Spiral scattering of nonstationary sounds



Vincent Lostanlen, Stéphane Mallat École normale supérieure de Paris

2015-06-19, Marseille, FR. Ce travail est financé par la bourse ERC InvariantClass.

- musical instrument recognition
- query by example
- automated transcription

rely on the characterization of musical notes.

NATURAL REGULARITY ALONG TIME

The physics of player-instrument interaction brings **regularity** in gestures:





GOAL

The TFR reveals short-term regularity (~50 ms).

Goal: to build a meaningful decomposition such that musical notes are regular (~500 ms).





TEMPLATE-BASED METHODS

- Hard constraints, e.g. Markovian.
 - see Kereliuk & Depalle (2008) on partial tracking.
 - useful for fine-grain audio effects...
 - ... but exposed to low-level detection errors.
- Loose constraints, e.g. Bayesian.
 - see Fuentes et al. (2013) on harmonic PLCA.
 - encompasses a broad range of meaningful priors...
 but timbral regularity remains challenging.

TEMPLATE-FREE APPROACHES

- In fact, most tasks do not require rigid templates.
- We advocate a progressive decomposition instead.
- ConvNets perform a data-driven decomposition, but they need a large annotated training set.
- We design a nonlinear scattering transform with
 - regularity
 - time-frequency localization
 - sparsity

in mind.

1. Nonstationary source-filter model

2. Leveraging harmonicity

3. Applications to classification and reconstruction

1. Nonstationary source-filter model



HARMONIC SOURCE

We define the harmonic source as a Dirac pulse train.



FILTER

We define the filter as a regular spectral envelope.



STATIONARY SOURCE-FILTER MODEL

The stationary source-filter model is x(t) = [e * h](t) i.e. $\hat{x}(\omega) = [\hat{e} \times \hat{h}](\omega)$



DEFORMATIONS

Let $\theta(t) \in C^3$ be a time warp function. $\dot{\theta}(t) > 0$ is the fundamental frequency of $e_{\theta}(t) = (e \circ \theta)(t)$.

 $\dot{\nu}(t) > 0$ is the position of the formant (spectral peak) $h_{\nu}(t) = (h \circ \nu)(t).$

The nonstationary source-filter model is defined as $x_{\theta,\nu}(t) = [e_\theta * h_\nu](t).$

EXAMPLE 1 : TROMBONE GLISSANDO



EXAMPLE 2 : TROMBONE CRESCENDO



A TIME-FREQUENCY PERSPECTIVE

We need a time-frequency representation that is 1. localized enough in time (a) to have $\dot{\theta}(t)$ approximately constant, (b) to have $\dot{\nu}(t)$ approximately constant,

2. localized enough in frequency (a) to distinguish the first peaks of $\hat{e}(\omega)$, (b) to have $\hat{h}(\omega)$ approximately constant.

We will use a constant-Q filter bank of wavelets.

ANALYTIC WAVELETS

Wavelets are oscillating, localized filters. By dilating a mother wavelet $\psi(t)$, we trade frequency resolution for time resolution.

 $\forall \lambda_1, \psi_{\lambda_1}(t) = \lambda_1 \psi(\lambda_1 t)$

Analyticity property: $\forall \omega \leq 0, \hat{\psi}(\omega) = 0$. Complex modulus improves regularity.

Wavelet spectrum is stable to deformations.

WAVELET FILTER BANK

In base 2:
$$\log \lambda_1 = j_1 + \frac{\chi_1}{Q}$$
 number of filters per octave

where the integer part j_1 is the octave index and $\chi_1 \in \{0 \dots (Q-1)\}$ is the chroma.



WAVELET RIDGE THEOREM

by Delprat, Escudié, Guillemain, Kronland-Martinet, Torrésani, and Tchamitchian (1992).

Let $f(t) = a(t) \cos \theta(t)$ and $\psi_{\lambda_1}(t) = \lambda_1 g(\lambda_1 t) \exp(i\lambda_1 t)$.

$$[f * \psi_{\lambda_1}](t)$$

= $a(t) \exp(i(\theta(t) - \lambda_1 t)) \times \left(\hat{g}\left(1 - \frac{\dot{\theta}(t)}{\lambda_1}\right) + \varepsilon(t, \lambda_1)\right)$

The corrective term $\varepsilon(t,\lambda_1)$ is small if

- amplitude modulation is slow: $\|\dot{a}/a\|_{\infty}^{2}, \|\ddot{a}/a\|_{\infty} \ll \lambda_{1}^{2}$,
- frequency modulation is slow: $\|\ddot{ heta}\|_{\infty} \ll \lambda_1^2$, and
- (t, λ_1) is near a ridge: $\dot{\theta}(t) \approx \lambda_1$.

