

Online spectrogram inversion for low-latency audio source separation

Paul Magron, Tuomas Virtanen

CNRS, IRIT, Université de Toulouse, France

ICASSP2021



Introduction

Multiple Input Spectrogram Inversion

Experiments

Introduction

Audio source separation

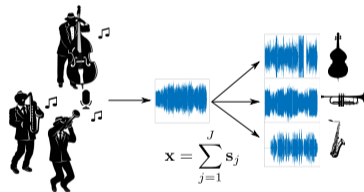
- . Audio signals are composed of several constitutive sounds: multiple speakers, background noise, domestic sounds, musical instruments...

Audio source separation

- . Audio signals are composed of several constitutive sounds: multiple speakers, background noise, domestic sounds, musical instruments...

Source separation = recovering the sources from the mixture.

- . Automatic speech recognition (clean speech vs. noise).
- . Rhythm analysis (drums vs. harmonic instruments).
- . Time-stretching (transients vs. partials).

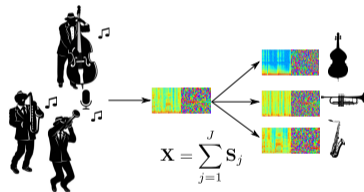


Audio source separation

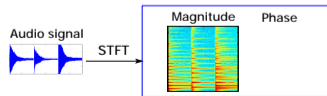
- . Audio signals are composed of several constitutive sounds: multiple speakers, background noise, domestic sounds, musical instruments...

Source separation = recovering the sources from the mixture.

- . Automatic speech recognition (clean speech vs. noise).
- . Rhythm analysis (drums vs. harmonic instruments).
- . Time-stretching (transients vs. partials).



Time-frequency separation = acts on the short-time Fourier transform (STFT).



1. Nonnegative representation, e.g., $\mathbf{V} = |j\text{STFT}(\mathbf{x})|^2$.

1. Nonnegative representation, e.g., $\mathbf{V} = |\text{STFT}(\mathbf{x})|^2$.
2. Structured model, e.g., nonnegative matrix factorization, deep neural networks.

1. Nonnegative representation, e.g., $\mathbf{V} = |\text{STFT}(\mathbf{x})|^2$.
2. Structured model, e.g., nonnegative matrix factorization, deep neural networks.
3. Nonnegative masking and synthesis: $\mathbf{s}_j = \text{STFT}^{-1}(\mathbf{M}_j \mathbf{X})$.

1. Nonnegative representation, e.g., $\mathbf{V} = |\text{STFT}(\mathbf{x})|^2$.
2. Structured model, e.g., nonnegative matrix factorization, deep neural networks.
3. Nonnegative masking and synthesis: $\mathbf{s}_j = \text{STFT}^{-1}(\mathbf{M}_j \mathbf{X})$.

The phase problem $\setminus \mathbf{S}_j = \setminus \mathbf{X}$

7 Issues in sound quality when sources overlap.

7 *Inconsistency*: $\hat{\mathbf{S}}_j \not\subseteq \text{STFT}(\mathbb{R}^N)$.

Multiple Input Spectrogram Inversion

Algorithm overview

Multiple Input Spectrogram Inversion (MISI) [Gunawan, 2010]:

- . Extends the Griffin-Lim algorithm to multiple sources in mixture models.
- . Iterate the following updates on top of initial estimates:

STFT	$\mathbf{S}_j = \text{STFT}(\mathbf{s}_j)$
Magnitude modification	$\mathbf{Y}_j = \mathbf{V}_j \frac{\mathbf{S}_j}{j \mathbf{S}_j }$
Inverse STFT	$\mathbf{y}_j = \text{iSTFT}(\mathbf{Y}_j)$
Mixing	$\mathbf{s}_j = \mathbf{y}_j + \frac{1}{j} \sum_{i=1}^j \mathbf{y}_i$

Algorithm overview

Multiple Input Spectrogram Inversion (MISI) [Gunawan, 2010]:

- . Extends the Griffin-Lim algorithm to multiple sources in mixture models.
- . Iterate the following updates on top of initial estimates:

