

Complex-valued and hybrid models for audio processing

PERCEPTION team seminar - INRIA Grenoble - Rhône-Alpes

Paul Magron

IRIT, Université de Toulouse, CNRS - Signal & Communications team

Academic background

- **PhD** – *Télécom ParisTech* (R. Badeau, B. David) 2013 - 2016
Phase recovery based on signal modeling: application to audio source separation.
- **Postdoc** – *Tampere University* (T. Virtanen) 2017 - 2019
Probabilistic phase modelling and real-time sound source separation.
- **Postdoc** – *IRIT* (C. Févotte) 2019 -
Content-aware music recommendation.

Publications since 2015:

- ▷ 4 journal articles.
 - ▷ 2 *IEEE Transactions on Audio, Speech, and Language Processing*
 - ▷ 1 *IEEE Signal Processing Letters*
 - ▷ 1 *IEEE Journal of Selected Topics in Signal Processing*
- ▷ 21 conference articles.
 - ▷ 1 best paper award (iWAENC 2018)
- ▷ 13 co-authors from 7 institutions.

Research field

Audio signal processing

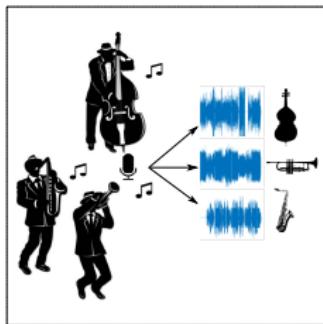
- ▷ Telecommunications: speech coding and synthesis.
- ▷ Medical devices: hearing aids.
- ▷ Educational softwares: language/music learning.
- ▷ Music information retrieval, music streaming.

Research field

Audio signal processing

- ▷ Telecommunications: speech coding and synthesis.
- ▷ Medical devices: hearing aids.
- ▷ Educational softwares: language/music learning.
- ▷ Music information retrieval, music streaming.

Main applications (of my work)



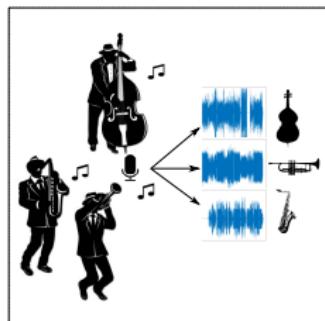
Source separation

Research field

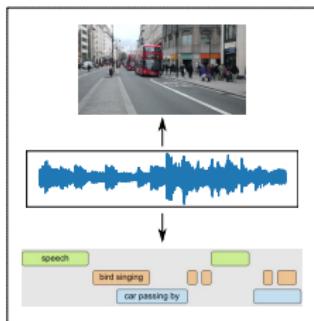
Audio signal processing

- ▷ Telecommunications: speech coding and synthesis.
- ▷ Medical devices: hearing aids.
- ▷ Educational softwares: language/music learning.
- ▷ Music information retrieval, music streaming.

Main applications (of my work)



Source separation



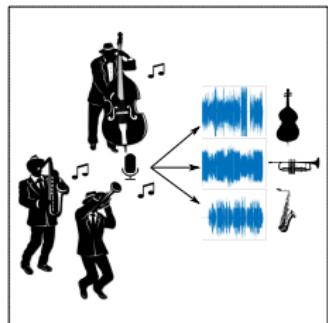
Acoustic scene analysis

Research field

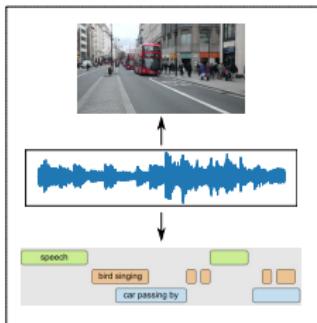
Audio signal processing

- ▷ Telecommunications: speech coding and synthesis.
- ▷ Medical devices: hearing aids.
- ▷ Educational softwares: language/music learning.
- ▷ Music information retrieval, music streaming.

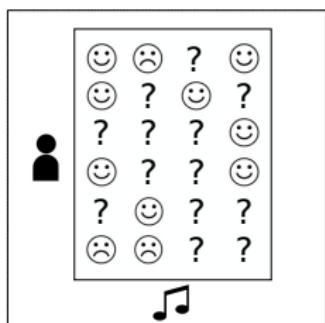
Main applications (of my work)



Source separation



Acoustic scene analysis



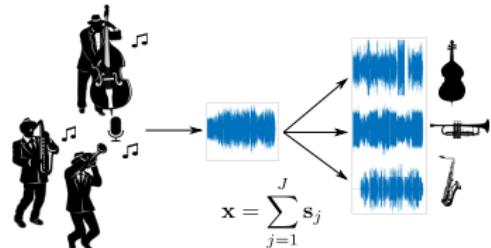
Music recommendation

Audio source separation

Audio source separation

Audio content is usually composed of several constitutive sounds:

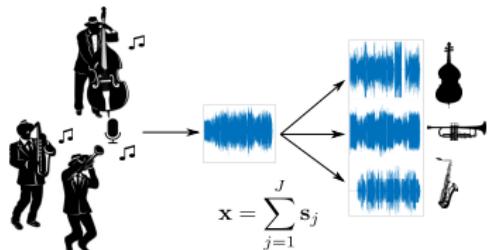
- ▷ One or several speakers + noise.
- ▷ Environmental / domestic sounds.
- ▷ Musical instruments.



Audio source separation

Audio content is usually composed of several constitutive sounds:

- ▷ One or several speakers + noise.
- ▷ Environmental / domestic sounds.
- ▷ Musical instruments.



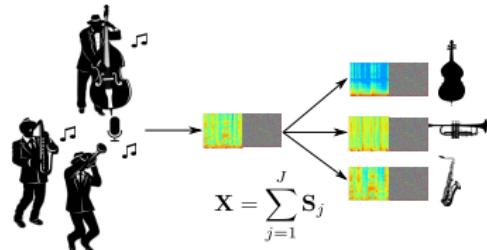
Source separation = recovering the sources from the mixture.

- ▷ Automatic speech recognition → clean speech vs. noise.
- ▷ Rhythm analysis → drums vs. harmonic instruments.
- ▷ Stationary / transient decomposition → time-stretching.

Audio source separation

Audio content is usually composed of several constitutive sounds:

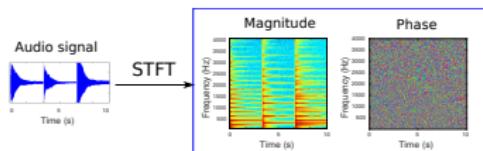
- ▷ One or several speakers + noise.
- ▷ Environmental / domestic sounds.
- ▷ Musical instruments.



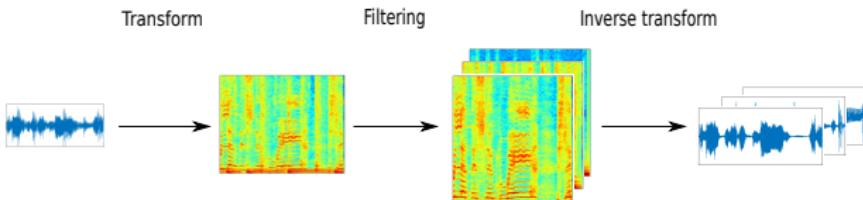
Source separation = recovering the sources from the mixture.

- ▷ Automatic speech recognition → clean speech vs. noise.
- ▷ Rhythm analysis → drums vs. harmonic instruments.
- ▷ Stationary / transient decomposition → time-stretching.

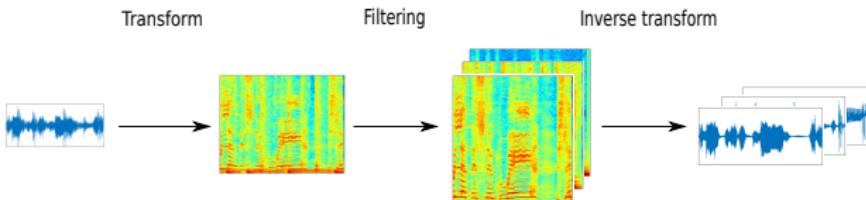
Time-frequency separation = acts on the short-time Fourier transform (STFT).



Source separation general framework

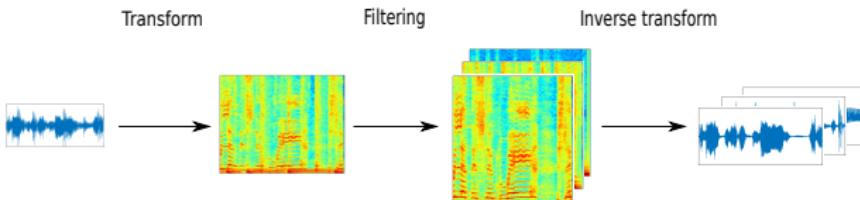


Source separation general framework



1. Nonnegative representation: $V = |\text{STFT}(x)|^2$.

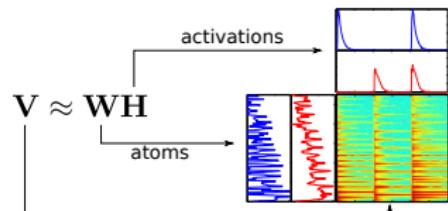
Source separation general framework



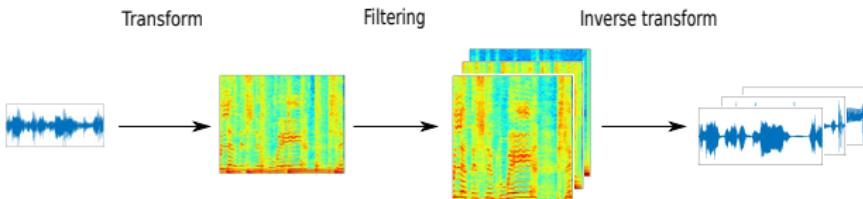
1. Nonnegative representation: $V = |\text{STFT}(x)|^2$.

2. Structured model:

- ▷ Nonnegative matrix factorization (NMF).
- ▷ Deep neural network (DNN).