FACTORIZATION IN THE SCALOGRAM

For Q between 12 and 24, we have $|x_{\theta,\nu} * \psi_{\lambda_1}| = \widehat{\psi}_{\lambda_1}(k\dot{\theta}(t)) \times \hat{h}\left(\frac{\lambda_1}{\dot{\nu}(t)}\right)$ where k is such that $\lambda_1 \approx k \theta(t)$. In log-frequency and after logarithmic compression: $\mathbf{U}_{\mathbf{1}} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}(t, \log \lambda_1) \stackrel{\text{def.}}{=} \log |x_{\boldsymbol{\theta}, \boldsymbol{\nu}} \ast \psi_{\lambda_1}|$ $= \mathbf{U}_{\mathbf{1}} \mathbf{e} (\log \lambda_1 - \log \dot{\mathbf{\theta}}(t)) + \mathbf{U}_{\mathbf{1}} \mathbf{h} (\log \lambda_1 - \log \dot{\mathbf{\nu}}(t))$ translated source translated filter

2. Leveraging harmonicity



Fig. from Shepard (1964).

HARMONICITY PROPERTY

The harmonic comb is self-similar: $\hat{e}(\omega) = \hat{e}(2^{j}\omega)$ for all $\omega > 1$ and $j \in \mathbb{N}$.

Regularity across octaves for a given chroma:

SPECTRAL SMOOTHNESS PROPERTY

The spectral envelope is regular across semitones:

Regularity along chromas within an octave:

PARTIAL DERIVATIVE ALONG TIME

By linearity and chain rule formula:

$$\frac{\partial \mathbf{U}_{\mathbf{1}} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}}{\partial t}(t, \log \lambda_{1}) = \frac{\ddot{\boldsymbol{\theta}}(t)}{\dot{\boldsymbol{\theta}}(t)} \frac{\mathrm{d} \mathbf{U}_{\mathbf{1}} \boldsymbol{e}}{\mathrm{d}(\log \lambda_{1})} (\log \lambda_{1} - \log \dot{\boldsymbol{\theta}}(t)) + \frac{\ddot{\boldsymbol{\nu}}(t)}{\dot{\boldsymbol{\nu}}(t)} \frac{\mathrm{d} \mathbf{U}_{\mathbf{1}} \boldsymbol{h}}{\mathrm{d}(\log \lambda_{1})} (\log \lambda_{1} - \log \dot{\boldsymbol{\nu}}(t)).$$

PARTIAL DERIVATIVE ALONG CHROMAS

By linearity:

$$\frac{\partial \mathbf{U}_{1} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}}{\partial (\log \lambda_{1})} (t, \log \lambda_{1}) = \frac{\mathrm{d} \mathbf{U}_{1} \boldsymbol{e}}{\mathrm{d} (\log \lambda_{1})} (\log \lambda_{1} - \log \dot{\boldsymbol{\theta}}(t)) + \frac{\mathrm{d} \mathbf{U}_{1} h}{\mathrm{d} (\log \lambda_{1})} (\log \lambda_{1} - \log \dot{\boldsymbol{\nu}}(t)).$$

neglected because of **spectral smoothness**

PARTIAL DERIVATIVE ACROSS OCTAVES

By linearity:

$$\frac{\Delta \mathbf{U}_{1} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}}{\Delta j_{1}} (t, \log \lambda_{1}) = \underbrace{\frac{\Delta \mathbf{U}_{1} \boldsymbol{e}}{\Delta j_{1}} (\log \lambda_{1} - \log \dot{\boldsymbol{\theta}}(t))}_{\Delta j_{1}}$$

neglected because of harmonicity

$$+\frac{\Delta \mathbf{U_1}h}{\Delta j_1}(\log \lambda_1 - \log \dot{\boldsymbol{\nu}}(t)).$$

OPTICAL FLOW EQUATION

 U_1x behaves like an object in rigid motion, hence an optical flow equation in t, $\log \lambda_1$, and j_1 .

$$\frac{\partial \mathbf{U}_{\mathbf{1}} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}}{\partial t} (t, \log \lambda_{1}) = \frac{\ddot{\boldsymbol{\theta}}(t)}{\dot{\boldsymbol{\theta}}(t)} \frac{\partial \mathbf{U}_{\mathbf{1}} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}}{\partial (\log \lambda_{1})} (t, \log \lambda_{1}) + \frac{\ddot{\boldsymbol{\nu}}(t)}{\dot{\boldsymbol{\nu}}(t)} \frac{\Delta \mathbf{U}_{\mathbf{1}} x_{\boldsymbol{\theta}, \boldsymbol{\nu}}}{\Delta j_{1}} (t, \log \lambda_{1}).$$

 $- \int Motion in (\log \lambda_1, j_1)$ is best expressed on a **spiral**: the chroma is angular, the octave is radial.