STFT	$\mathbf{S}_j = \text{STFT}(\mathbf{s}_j)$
Magnitude modification	$\mathbf{Y}_j = \mathbf{V}_j \frac{\mathbf{S}_j}{\ \mathbf{S}_j\ }$
Inverse STFT	$\mathbf{y}_j = \text{iSTFT}(\mathbf{Y}_j)$
Mixing	$\mathbf{s}_j = \mathbf{y}_j + \frac{1}{J} \sum_{i=1}^J \mathbf{y}_i$

- 3 Performance (post-processing, unfolded within end-to-end networks).
- 7 Convergence is only observed (no guarantee).
- 7 Offline processing, not applicable in real-time.

MISI derivation

Time-frequency formulation

- Main objective: reduce the magnitude mismatch $\sum_j \|\mathbf{S}_j\|_2^2 \approx \|\mathbf{X}\|_2^2$.
- Enforce consistency: $\mathbf{S}_j = \text{STFT}(\text{ISTFT}^{-1}(\mathbf{S}_j))$.
- Enforce a mixing constraint: $\mathbf{X} = \sum_{j=1}^J \mathbf{S}_j$

MISI derivation

Time-frequency formulation

- Main objective: reduce the magnitude mismatch $\sum_j \|S_j\|_2^2 \approx \|\mathbf{X}\|_2^2$.
 - Enforce consistency: $S_j = \text{STFT}(\text{ISTFT}^{-1}(S_j))$.
 - Enforce a mixing constraint: $\mathbf{X} = \sum_{j=1}^J S_j$
- 7 An ill-posed problem.

MISI derivation

Time-frequency formulation

- Main objective: reduce the magnitude mismatch $\sum_{j=1}^J \|\mathbf{S}_j\|_2 \approx \|\mathbf{V}\|_2^2$.
 - Enforce consistency: $\mathbf{S}_j = \text{STFT}^{-1}(\text{STFT}(\mathbf{S}_j))$.
 - Enforce a mixing constraint: $\mathbf{X} = \sum_{j=1}^J \mathbf{S}_j$
- 7 An ill-posed problem.

Time-domain formulation $\min_{\mathbf{s}_j} \sum_{j=1}^J \|\mathbf{V}_j \text{STFT}(\mathbf{s}_j)\|_2^2$ s.t. $\sum_{j=1}^J \mathbf{s}_j = \mathbf{x}$.

MISI derivation

Time-frequency formulation

- Main objective: reduce the magnitude mismatch $\sum_{j=1}^J \|\mathbf{S}_j\|_2 \approx \|\mathbf{X}\|_2$.
 - Enforce consistency: $\mathbf{S}_j = \text{STFT}^{-1}(\text{STFT}(\mathbf{S}_j))$.
 - Enforce a mixing constraint: $\mathbf{X} = \sum_{j=1}^J \mathbf{S}_j$
- 7 An ill-posed problem.

Time-domain formulation $\min_{\mathbf{s}_j} \sum_{j=1}^J \|\mathbf{V}_j \text{STFT}(\mathbf{s}_j)\|_2^2$ s.t. $\sum_{j=1}^J \mathbf{s}_j = \mathbf{x}$.

Majorization-minimization algorithm:

- Majorize the data fitting term:

$$\|\mathbf{V}_j \text{STFT}(\mathbf{s}_j)\|_2^2 \leq \|\mathbf{Y}_j \text{STFT}(\mathbf{s}_j)\|_2^2 \text{ with } \|\mathbf{Y}_j\|_2 = \|\mathbf{V}_j\|_2$$

- Incorporate the constraints using Lagrange multipliers.
- Find a saddle point for the majorizing function: 3 MISI with a convergence guarantee.

MISI derivation

Time-frequency formulation

- Main objective: reduce the magnitude mismatch $\sum_{j=1}^J \|\mathbf{S}_j\|_2 \approx \|\mathbf{X}\|_2$.
 - Enforce consistency: $\mathbf{S}_j = \text{STFT}^{-1}(\text{STFT}(\mathbf{S}_j))$.
 - Enforce a mixing constraint: $\mathbf{X} = \sum_{j=1}^J \mathbf{S}_j$
- 7 An ill-posed problem.