Source separation general framework



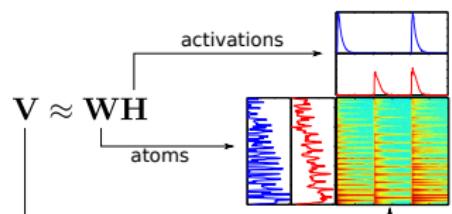
1. Nonnegative representation: $V = |\text{STFT}(x)|^2$.

2. Structured model:

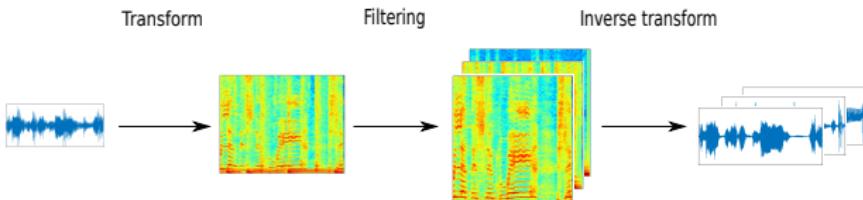
- ▷ Nonnegative matrix factorization (NMF).
- ▷ Deep neural network (DNN).

3. Estimation:

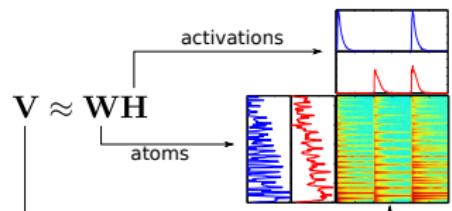
- ▷ $\min_{W,H} \mathcal{D}(V, WH)$
- ▷ $\min_{\theta} \mathcal{L}(\{V_j\}, \phi_{\theta}(V))$



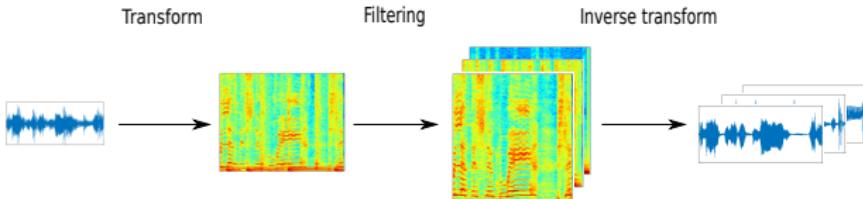
Source separation general framework



1. Nonnegative representation: $V = |\text{STFT}(x)|^2$.
2. Structured model:
 - ▷ Nonnegative matrix factorization (NMF).
 - ▷ Deep neural network (DNN).
3. Estimation:
 - ▷ $\min_{W,H} \mathcal{D}(V, WH)$
 - ▷ $\min_{\theta} \mathcal{L}(\{V_j\}, \phi_{\theta}(V))$
4. Set of nonnegative masks M_j .



Source separation general framework



1. Nonnegative representation: $V = |\text{STFT}(x)|^2$.

2. Structured model:

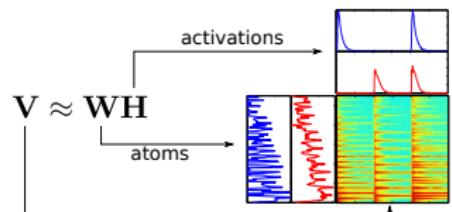
- ▷ Nonnegative matrix factorization (NMF).
- ▷ Deep neural network (DNN).

3. Estimation:

- ▷ $\min_{W,H} \mathcal{D}(V, WH)$
- ▷ $\min_{\theta} \mathcal{L}(\{V_j\}, \phi_{\theta}(V))$

4. Set of nonnegative masks M_j .

5. Synthesis: $\tilde{s}_j = \text{STFT}^{-1}(M_j \odot X)$.



The phase problem

Wiener filter

$$\hat{S}_j = \frac{\hat{V}_j}{\sum_{k=1}^J \hat{V}_k} \odot X$$

Source STFT (Estimated) Mask Mixture STFT

$\angle S_j = \angle X \rightarrow$ Issues in sound quality when sources overlap in the TF domain.

The phase problem

Wiener filter

$$\hat{S}_j = \frac{\hat{V}_j}{\sum_{k=1}^J \hat{V}_k} \odot X$$

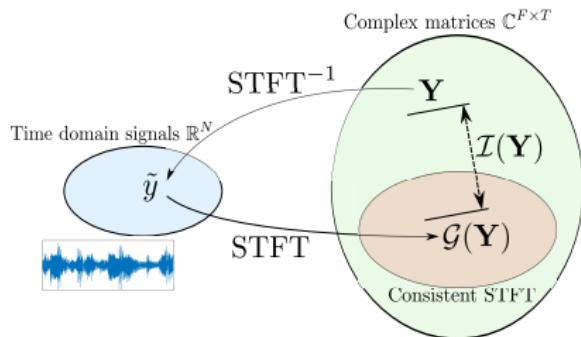
Source STFT (Estimated) Mask Mixture STFT

$\angle S_j = \angle X \rightarrow$ Issues in sound quality when sources overlap in the TF domain.

Inconsistency

$$\mathcal{I}(Y) = \|Y - \mathcal{G}(Y)\|^2$$

$$\mathcal{G} = \text{STFT} \circ \text{STFT}^{-1}$$



The phase problem

Wiener filter

$$\hat{S}_j = \frac{\hat{V}_j}{\sum_{k=1}^J \hat{V}_k} \odot X$$

Source STFT (Estimated) Mask Mixture STFT

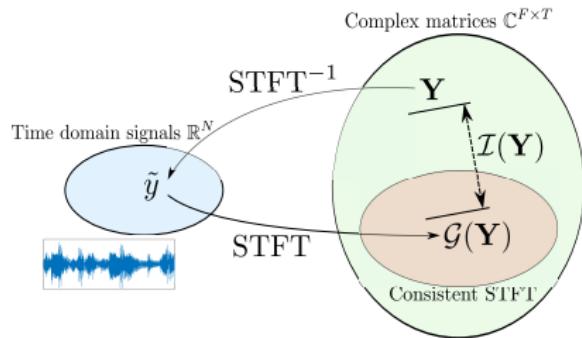
$\angle S_j = \angle X \rightarrow$ Issues in sound quality when sources overlap in the TF domain.

Inconsistency

$$\mathcal{I}(Y) = \|Y - \mathcal{G}(Y)\|^2$$

$$\mathcal{G} = \text{STFT} \circ \text{STFT}^{-1}$$

- ▷ Griffin-Lim (GL): minimization of \mathcal{I} .
- ▷ Extension to source separation.
- ▷ Combination with Wiener filtering.



The phase problem

Comparative study [ICASSP 15]

- ▷ Room for improvement for phase recovery.
- ▷ Test alternative separation methods accounting for the phase:
 - ▷ Consistency-based.
 - ▷ Signal model-based.
- ▷ The most promising uses a signal model for structuring the phase / STFT.

The phase problem

Comparative study [ICASSP 15]

- ▷ Room for improvement for phase recovery.
- ▷ Test alternative separation methods accounting for the phase:
 - ▷ Consistency-based.
 - ▷ Signal model-based.
- ▷ The most promising uses a signal model for structuring the phase / STFT.

My approach

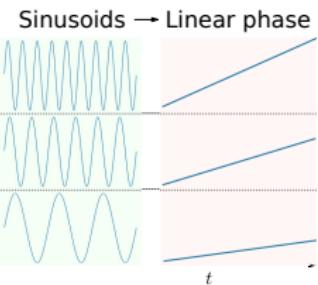
- ▷ Leveraging model-based phase properties in source separation.
- ▷ A phase-aware probabilistic framework.
- ▷ Joint estimation of magnitude and phase.

Model-based phase recovery

Sinusoidal phase model

For a mixture of sinusoids, the phase is:

$$\mu_{f,t} = \mu_{f,t-1} + 2\pi \underbrace{\nu_{f,t}}_{\text{normalized frequency}}$$



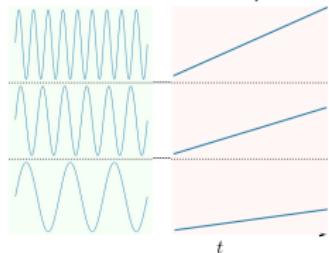
Sinusoidal phase model

For a mixture of sinusoids, the phase is:

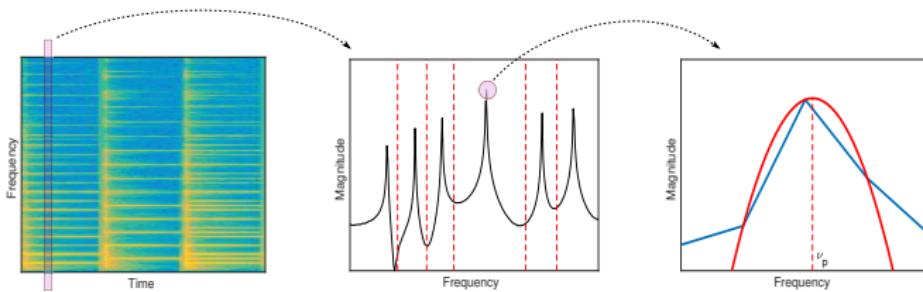
$$\mu_{f,t} = \mu_{f,t-1} + 2\pi$$

$$\underbrace{\nu_{f,t}}_{\text{normalized frequency}}$$

Sinusoids \rightarrow Linear phase

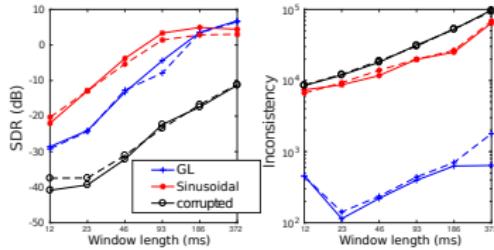


Recursive estimation in each time frame:



Sinusoidal phase model

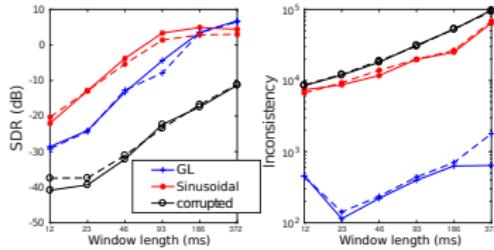
Restoration of piano pieces:



- ▷ Better performance than the consistency-based GL algorithm.
- ▷ But overall low SDR: error propagates over time frames.