Shepard-Risset Helix

ROTATING MOTION IN THE SPIRAL A chirp has a rotating motion in the spiral.

RADIAL MOTION IN THE SPIRAL A formantic change has a radial motion in the spiral.

Two degrees of freedom

Musical transients are not regular in time-frequency... ... but in **time-chroma-octave**.

3. Spiral scattering

SPIRAL WAVELETS

Accepted at GRETSI 2015 with Stéphane Mallat: *Transformée de scattering en spirale temps-chroma-octave*.

 $\Psi_{(\alpha,\beta,\gamma)}(t) = \psi_{\alpha}(t) \times \psi_{\beta}(\log \lambda_1) \times \psi_{\gamma}(\lfloor \log \lambda_1 \rfloor)$

FREQUENCIES VS. QUEFRENCIES

 $\Psi_{(\alpha,\beta,\gamma)}(t) = \psi_{\alpha}(t) \times \psi_{\beta}(\log \lambda_1) \times \psi_{\gamma}(\lfloor \log \lambda_1 \rfloor).$

- α is a modulation frequency along time, in Hertz. α^{-1} is typically between 1 ms and 100 ms.
- β is a « quefrency » along chromas, in cycles per octaves. $|\beta^{-1}|$ is typically between 1 semitone and 1 octave.
- γ is a quefrency across octaves, in cycles per octaves. $|\gamma^{-1}|$ is typically between 1 and 4 octaves.

We define the multiindex frequency variable $\lambda_2 = (\alpha, \beta, \gamma)$. By convention, $\log \lambda_2 = (\log \alpha, \log \beta, \operatorname{sign} \beta, \log \gamma, \operatorname{sign} \gamma)$.

Scattering cascade

Wavelet filter banks *scatter* the energy from U_m to U_{m+1} . Complex modulus improves regularity and phase invariance.

$$\mathbf{U}_{1}x(t, \log \lambda_{1}) = |x \stackrel{t}{*} \psi_{\lambda_{1}}|(t)$$
$$\mathbf{U}_{2}x(t, \log \lambda_{1}, \log \lambda_{2}) = |\mathbf{U}_{1}x \circledast \Psi_{\lambda_{2}}|(t, \log \lambda_{1})$$
$$\mathbf{U}_{2}x(t, \log \lambda_{1}, \log \lambda_{2}) = |\mathbf{U}_{1}x \circledast \Psi_{\lambda_{2}}|(t, \log \lambda_{1})$$
$$\mathbf{U}_{2}x(t, \log \lambda_{1}, \log \lambda_{2}) = |\mathbf{U}_{1}x \circledast \Psi_{\lambda_{2}}|(t, \log \lambda_{1})$$

A lowpass filter $\phi(t)$ enforces the amount of translation invariance that is required by the classification task.

$$\mathbf{S_1} x(t, \log \lambda_1) = \mathbf{U_1} x \overset{t}{*} \phi$$
$$\mathbf{S_2} x(t, \log \lambda_1) = \mathbf{U_2} x \overset{t}{*} \phi$$

Source-filter properties

- Vanishing moment property: Convolving a wavelet with a linear function yields almost zero.
- Harmonicity and spectral smoothness rewrite as

$$\left| \mathbf{U}_{\mathbf{1}} e_{\theta} \overset{j_{1}}{*} \psi_{\gamma} \right| = 0 \text{ and } \left| \mathbf{U}_{\mathbf{1}} h_{\nu} \overset{\chi_{1}}{*} \psi_{\beta} \right| \approx 0.$$

• The spiral scattering transform boils down to

$$\begin{aligned} \mathbf{U}_{1} x_{\theta,\nu} & \overset{t,\chi_{1},j_{1}}{\circledast} \Psi_{\lambda_{2}} \\ &= \left[\left(\mathbf{U}_{1} e_{\theta} \overset{\chi_{1}}{\ast} \psi_{\beta} \right) \times \left(\mathbf{U}_{1} h_{\nu} \overset{j_{1}}{\ast} \psi_{\gamma} \right) \right] \overset{t}{\ast} \psi_{\alpha} \end{aligned}$$

