- Convolutive mixture model in the time domain:
  - Appropriate for highly reverberant mixtures;
  - Simple priors on the mixing filters, in the time domain;
  - Multi-resolution source modeling.

# Perspectives (1)

#### Supervised source models:

- Learn NMF dictionaries on an external dataset;
- Neural network model for the source scale parameters [13];
- Variational autoencoder for learning the source prior distributions [14].

### Speech enhancement and recognition:

- Robustness to interfering noise and reverberant conditions is still an open issue for automatic speech recognition (ASR);
- ► In ASR, features are computed from the enhanced signal → uncertainty propagation from the speech posterior distribution [15].

[13] A. Nugraha, A. Liutkus and E. Vincent, "Multichannel music separation with deep neural networks", European Signal Processing Conference (EUSIPCO), 2016.

[14] Y. Bando, M. Mimura, K. Itoyama, K. Yoshii and T. Kawahara, "Statistical speech enhancement based on probabilistic integration of variational autoencoder and non-negative matrix factorization", *arXiv:1710.11439 [cs.SD]*, 2017.

[15] K. Adiloğlu and E. Vincent, "Variational Bayesian inference for source separation and robust feature extraction", IEEE Transactions on Audio, Speech, and Language Processing, 2016.

STFT-domain mixture model

Time-domain mixture model

Conclusion ○○●○

## Perspectives (2)

#### Deconvolution:

- Can we recover s(t) from  $x(t) = [a \star s](t)$ ?
- Dictionary learning problem [16]:

$$x(t) = \sum_{f,n} s_{fn} \left[ a \star \psi_{fn} \right](t).$$

#### Source localization:

- The latent variables u<sub>ij</sub>(t) are supposed to "encode" the early contributions;
- Can we use features computed from their posterior distribution to perform source localization?
- Could we extend the hierarchical model so that they depend on the source location?

<sup>[16]</sup> D. Barchiesi and M. Plumbley, "Dictionary learning of convolved signals", IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), 2011.

### **Publications**

#### International peer-reviewed journals:

- S. L., R. Badeau, G. Richard, "Student's t source and mixing models for multichannel audio source separation", *IEEE Transactions on Audio, Speech and Language Processing*, 2017 (submitted).
- S. L., R. Badeau, G. Richard, "Multichannel audio source separation with probabilistic reverberation priors", *IEEE Transactions on Audio, Speech and Language Processing*, 2016.

#### International peer-reviewed conferences and workshops:

- S. L., R. Badeau, G. Richard, "Separating time-frequency sources from time-domain convolutive mixtures using non-negative matrix factorization", *IEEE WASPAA*, 2017.
- S. L., R. Badeau, G. Richard, "Semi-blind Student's t source separation for multichannel audio convolutive mixtures", EUSIPCO, 2017.
- S. L., U. Şimşekli, A. Liutkus, R. Badeau, G. Richard, "Alpha-stable multichannel audio source separation", IEEE ICASSP, 2017.
- S. L., R. Badeau, G. Richard, "Multichannel audio source separation: Variational inference of time-frequency sources from time-domain observations", *IEEE ICASSP*, 2017.
- S. L., R. Badeau, G. Richard, "Autoregressive moving average modeling of late reverberation in the frequency domain", EUSIPCO, 2016.
- S. L., R. Badeau, G. Richard, "Multichannel audio source separation with probabilistic reverberation modeling", *IEEE WASPAA*, 2015.

# Thank you

Audio examples and Matlab code available at:

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