Sinusoidal phase model

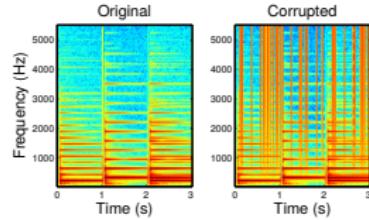
Restoration of piano pieces:



- ▷ Better performance than the consistency-based GL algorithm.
- ▷ But overall low SDR: error propagates over time frames.

Applications

- ▷ Click removal [EUSIPCO 15].
- ▷ Source separation.



Repeating attacks

Why treating attacks?

- ▷ Perceptive quality of the sound.
- ▷ Initialize the sinusoidal model.

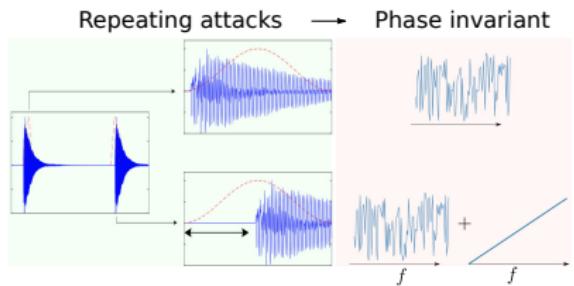
Repeating attacks

Why treating attacks?

- ▷ Perceptive quality of the sound.
- ▷ Initialize the sinusoidal model.

Model within onset frames:

$$\mu_{f,t} = \underbrace{\psi_f}_{\text{invariant}} + \underbrace{\eta_t f}_{\text{offset}}$$



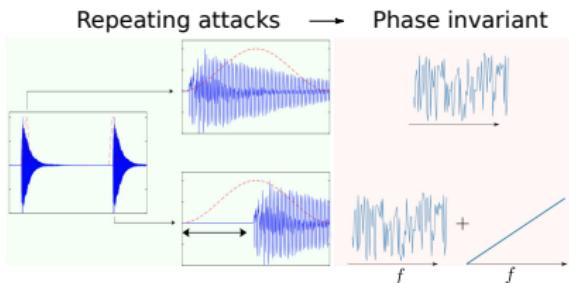
Repeating attacks

Why treating attacks?

- ▷ Perceptive quality of the sound.
- ▷ Initialize the sinusoidal model.

Model within onset frames:

$$\mu_{f,t} = \underbrace{\psi_f}_{\text{invariant}} + \underbrace{\eta_t}_{\text{offset}} f$$



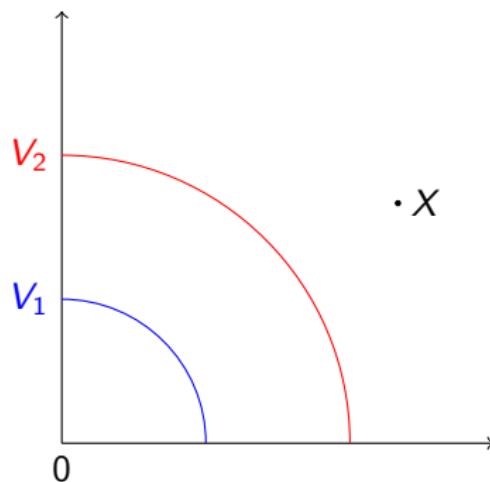
Incorporation in a mixture model

- ▷ Estimation with coordinate descent.
- ▷ Slight improvement over using the mixture's phase [WASPAA 15].

Source separation iterative algorithm

Problem

$$\min_{\mu} \left| \left| \mathbf{X} - \sum_{j=1}^J V_j e^{i\mu_j} \right| \right|^2$$



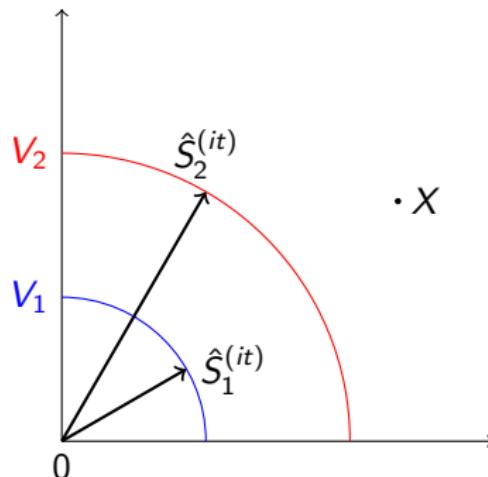
Source separation iterative algorithm

Problem

$$\min_{\mu} \left| \left| \mathbf{X} - \sum_{j=1}^J V_j e^{i\mu_j} \right| \right|^2$$

Optimization

- ▷ Use a phase model for initialization.



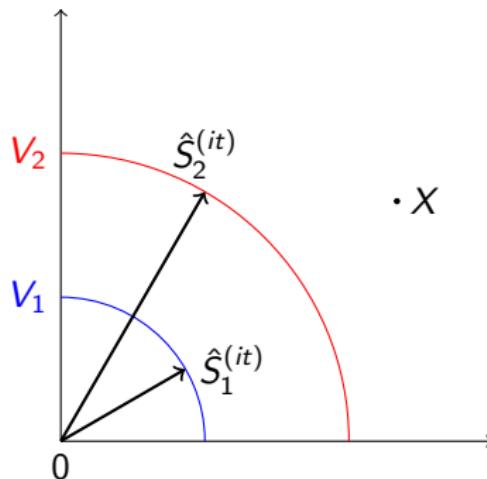
Source separation iterative algorithm

Problem

$$\min_{\mu} \left| \left| \mathbf{X} - \sum_{j=1}^J V_j e^{i\mu_j} \right| \right|^2$$

Optimization

- ▷ Use a phase model for initialization.
- ▷ Auxiliary function method: iterative procedure.



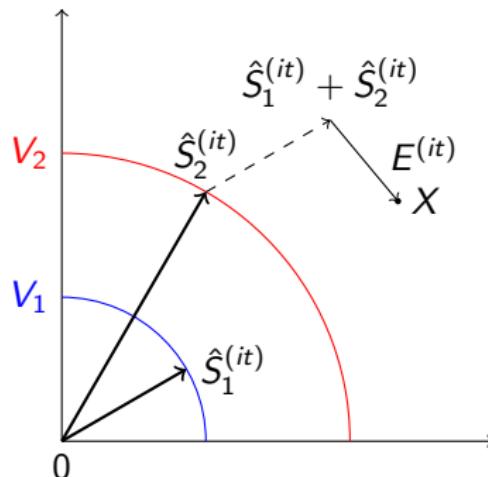
Source separation iterative algorithm

Problem

$$\min_{\mu} \left| \left| \mathbf{X} - \sum_{j=1}^J V_j e^{i\mu_j} \right| \right|^2$$

Optimization

- ▷ Use a phase model for initialization.
- ▷ Auxiliary function method: iterative procedure.



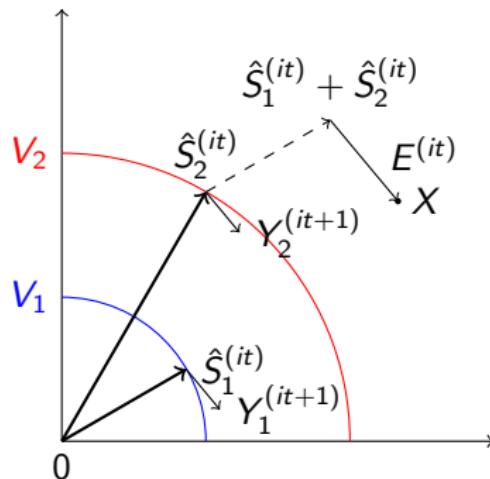
Source separation iterative algorithm

Problem

$$\min_{\mu} \left| \left| \mathbf{X} - \sum_{j=1}^J V_j e^{i\mu_j} \right| \right|^2$$

Optimization

- ▷ Use a phase model for initialization.
- ▷ Auxiliary function method: iterative procedure.



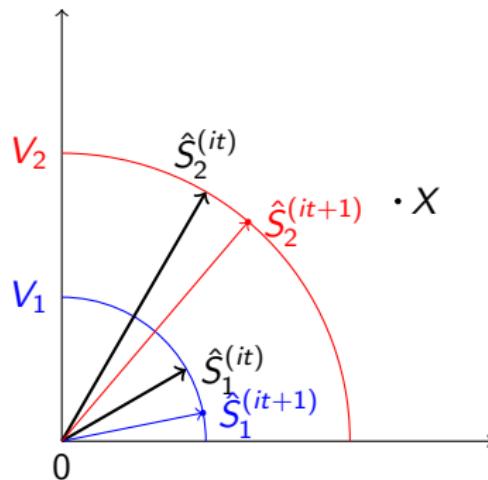
Source separation iterative algorithm

Problem

$$\min_{\mu} \left| \left| \mathbf{X} - \sum_{j=1}^J V_j e^{i\mu_j} \right| \right|^2$$

Optimization

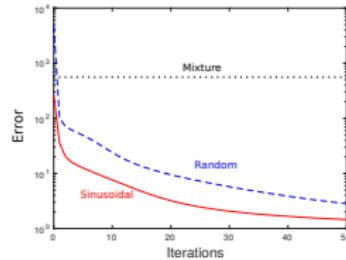
- ▷ Use a phase model for initialization.
- ▷ Auxiliary function method: iterative procedure.