Time-domain formulation $\min_{\mathbf{s}_j} \sum_{j=1}^J \|\mathbf{V}_j \text{STFT}(\mathbf{s}_j)\|_2^2$ s.t. $\sum_{j=1}^J \mathbf{s}_j = \mathbf{x}$.

Majorization-minimization algorithm:

- Majorize the data fitting term:

$$\|\mathbf{V}_j \text{STFT}(\mathbf{s}_j)\|_2^2 \leq \|\mathbf{Y}_j \text{STFT}(\mathbf{s}_j)\|_2^2 \text{ with } \|\mathbf{Y}_j\|_2 = \|\mathbf{V}_j\|_2$$

- Incorporate the constraints using Lagrange multipliers.
- Find a saddle point for the majorizing function: 3 MISI with a convergence guarantee.

Online MISI (oMISI)

Problem: MISI involves the inverse STFT, which does not operate online:

$$\mathbf{s}_{j;t}^0 = \text{iDFT}(\mathbf{S}_{j;t}) \quad \mathbf{w} \quad \text{and} \quad \mathbf{s}_j(n) = \sum_{t=0}^{T-1} \mathbf{s}_{j;t}^0(n - t)$$

Online MISI (oMISI)

Problem: MISI involves the inverse STFT, which does not operate online:

$$\mathbf{s}_{j;t}^0 = \text{iDFT}(\mathbf{S}_{j;t}) \quad \mathbf{w} \quad \text{and} \quad \mathbf{s}_j(n) = \sum_{t=0}^{\infty} \mathbf{s}_{j;t}^0(n-t)$$

Approach: Only account for a limited amount of future time frames [Zhu, 2007]

Online MISI (oMISI)

Problem: MISI involves the inverse STFT, which does not operate online:

$$\mathbf{s}_{j;t}^{\circ} = \text{iDFT}(\mathbf{S}_{j;t}) \quad \mathbf{w} \quad \text{and} \quad \mathbf{s}_j(n) = \sum_{t=0}^{\infty} \mathbf{s}_{j;t}^{\circ}(n-t)$$

Approach: Only account for a limited amount of future time frames [Zhu, 2007]

- Split the overlap-add around the current frame:

$$\mathbf{s}_j(n) = \sum_{k=0}^{\infty} \mathbf{s}_{j;k}^{\circ}(n-t) + \sum_{k=t}^{\infty} \mathbf{s}_{j;k}^{\circ}(n-t)$$

\uparrow $\{Z\}$ \uparrow $\{Z\}$
past frames present and future frames

Online MISI (oMISI)

Problem: MISI involves the inverse STFT, which does not operate online:

$$\mathbf{s}_{j;t}^0 = \text{iDFT}(\mathbf{S}_{j;t}) \quad \mathbf{w} \quad \text{and} \quad \mathbf{s}_j(n) = \sum_{t=0}^{T-1} \mathbf{s}_{j;t}^0(n-t)$$

Approach: Only account for a limited amount of future time frames [Zhu, 2007]

- Split the overlap-add around the current frame:

$$\mathbf{s}_j(n) = \sum_{k=0}^{t-1} \mathbf{s}_{j;k}^0(n-t) + \sum_{k=t}^{K-1} \mathbf{s}_{j;k}^0(n-t)$$

$\underbrace{\quad}_{\text{past frames}} \quad \underbrace{\quad}_{\text{present and future frames}}$

- Only use K look-ahead future frames.

Initialization with the sinusoidal phase

oMISI allows for using alternative initialization schemes.

Initialization with the sinusoidal phase

oMISI allows for using alternative initialization schemes.

Sinusoidal model

- . Model each source as a sum of sinusoids.
- . The phase is given by:

$$f;t = f;t - 1 + 2 \left| \frac{f;t}{Z} \right|$$

normalized frequency

Initialization with the sinusoidal phase

oMISI allows for using alternative initialization schemes.

Sinusoidal model

- . Model each source as a sum of sinusoids.
- . The phase is given by:

$$f;t = f;t - 1 + 2 \left| \frac{f;t}{Z} \right|$$

normalized frequency

Frequency estimation with quadratic interpolation around each frequency peak.