SPIRAL WAVELET RIDGES

• Applying the wavelet ridge theorem three times yields:

$$\begin{aligned} \mathbf{U}_{\mathbf{1}} x_{\theta,\nu} & \stackrel{t,\chi_{1},j_{1}}{\circledast} \Psi_{\lambda_{2}} \\ &= \left| \mathbf{U}_{\mathbf{1}} e_{\theta} \stackrel{\chi_{1}}{\ast} \psi_{\beta} \right| \left| \mathbf{U}_{\mathbf{1}} h_{\nu} \stackrel{j_{1}}{\ast} \psi_{\gamma} \right| \left| \hat{\psi}_{\alpha} \left(-\frac{\ddot{\theta}(t)}{\dot{\theta}(t)} \beta - \frac{\ddot{\nu}(t)}{\dot{\nu}(t)} \gamma \right) \end{aligned}$$

• Ridges are on a plane whose Cartesian equation is

$$\alpha + \frac{\ddot{\theta}(t)}{\dot{\theta}(t)}\beta + \frac{\ddot{\nu}(t)}{\dot{\nu}(t)}\gamma = 0.$$

- The same holds for averaged coefficients ${f S_2} x_{{m heta},{m
u}}$ over T if

$$\frac{\ddot{\theta}(t)}{\ddot{\theta}(t)} - \frac{\ddot{\theta}(t)}{\dot{\theta}(t)} \bigg| \ll T^{-1} \text{ and } \left| \frac{\dddot{\nu}(t)}{\dddot{\nu}(t)} - \frac{\dddot{\nu}(t)}{\dot{\nu}(t)} \right| \ll T^{-1}.$$

SPATIAL LOCALIZATION OF RIDGES

Below: $U_1 x$ of the word « lion ».

Opposite: $\mathbf{U_2} x$ slices for

- $\alpha^{-1} = 120 \text{ ms}$
- $|\beta^{-1}| = \pm 1$ octave
- $|\gamma^{-1}| = \pm 1$ octave

(α, β, γ) localization of Ridges

- Top: $\mathbf{U_1} x$ of a trombone note.
- (a) **attack** part with upwards glissando
- (b) **release** part with downwards glissando.

- Bottom: $\mathbf{U_2}x$ slices for
- $\alpha^{-1} = 46 \text{ ms}$
- fixed t and $\log \lambda_1$.

PHONEME CLASSIFICATION

Submitted to MLSP 2015 with Joakim Andén and Stéphane Mallat: Joint Time-frequency Scattering for Audio Classification.

Phoneme error rate on the TIMIT dataset [Fisher et al. 1986].

MFCC and SVM	18,3 %
MFCC and GMM commitee [Chang & Glass, 2007]	16.7 %
α scattering and SVM	17,3 %
(α, β) scattering and SVM	15,8 %
(α, β, γ) scattering	coming next

INVARIANT RECONSTRUCTION

Problem:

« find a translation-invariant representation that allows the most plausible signal reconstruction ».

Given a target $\mathbf{S}^{\infty} = (\mathbf{S}_{1}^{\infty}, \mathbf{S}_{2}^{\infty})$, the gradient descent $\Delta \mathbf{U}_{1}x = \Delta \mathbf{S}_{1}x \overset{t}{*} \overline{\phi}$ $+ \Re \left[\sum_{\lambda_{2}} \left(\frac{\mathbf{U}_{1}x \circledast \Psi_{\lambda_{2}}}{|\mathbf{U}_{1}x \circledast \Psi_{\lambda_{2}}|} \times \Delta \mathbf{U}_{2}x \right) \circledast \overline{\Psi}_{\lambda_{2}} \right].$

converges to a local minimum of the loss function

$$E(x) = \|\mathbf{S}_1 x - \mathbf{S}_1^{\infty}\|_2 + \|\mathbf{S}_2 x - \mathbf{S}_2^{\infty}\|_2.$$

original

first-order only

McDermott & Simoncelli

time scattering

time-frequency scattering

spiral scattering

CONCLUSIONS

- Natural sounds are nonstationary, but physically regular.
- In the pitch spiral, source-filter transients become translations.
- Spiral scattering yields source-filter velocities without detection.
- Encouraging results in classification and invariant reconstruction.

Experiments can be reproduced at: <u>www.github.com/lostanlen/</u>