Source separation iterative algorithm

Initialization impact (DSD100 dataset):

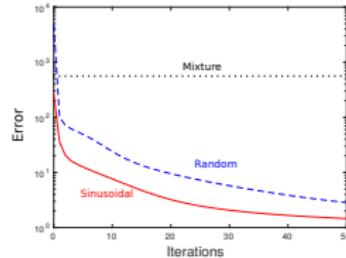
	SDR	SIR	SAR
Mixture	7.5	13.7	8.9
Random	9.5	22.8	9.7
Sinusoidal	13.6	31.0	13.7



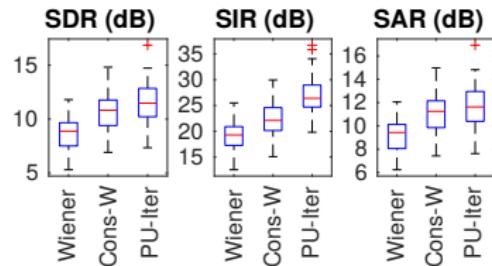
Source separation iterative algorithm

Initialization impact (DSD100 dataset):

	SDR	SIR	SAR
Mixture	7.5	13.7	8.9
Random	9.5	22.8	9.7
Sinusoidal	13.6	31.0	13.7



Comparison with Wiener filters:



How to re-introduce consistency?

Time-domain formulation

$$\min_{\{\tilde{s}_j\}} \sum_j |||\text{STFT}(\tilde{s}_j)| - V_j||^2 \text{ s. t. } \sum_j \tilde{s}_j = \tilde{x}$$

How to re-introduce consistency?

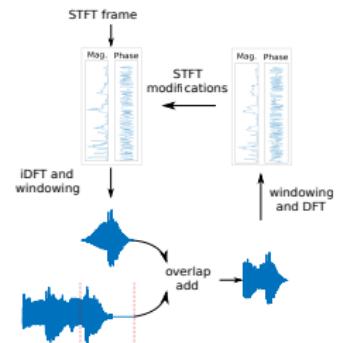
Time-domain formulation

$$\min_{\{\tilde{s}_j\}} \sum_j |||\text{STFT}(\tilde{s}_j)| - V_j||^2 \text{ s. t. } \sum_j \tilde{s}_j = \tilde{x}$$

- ▷ Auxiliary function method: the (well-known) MISI algorithm, but convergence-guaranteed.

Online implementation

- ▷ Useful for real-time applications (hearing aids): for a latency of 16 ms, same results as the offline counterpart.
- ▷ Allows initialization with the sinusoidal phase: better results in some cases.



Alternative divergences

- ▷ Euclidean distance: not the most appropriate in audio.
- ▷ Popular alternative: the beta-divergences.

$\beta = 2$ Euclidean Emphasis on high-energy components	$\beta = 1$ Kullback-Leibler (KL) ← In between →	$\beta = 0$ Itakura-Saito (IS) Scale invariance
--	--	---

P. Magron, P.-H. Vial, T. Oberlin, C. Févotte, "Phase recovery with Bregman divergences for audio source separation", *Proc. IEEE International Conference on Audio, Speech and Signal Processing (ICASSP)*, June 2021.

P.-H. Vial, P. Magron, T. Oberlin, C. Févotte, "Phase retrieval with Bregman divergences and application to audio signal recovery", *IEEE Journal of Selected Topics in Signal Processing*, January 2021.

Alternative divergences

- ▷ Euclidean distance: not the most appropriate in audio.
- ▷ Popular alternative: the beta-divergences.

$\beta = 2$ Euclidean Emphasis on high-energy components	$\beta = 1$ Kullback-Leibler (KL) ← In between →	$\beta = 0$ Itakura-Saito (IS) Scale invariance
--	--	---

Phase retrieval with beta-divergences

$$\min_{\{\tilde{s}_j\}} \sum_j D_\beta(|\text{STFT}(\tilde{s}_j)|^d, V_j) \text{ s. t. } \sum_j \tilde{s}_j = \tilde{x}$$

- ▷ Optimization with accelerated gradient descent or ADMM.
- ▷ Experimentally: alternative divergences (e.g., KL or $\beta = 0.5$) > Euclidean.

P. Magron, P.-H. Vial, T. Oberlin, C. Févotte, "Phase recovery with Bregman divergences for audio source separation", *Proc. IEEE International Conference on Audio, Speech and Signal Processing (ICASSP)*, June 2021.

P.-H. Vial, P. Magron, T. Oberlin, C. Févotte, "Phase retrieval with Bregman divergences and application to audio signal recovery", *IEEE Journal of Selected Topics in Signal Processing*, January 2021.

Probabilistic phase modelling

Probabilistic framework

Why?

- ▷ Modeling uncertainty.
- ▷ Incorporating prior information.
- ▷ Obtaining estimators with nice statistical properties.
- ▷ Deriving inference schemes with convergence guarantees.

Probabilistic framework

Why?

- ▷ Modeling uncertainty.
- ▷ Incorporating prior information.
- ▷ Obtaining estimators with nice statistical properties.
- ▷ Deriving inference schemes with convergence guarantees.

Traditionally

Circularly-symmetric (or **isotropic**) sources \iff Uniform phase
 \Rightarrow Phase-unaware estimators.

Probabilistic framework

Why?

- ▷ Modeling uncertainty.
- ▷ Incorporating prior information.
- ▷ Obtaining estimators with nice statistical properties.
- ▷ Deriving inference schemes with convergence guarantees.

Traditionally

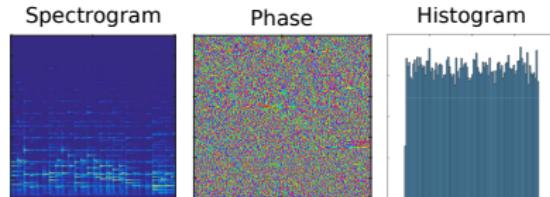
Circularly-symmetric (or **isotropic**) sources \iff Uniform phase
 \Rightarrow Phase-unaware estimators.

My approach

A phase-aware probabilistic framework for source separation.

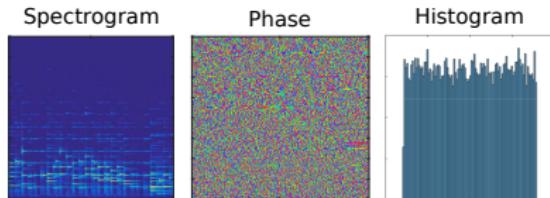
Is the phase really uniform?

A simple example (piano piece), where the phase appears uniformly-distributed.



Is the phase really uniform?

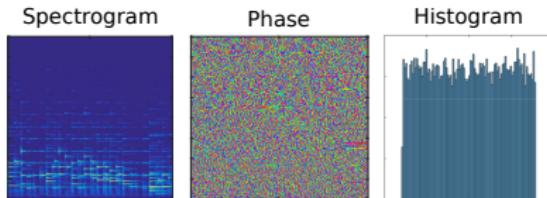
A simple example (piano piece), where the phase appears uniformly-distributed.



But: is that consistent with, e.g., the sinusoidal model?

Is the phase really uniform?

A simple example (piano piece), where the phase appears uniformly-distributed.



But: is that consistent with, e.g., the sinusoidal model?

Interpretation

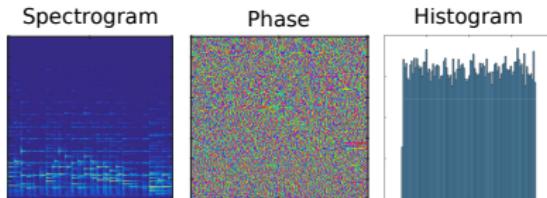
- ▷ The histogram validates an iid assumption on $\{\phi_{f,t}\}$:

$$\phi_{f,t} \sim \mathcal{D} \text{ and independent } \rightarrow \mathcal{D} = \mathcal{U}_{[0,2\pi]}$$

- ▷ This model only conveys a **global** information.

Is the phase really uniform?

A simple example (piano piece), where the phase appears uniformly-distributed.



But: is that consistent with, e.g., the sinusoidal model?

Interpretation

- ▷ The histogram validates an iid assumption on $\{\phi_{f,t}\}$:

$$\phi_{f,t} \sim \mathcal{D} \text{ and independent } \rightarrow \mathcal{D} = \mathcal{U}_{[0,2\pi]}$$

- ▷ This model only conveys a **global** information.

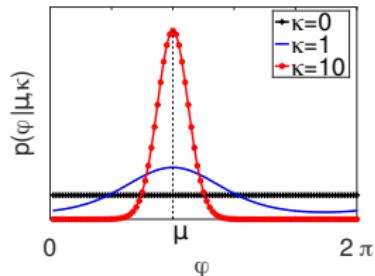
What about the **local structure** of the phase?

Von Mises phase model

Von Mises distribution

$$\phi_{ft} \sim \mathcal{VM}(\mu_{ft}, \kappa)$$

- ▷ $\mu_{f,t}$ = phase location.
- ▷ κ = concentration (quantifies non-uniformity).

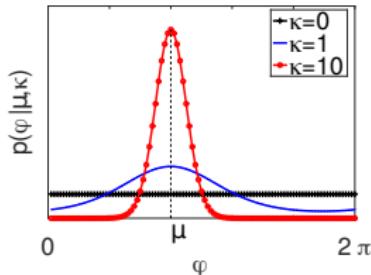


Von Mises phase model

Von Mises distribution

$$\phi_{ft} \sim \mathcal{VM}(\mu_{ft}, \kappa)$$

- ▷ $\mu_{f,t}$ = phase location.
- ▷ κ = concentration (quantifies non-uniformity).