Experiments

Task

- . Speech separation ($J = 2$) from the Danish HINT dataset.
- . Three speaker pairs (male+male, female+female, and male+female).

Two scenarios

- . \Oracle": ground truth magnitudes.
- . \Estim": magnitudes are estimated using a DNN [Naithani, 2017].

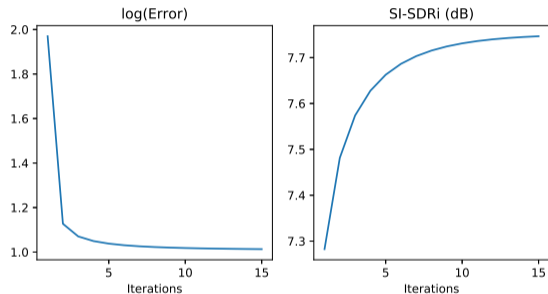
Baselines

- . Amplitude mask (AM).
- . MISI (online).

Metric: Scale-invariant signal-to-distortion ratio improvement (SI-SDRi, higher is better).

MISI convergence

In the Estim scenario:



- . Convergence is confirmed experimentally.
- . Performance (SI-SDRi) saturates at around 15 iterations (but further increases in the Oracle scenario).
- . oMISI will use $15=(K + 1)$ iterations for a fair comparison.

oMISI performance

With 50 % overlap:

	Latency	Male+Female		Male+Male		Female+Female	
		Estim	Oracle	Estim	Oracle	Estim	Oracle
AM	16 ms	7:5	8:8	5:7	7:3	5:1	7:5
MISI	online	7:9	23:8	6:2	22:3	5:4	22:9

. MISI > AM ! room for improvement for phase recovery.

oMISI performance

With 50 % overlap:

	Latency	Male+Female		Male+Male		Female+Female	
		Estim	Oracle	Estim	Oracle	Estim	Oracle
AM	16 ms	7:5	8:8	5:7	7:3	5:1	7:5
MISI	online	7:9	23:8	6:2	22:3	5:4	22:9
oMISI - mix	16 ms (K=0)	7:7	16:4	6:1	15:8	5:4	16:9
	24 ms (K=1)	7:9	20:2	6:2	19:4	5:4	19:6
	32 ms (K=2)	7:9	21:4	6:2	20:4	5:4	20:6

- . MISI > AM ! room for improvement for phase recovery.
- . oMISI with $K = 1$ performs as well as MISI (in the Estim. scenario).
 - . The optimal K depends on the overlap ratio (e.g., $K = 3$ for 75 %).

oMISI performance

With 50 % overlap:

	Latency	Male+Female		Male+Male		Female+Female	
		Estim	Oracle	Estim	Oracle	Estim	Oracle
AM	16 ms	7:5	8:8	5:7	7:3	5:1	7:5
MISI	online	7:9	23:8	6:2	22:3	5:4	22:9
oMISI - mix	16 ms (K=0)	7:7	16:4	6:1	15:8	5:4	16:9
	24 ms (K=1)	7:9	20:2	6:2	19:4	5:4	19:6
	32 ms (K=2)	7:9	21:4	6:2	20:4	5:4	20:6
oMISI - sin	24 ms (K=1)	7:8	15:2	6:2	14:6	5:4	20:7

- . MISI > AM ! room for improvement for phase recovery.
- . oMISI with $K = 1$ performs as well as MISI (in the Estim. scenario).
 - . The optimal K depends on the overlap ratio (e.g., $K = 3$ for 75 %).
- . The sinusoidal initialization is only interesting in a specific setting.

Conclusion

Contributions

- . A rigorous derivation of MISI with a convergence guarantee.
- . An online implementation with competitive separation performance and reduced latency.

Perspectives

- . Alternative loss functions (see our other ICASSP paper!)
- . Inclusion within deep learning for end-to-end separation.



<https://github.com/magronp/omisi>



https://magronp.github.io/demos/spl20_omisi.html

P. Magron, T. Virtanen, "Online spectrogram inversion for low-latency audio source separation", IEEE Signal Processing Letters, January 2020.