Model

- ▷ $\mu_{f,t}$ = sinusoidal phase.
- ▷ Center the phases: $\psi_{f,t} = \phi_{f,t} - \mu_{f,t}$

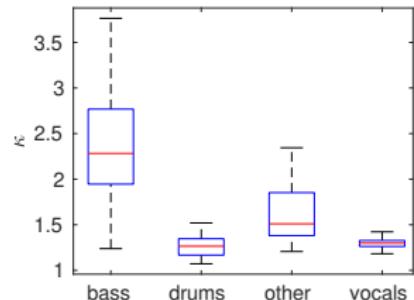
Distribution	Uniform	VM	Centered VM
iid	$\phi_{f,t} \sim \mathcal{U}_{[0,2\pi]}$	$\phi_{ft} \sim \mathcal{VM}(\mu_{f,t}, \kappa)$	$\psi_{f,t} \sim \mathcal{VM}(0, \kappa)$
Local structure	✗	✓	✓

Von Mises phase model [iWAENC 18]

Estimation of κ (maximum likelihood):

$$\frac{l_1(\kappa)}{l_0(\kappa)} = \frac{1}{FT} \sum_{f,t} \cos(\psi_{ft})$$

- ▷ Solved with fast numerical schemes.
- ▷ κ quantifies the “sinusoidality” of the sources.

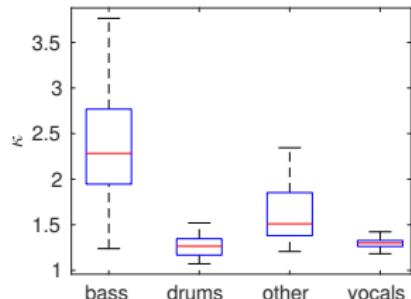


Von Mises phase model [iWAENC 18]

Estimation of κ (maximum likelihood):

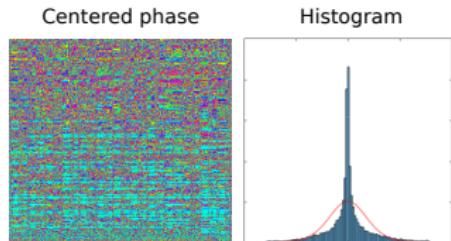
$$\frac{I_1(\kappa)}{I_0(\kappa)} = \frac{1}{FT} \sum_{f,t} \cos(\psi_{ft})$$

- ▷ Solved with fast numerical schemes.
- ▷ κ quantifies the “sinusoidality” of the sources.



Validation

- ▷ Both uniform and VM models are statistically relevant.
- ▷ They convey different information about the phase (global vs. local).



Multiple sources model

In each TF bin, $x = \sum_{j=1}^J s_j$.

Isotropic gaussian model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.

- ▷ m_j : mean (location) / γ_j : variance (energy).
- ▷ c_j : relation term, usually 0.

Multiple sources model

In each TF bin, $x = \sum_{j=1}^J s_j$.

Isotropic gaussian model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.

- ▷ m_j : mean (location) / γ_j : variance (energy).
- ▷ c_j : relation term, usually 0.

Polar coordinate equivalent model: $s_j = r_j e^{i\phi_j}$ where:

- ▷ $r_j \sim \mathcal{R}(v_j)$ (Rayleigh magnitude).
- ▷ $\phi_j \sim \mathcal{U}_{[0, 2\pi[}$ (uniform phase).

Multiple sources model

In each TF bin, $x = \sum_{j=1}^J s_j$.

Isotropic gaussian model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.

- ▷ m_j : mean (location) / γ_j : variance (energy).
- ▷ c_j : relation term, usually 0.

Polar coordinate equivalent model: $s_j = r_j e^{i\phi_j}$ where:

- ▷ $r_j \sim \mathcal{R}(v_j)$ (Rayleigh magnitude).
- ▷ $\phi_j \sim \mathcal{U}_{[0, 2\pi[}$ (uniform phase).

Naive idea

Consider a VM phase instead → Phase-aware

Multiple sources model

In each TF bin, $x = \sum_{j=1}^J s_j$.

Isotropic gaussian model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.

- ▷ m_j : mean (location) / γ_j : variance (energy).
- ▷ c_j : relation term, usually 0.

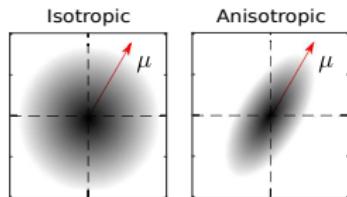
Polar coordinate equivalent model: $s_j = r_j e^{i\phi_j}$ where:

- ▷ $r_j \sim \mathcal{R}(v_j)$ (Rayleigh magnitude).
- ▷ $\phi_j \sim \mathcal{U}_{[0, 2\pi[}$ (uniform phase).

Naive idea

Consider a VM phase instead → Phase-aware... but not tractable.

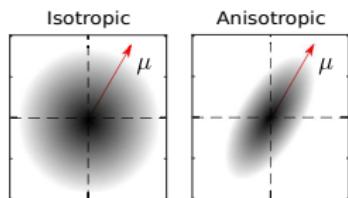
Anisotropic Gaussian model



	Phase-aware	Tractable
Isotropic Gaussian	✗	✓
Rayleigh + von Mises	✓	✗
Anisotropic Gaussian	✓	✓

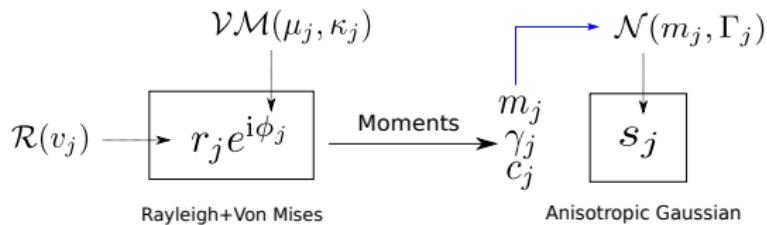
Source model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.

Anisotropic Gaussian model

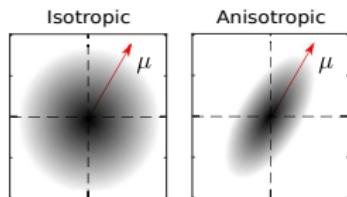


	Phase-aware	Tractable
Isotropic Gaussian	✗	✓
Rayleigh + von Mises	✓	✗
Anisotropic Gaussian	✓	✓

Source model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.

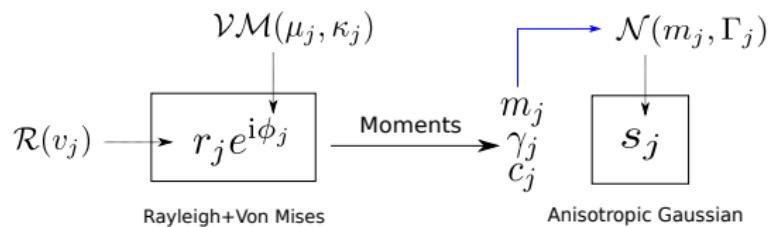


Anisotropic Gaussian model



	Phase-aware	Tractable
Isotropic Gaussian	✗	✓
Rayleigh + von Mises	✓	✗
Anisotropic Gaussian	✓	✓

Source model: $s_j \sim \mathcal{N}_{\mathbb{C}}(m_j, \Gamma_j)$ with $\Gamma_j = \begin{pmatrix} \gamma_j & c_j \\ \bar{c}_j & \gamma_j \end{pmatrix}$.



- Parameters: v_j (energy), μ_j (phase location), κ (non-uniformity).

Anisotropic Wiener filter

- ▷ Posterior mean of the sources: anisotropic Wiener filter [ICASSP 17].
- ▷ Performance (oracle separation results):

	SDR	SIR	SAR
Wiener	8.5	19.1	9.1
Anisotropic Wiener	9.7	21.9	10.1

Main message

- ▷ Including phase information in the filter improves the separation quality.
- ▷ Potential of a phase-aware statistical framework.

Again... consistency?

Problem

- ▷ The (anisotropic) Wiener filter produces inconsistent matrices.

Again... consistency?

Problem

- ▷ The (anisotropic) Wiener filter produces inconsistent matrices.

Consistent anisotropic Wiener [WASPAA 17]

- ▷ Consider the loss function:
 $\{\text{posterior distribution of the (anisotropic) sources}\} + \{\text{consistency constraint}\}$
- ▷ Minimization with preconditioned conjugate gradient descent.

Again... consistency?

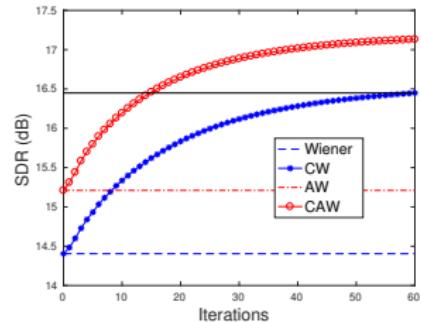
Problem

- ▷ The (anisotropic) Wiener filter produces inconsistent matrices.

Consistent anisotropic Wiener [WASPAA 17]

- ▷ Consider the loss function:
 $\{\text{posterior distribution of the (anisotropic) sources}\} + \{\text{consistency constraint}\}$
- ▷ Minimization with preconditioned conjugate gradient descent.

	Sinusoidal model	Consistent
Wiener	✗	✗
CW	✗	✓
AW	✓	✗
CAW	✓	✓

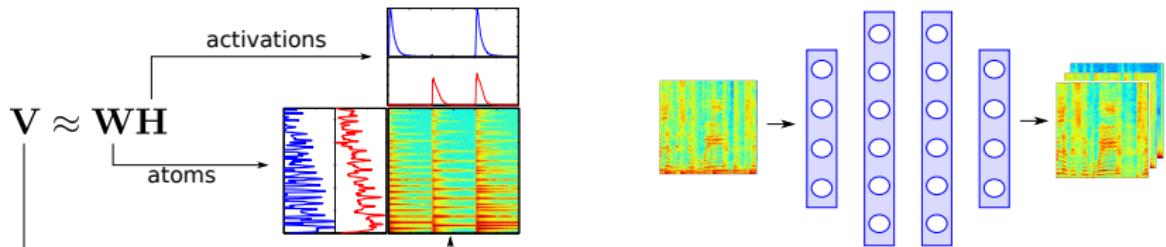


Joint estimation of magnitude and phase

Realistic source separation

Goal: estimate the magnitude **and** the phase of the sources.

- ▷ Needs an additional spectrogram-like model and estimation technique.



Approaches

- ▷ Two-stage: first estimate the magnitude, and then recover the phase.
- ▷ One-stage: jointly estimate the magnitude and the phase.

Two-stage approaches

NMF + phase recovery

- ▷ Phase recovery induces a slight improvement (interference reduction).

P. Magron, K. Drossos, S. I. Mimalakis, T. Virtanen, "Reducing interference with phase recovery in DNN-based monaural singing voice separation", *Proc. Interspeech*. September 2018.

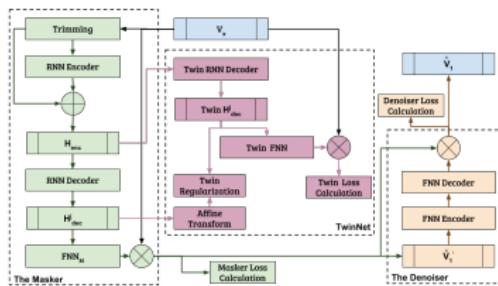
K. Drossos, P. Magron, S. I. Mimalakis, T. Virtanen, "Harmonic-percussive source separation with deep neural networks and phase recovery", *Proc. International Workshop on Acoustic Signal Enhancement (IWAENC)*, September 2018.

Two-stage approaches

NMF + phase recovery

- ▷ Phase recovery induces a slight improvement (interference reduction).

DNN + phase recovery [Interspeech 18, iWAENC 18]



- ▷ More significant results (DNNs > NMF).
- ▷ Phase recovery makes sense on top of good magnitude estimates.

P. Magron, K. Drossos, S. I. Mimalakis, T. Virtanen, "Reducing interference with phase recovery in DNN-based monaural singing voice separation", *Proc. Interspeech*. September 2018.

K. Drossos, P. Magron, S. I. Mimalakis, T. Virtanen, "Harmonic-percussive source separation with deep neural networks and phase recovery", *Proc. International Workshop on Acoustic Signal Enhancement (IWAENC)*, September 2018.

Complex NMF

Baseline NMF $|X| \approx \hat{V} = WH$

- ▷ Assumes the additivity of the sources' magnitudes → phase?

Complex NMF

Baseline NMF $|X| \approx \hat{V} = WH$

- ▷ Assumes the additivity of the sources' magnitudes \rightarrow phase?

Complex NMF $X \approx \hat{X} = \sum_j [\underbrace{W_j H_j}_{\text{NMF}}] e^{i\mu_j}$

- ▷ Regularize the phases with model-based properties [ICASSP 16]:

$$\min ||X - \hat{X}||^2 + \underbrace{\lambda_{sp} \mathcal{C}_{sp}(H)}_{\text{sparsity}} + \underbrace{\lambda_{sin} \mathcal{C}_{sin}(\mu)}_{\text{sinusoidal phase}} + \underbrace{\lambda_{rep} \mathcal{C}_{rep}(\mu, \psi, \eta)}_{\text{repeating attacks}}$$

- ▷ Optimization: coordinate descent or auxiliary function method.
- ▷ Tuning λ_{sin} and λ_{rep} : trade-off between interference and artifact reduction.

Complex ISNMF

Hard to extend complex NMF to non-Euclidean metrics (e.g., beta-divergences).

Complex ISNMF

Hard to extend complex NMF to non-Euclidean metrics (e.g., beta-divergences).

Probabilistic view

(Real) Gaussian	$r \sim \mathcal{N}(m, \sigma^2)$	$m = wh$	Euclidean NMF
Poisson	$r \sim \mathcal{P}(v)$	$v = wh$	KLNMF
Isotropic Gaussian	$x \sim \mathcal{N}_{\mathbb{C}}(0, v^2 I)$	$v^2 = wh$	ISNMF
Isotropic Gaussian	$x \sim \mathcal{N}_{\mathbb{C}}(m, \sigma^2 I)$	$m = whe^{i\mu}$	Complex NMF

Complex ISNMF

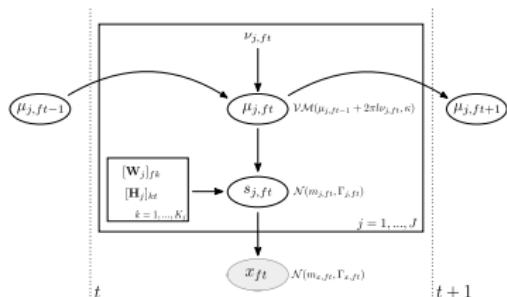
Hard to extend complex NMF to non-Euclidean metrics (e.g., beta-divergences).

Probabilistic view

(Real) Gaussian	$r \sim \mathcal{N}(m, \sigma^2)$	$m = wh$	Euclidean NMF
Poisson	$r \sim \mathcal{P}(v)$	$v = wh$	KLNMF
Isotropic Gaussian	$x \sim \mathcal{N}_{\mathbb{C}}(0, v^2 I)$	$v^2 = wh$	ISNMF
Isotropic Gaussian	$x \sim \mathcal{N}_{\mathbb{C}}(m, \sigma^2 I)$	$m = whe^{i\mu}$	Complex NMF

Complex ISNMF [TASLP 19]

- ▷ Anisotropic Gaussian model (NMF variance).
- ▷ Markov chain prior on the phase parameter.
- ▷ Estimation: EM algorithm.
- ▷ Better results than Complex NMF and ISNMF.

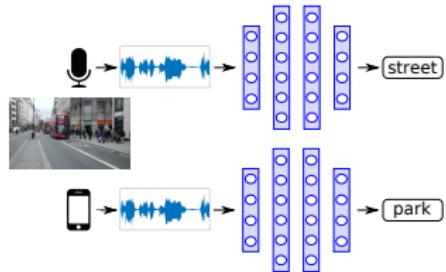


Acoustic scene classification

Domain adaptation for acoustic scene classification

Problem

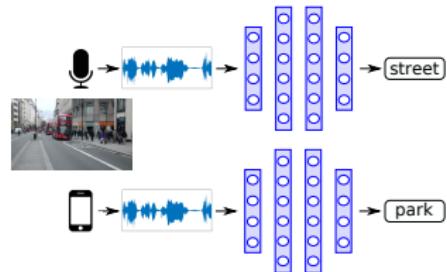
- ▷ Classify an audio segment into classes.
- ▷ Mismatch between recording conditions.



Domain adaptation for acoustic scene classification

Problem

- ▷ Classify an audio segment into classes.
- ▷ Mismatch between recording conditions.



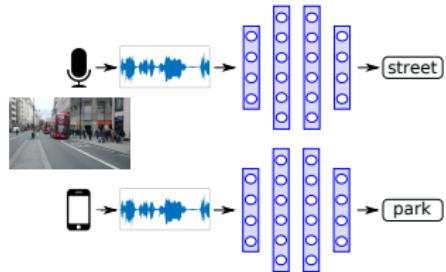
Domain adaptation [WASPAA 19]

- ▷ *Source* and *target* domain data → same latent distribution.
- ▷ Wasserstein GAN: minimize scene classification error / maximize domain classification error.

Domain adaptation for acoustic scene classification

Problem

- ▷ Classify an audio segment into classes.
- ▷ Mismatch between recording conditions.



Domain adaptation [WASPAA 19]

- ▷ *Source* and *target* domain data → same latent distribution.
- ▷ Wasserstein GAN: minimize scene classification error / maximize domain classification error.

	airport	bus	metro	park	public_square	shopping_mall	street_pedestrian	street_traffic	train	airport	bus	metro	park	public_square	shopping_mall	street_pedestrian	street_traffic	train
airport	0.06 0.00 0.19	0.61	0.00 0.00 0.00 0.00 0.06	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.00 0.58 0.00 0.11 0.03 0.08 0.08 0.00 0.14 0.03 0.03	0.00 0.58 0.00 0.11 0.03 0.08 0.08 0.00 0.14 0.03 0.03										
bus	0.00 0.22 0.28	0.39	0.03 0.00 0.00 0.00 0.08	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.03 0.14 0.28 0.17 0.00 0.06 0.06 0.00 0.17 0.11 0.06	0.03 0.14 0.28 0.17 0.00 0.06 0.06 0.00 0.17 0.11 0.06										
metro	0.00 0.00 0.00	0.36	0.61	0.00 0.00 0.00 0.00 0.03	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.00 0.18 0.06 0.44 0.00 0.08 0.08 0.00 0.06 0.00	0.00 0.18 0.06 0.44 0.00 0.08 0.08 0.00 0.06 0.00									
park	0.00 0.06 0.22	0.61	0.00 0.00 0.00 0.00 0.06	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.00 0.03 0.00 0.03 0.58 0.11 0.00 0.06 0.14 0.06	0.00 0.03 0.00 0.03 0.58 0.11 0.00 0.06 0.14 0.06										
public_square	0.03 0.00 0.33	0.25	0.61	0.00 0.00 0.00 0.00 0.28	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.00 0.14 0.06 0.14 0.36 0.00 0.08 0.08 0.00 0.03	0.00 0.14 0.06 0.14 0.36 0.00 0.08 0.08 0.00 0.03									
shopping_mall	0.00 0.00 0.08	0.81	0.00 0.00 0.00 0.00 0.05	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.06 0.00 0.00 0.00 0.06 0.87 0.11 0.11 0.00	0.06 0.00 0.00 0.00 0.06 0.87 0.11 0.11 0.00										
street_pedestrian	0.00 0.00 0.39	0.61	0.00 0.00 0.03 0.00 0.11	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.06 0.06 0.00 0.08 0.11 0.17 0.27 0.31 0.00 0.06	0.06 0.06 0.00 0.08 0.11 0.17 0.27 0.31 0.00 0.06										
street_traffic	0.00 0.00 0.03	0.10	0.06 0.00 0.00 0.00 0.03	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.00 0.02 0.00 0.00 0.06 0.07 0.20 0.03 0.72 0.00	0.00 0.02 0.00 0.00 0.06 0.07 0.20 0.03 0.72 0.00										
train	0.00 0.07 0.36	0.54	0.00 0.00 0.08 0.00 0.04	0.00	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	-0.36 0.03 0.00 0.00 0.03 0.17 0.28 0.11 0.08	0.00 0.28 0.14 0.25 0.09 0.04 0.07 0.04 0.00 0.18	0.00 0.28 0.14 0.25 0.09 0.04 0.07 0.04 0.00 0.18										

K. Drossos, P. Magron, T. Virtanen, "Unsupervised adversarial domain adaptation based on the Wasserstein distance for acoustic scene classification", Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA), October 2019.

Music recommendation

Content-aware music recommendation

Collaborative filtering with matrix factorization: $Y \approx WH$.

- ▷ $Y / W / H = \text{interactions} / \text{users preferences} / \text{songs attributes}$.
- ▷ Cannot recommend novel items: cold-start problem.

P. Magron, C. Févotte, "Leveraging the structure of musical preference in content-aware music recommendation", *Proc. IEEE International Conference on Audio, Speech and Signal Processing (ICASSP)*, June 2021.

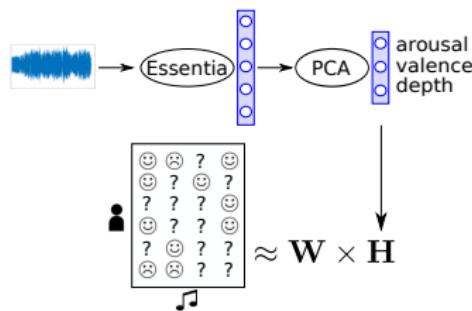
P. Magron, C. Févotte, "Neural content-aware collaborative filtering for cold-start music recommendation", to be submitted in the *ACM Transactions on Information Systems*.

Content-aware music recommendation

Collaborative filtering with matrix factorization: $Y \approx WH$.

- ▷ $Y / W / H = \text{interactions} / \text{users preferences} / \text{songs attributes}$.
- ▷ Cannot recommend novel items: cold-start problem.

Content-aware recommendation



P. Magron, C. Févotte, "Leveraging the structure of musical preference in content-aware music recommendation", *Proc. IEEE International Conference on Audio, Speech and Signal Processing (ICASSP)*, June 2021.

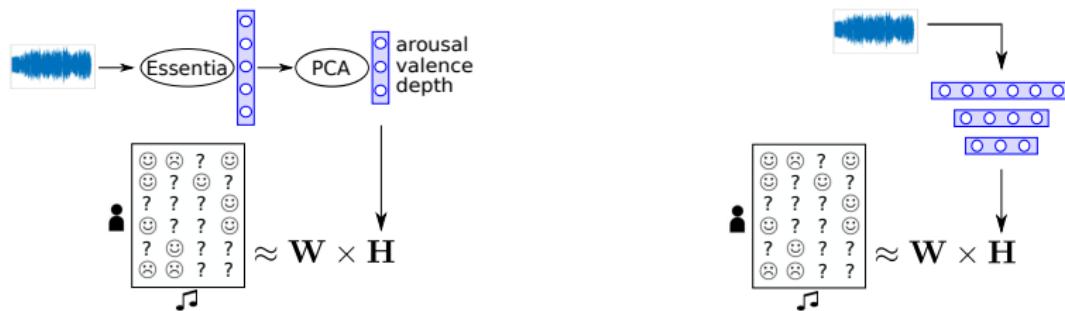
P. Magron, C. Févotte, "Neural content-aware collaborative filtering for cold-start music recommendation", to be submitted in the *ACM Transactions on Information Systems*.

Content-aware music recommendation

Collaborative filtering with matrix factorization: $Y \approx WH$.

- ▷ $Y / W / H = \text{interactions} / \text{users preferences} / \text{songs attributes}$.
- ▷ Cannot recommend novel items: cold-start problem.

Content-aware recommendation



P. Magron, C. Févotte, "Leveraging the structure of musical preference in content-aware music recommendation", *Proc. IEEE International Conference on Audio, Speech and Signal Processing (ICASSP)*, June 2021.

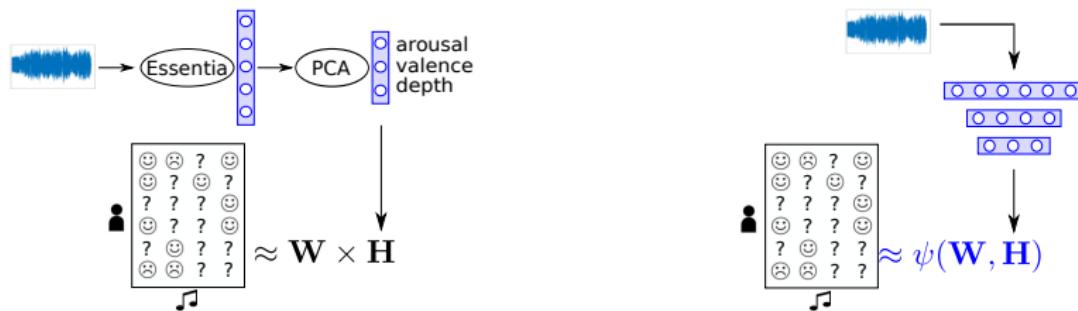
P. Magron, C. Févotte, "Neural content-aware collaborative filtering for cold-start music recommendation", to be submitted in the *ACM Transactions on Information Systems*.

Content-aware music recommendation

Collaborative filtering with matrix factorization: $Y \approx WH$.

- ▷ $Y / W / H = \text{interactions} / \text{users preferences} / \text{songs attributes}$.
- ▷ Cannot recommend novel items: cold-start problem.

Content-aware recommendation



P. Magron, C. Févotte, "Leveraging the structure of musical preference in content-aware music recommendation", *Proc. IEEE International Conference on Audio, Speech and Signal Processing (ICASSP)*, June 2021.

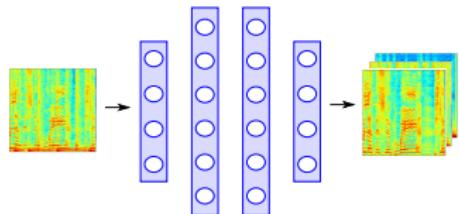
P. Magron, C. Févotte, "Neural content-aware collaborative filtering for cold-start music recommendation", to be submitted in the *ACM Transactions on Information Systems*.

Research project

Current context

For most audio processing tasks: (supervised) deep learning.

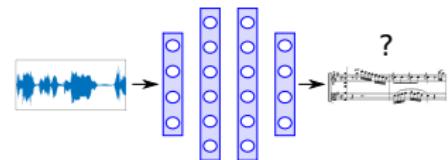
- ✓ Performances in controlled conditions.



Current context

For most audio processing tasks: (supervised) deep learning.

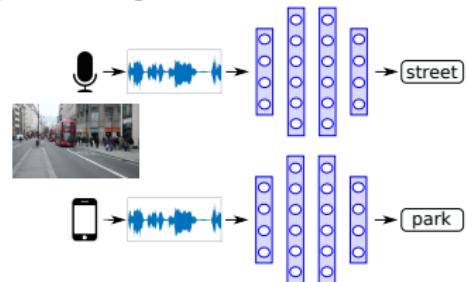
- ✓ Performances in controlled conditions.
- ✗ Cost of gathering annotated data.



Current context

For most audio processing tasks: (supervised) deep learning.

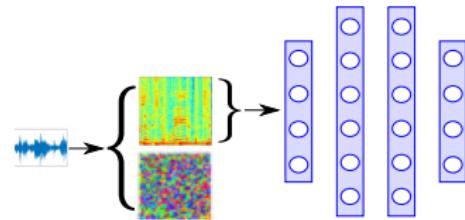
- ✓ Performances in controlled conditions.
- ✗ Cost of gathering annotated data.
- ✗ Adaptation to unseen conditions.



Current context

For most audio processing tasks: (supervised) deep learning.

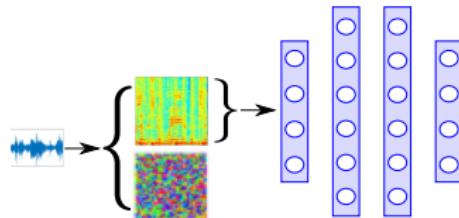
- ✓ Performances in controlled conditions.
- ✗ Cost of gathering annotated data.
- ✗ Adaptation to unseen conditions.
- ✗ Nonnegative representation processing.



Current context

For most audio processing tasks: (supervised) deep learning.

- ✓ Performances in controlled conditions.
- ✗ Cost of gathering annotated data.
- ✗ Adaptation to unseen conditions.
- ✗ Nonnegative representation processing.



My approach

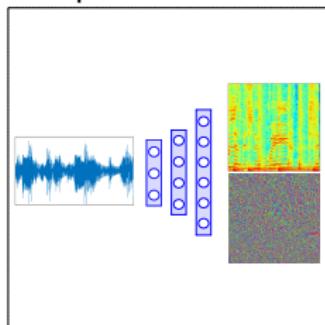
- ▷ Complex-valued representations processing.
- ▷ Interfacing factorization models and deep learning.
 - ▷ **Exhaustive** data processing.
 - ▷ **Light** and **interpretable** architectures.
 - ▷ **Reduced** supervision and **data-efficient** systems.

Research plan

Complex-valued and hybrid models for audio processing

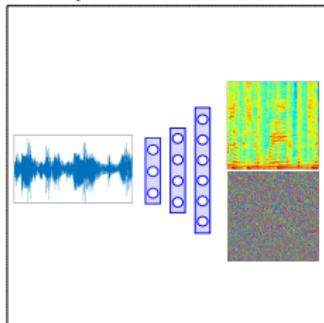
Complex-valued and hybrid models for audio processing

I. Complex-valued representations

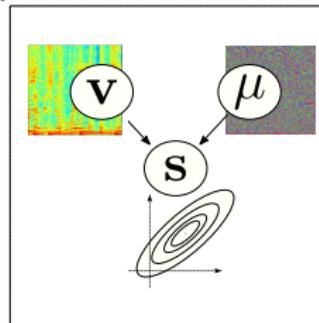


Complex-valued and hybrid models for audio processing

I. Complex-valued representations

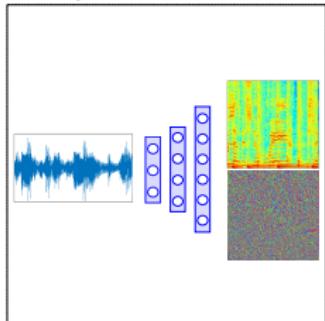


II. Anisotropic probabilistic modeling

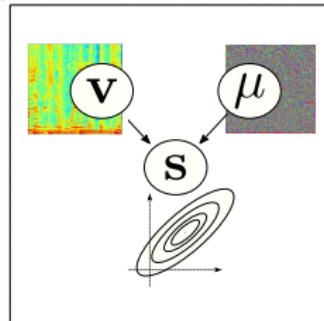


Complex-valued and hybrid models for audio processing

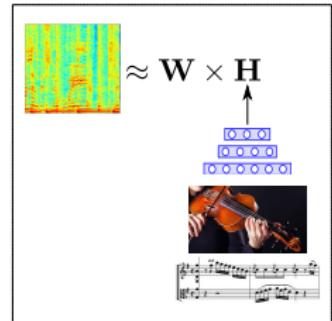
I. Complex-valued representations



II. Anisotropic probabilistic modeling



III. Hybrid factorization models



I. Complex-valued representations

Current trend: time-domain approaches.

- ▷ State-of-the-art results (source separation, speech enhancement).
- ▷ Potential of processing learned (complex-valued) mi-level representations.

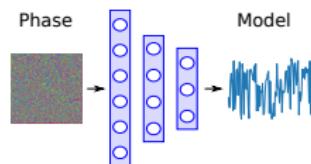
I. Complex-valued representations

Current trend: time-domain approaches.

- ▷ State-of-the-art results (source separation, speech enhancement).
- ▷ Potential of processing learned (complex-valued) mi-level representations.

Deep phase processing

- ▷ Learning and including deep phase models in complex-valued neural networks.



M. Pariente et al., "Filterbank design for end-to-end speech separation", *Proc. ICASSP*, April 2020.

D. Ditter and T. Gerkmann, "A multi-phase gammatone filterbank for speech separation Via Tasnet", *Proc. ICASSP*, April 2020.

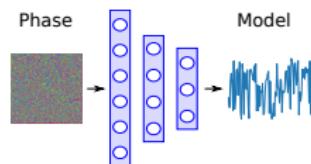
I. Complex-valued representations

Current trend: time-domain approaches.

- ▷ State-of-the-art results (source separation, speech enhancement).
- ▷ Potential of processing learned (complex-valued) mi-level representations.

Deep phase processing

- ▷ Learning and including deep phase models in complex-valued neural networks.



Algorithms for inverse problems

- ▷ Leveraging alternative divergences in optimization-based phase recovery.

M. Pariente et al., "Filterbank design for end-to-end speech separation", *Proc. ICASSP*, April 2020.

D. Ditter and T. Gerkmann, "A multi-phase gammatone filterbank for speech separation Via Tasnet", *Proc. ICASSP*, April 2020.

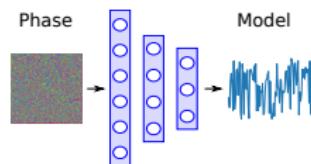
I. Complex-valued representations

Current trend: time-domain approaches.

- ▷ State-of-the-art results (source separation, speech enhancement).
- ▷ Potential of processing learned (complex-valued) mi-level representations.

Deep phase processing

- ▷ Learning and including deep phase models in complex-valued neural networks.

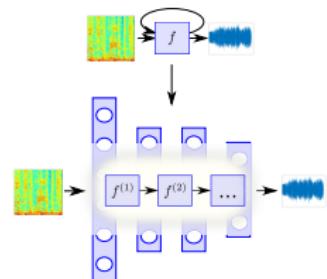


Algorithms for inverse problems

- ▷ Leveraging alternative divergences in optimization-based phase recovery.

Synthesis models

- ▷ Deep unfolding for time-domain synthesis.
- ▷ Generative time-frequency complex models.



M. Pariente et al., "Filterbank design for end-to-end speech separation", *Proc. ICASSP*, April 2020.

D. Ditter and T. Gerkmann, "A multi-phase gammatone filterbank for speech separation Via Tasnet", *Proc. ICASSP*, April 2020.

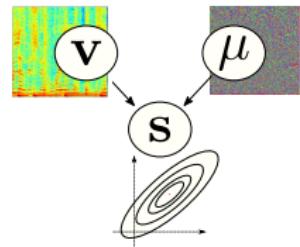
II. Anisotropic probabilistic modelling

II. Anisotropic probabilistic modelling

Bayesian anisotropic Gaussian models

$$s_{f,t} \sim \mathcal{N}_{\mathbb{C}}(m_{f,t}, \Gamma_{f,t})$$

- ▷ Model-based (e.g., DNN) structure for $m_{f,t}$.
- ▷ Bayesian prior for $\Gamma_{f,t} \rightarrow$ easy inference.



II. Anisotropic probabilistic modelling

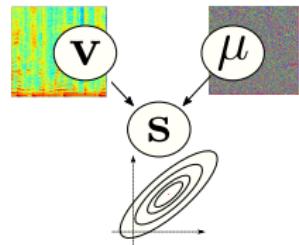
Bayesian anisotropic Gaussian models

$$s_{f,t} \sim \mathcal{N}_{\mathbb{C}}(m_{f,t}, \Gamma_{f,t})$$

- ▷ Model-based (e.g., DNN) structure for $m_{f,t}$.
- ▷ Bayesian prior for $\Gamma_{f,t} \rightarrow$ easy inference.

Extending anisotropy

- ▷ Alternative distributions (Tweedie, stable) are more appropriate for audio.

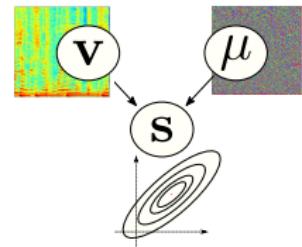


II. Anisotropic probabilistic modelling

Bayesian anisotropic Gaussian models

$$s_{f,t} \sim \mathcal{N}_{\mathbb{C}}(m_{f,t}, \Gamma_{f,t})$$

- ▷ Model-based (e.g., DNN) structure for $m_{f,t}$.
- ▷ Bayesian prior for $\Gamma_{f,t} \rightarrow$ easy inference.

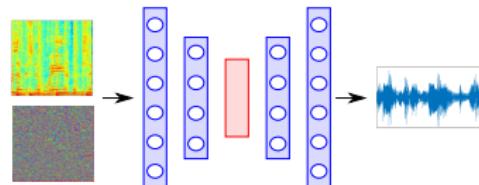


Extending anisotropy

- ▷ Alternative distributions (Tweedie, stable) are more appropriate for audio.

Deep generative models

- ▷ Anisotropic variational auto-encoders.
- ▷ Application to inter/multimodal synthesis.



III. Hybrid factorization models

NMF

- ▷ Worse performance than DNNs, additive assumption.
- ▷ But some potential: light, interpretable, unsupervised.

III. Hybrid factorization models

NMF

- ▷ Worse performance than DNNs, additive assumption.
- ▷ But some potential: light, interpretable, unsupervised.

Learn to factorize $\min_{W,H,\Phi} D(\Phi(x), WH)$

- ▷ Hand-craft, optimize, or deeply-learn a factorizable representation.

III. Hybrid factorization models

NMF

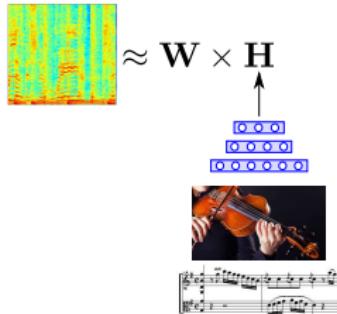
- ▷ Worse performance than DNNs, additive assumption.
- ▷ But some potential: light, interpretable, unsupervised.

Learn to factorize $\min_{W,H,\Phi} D(\Phi(x), WH)$

- ▷ Hand-craft, optimize, or deeply-learn a factorizable representation.

NMF/DNN interface

- ▷ NMF with deep priors \rightarrow self-supervised / multimodal feature learning.



III. Hybrid factorization models

NMF

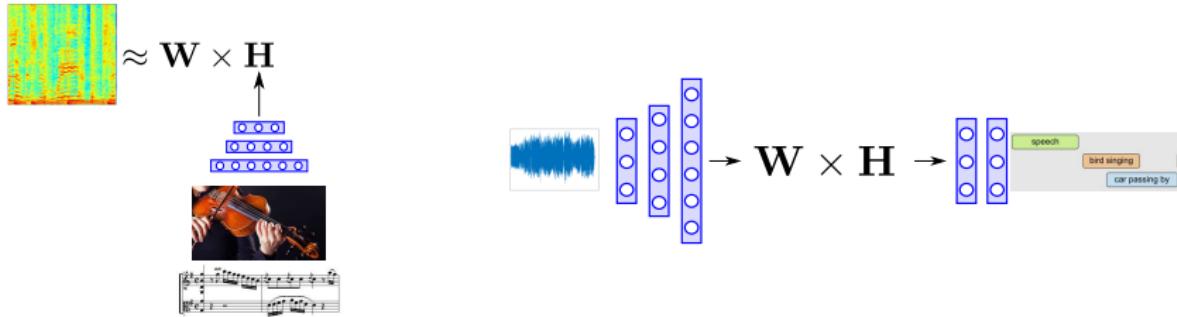
- ▷ Worse performance than DNNs, additive assumption.
- ▷ But some potential: light, interpretable, unsupervised.

Learn to factorize $\min_{W,H,\Phi} D(\Phi(x), WH)$

- ▷ Hand-craft, optimize, or deeply-learn a factorizable representation.

NMF/DNN interface

- ▷ NMF with deep priors → self-supervised / multimodal feature learning.
- ▷ Low-rank latent representations → light domain adaptation.



Thanks!

- ▷ Webpage: <https://magronp.github.io/>
- ▷ Code: <https://github.com/magronp>