

WAVELET SCATTERING ON THE PITCH SPIRAL

Vincent Lostanlen, Stéphane Mallat*

Dept. of Computer Science,
École normale supérieure
Paris, France
vincent.lostanlen@ens.fr

ABSTRACT

We present a new representation of sounds that linearizes the dynamics of pitch chroma and pitch height, while remaining stable to deformations in the time-frequency plane. It is an instance of the scattering transform, a generic operator which cascades wavelet convolutions and modulus nonlinearities. It is derived from the Shepard pitch spiral, in that convolutions are performed in time, log-frequency (correlated to pitch chroma) and octave index (correlated to pitch height).

1. INTRODUCTION

The spectro-temporal evolution of harmonic spectra conveys essential information to audio classification, blind source separation, automated transcription, as well as other processing tasks. This information is however difficult to capture in time-varying, polyphonic mixtures. On one hand, spectrogram-based pattern recognition algorithms are exposed to detection errors as they impose strong constraints on the shape of harmonic templates. On the other, time-varying generalizations of matrix factorization [10] are under-constrained and thus may fail to converge to a satisfying solution. In this article, we address the problem of characterizing harmonic structures without resorting to any detection nor learning step.

First, local continuity in frequency modulation cause a slow rigid motion along the log-frequency axis. Second, harmonic sounds exhibit a comb-like spectrum, which conveys the global evolution of the spectral envelope. This comb looks highly irregular: unlike frequency modulation, it cannot be captured efficiently with local convolutions in time and log-frequency.

To recover regularity across partials, we capitalize on the fact that power-of-two harmonics are distant from exactly one octave. By rolling up the log-frequency axis in a spiral, such that octave intervals correspond to full turns, these partials get aligned on a radius. Consequently, introducing the integer-valued octave variable reveals harmonic regularity that was not explicit in the plane of time and log-frequency. Once specified the variables of time, log-frequency, and octave index, our transform merely consists in cascading three wavelet decompositions along them and applying complex modulus.

Section 2 gives a formal definition of the spiral scattering transform. Section 3 introduces a nonstationary formulation of the source-filter model relying on time warps, and shows that its variabilities in pitch and spectral envelope are jointly linearized by the spiral scattering transform. Section 4 provides a visual inter-

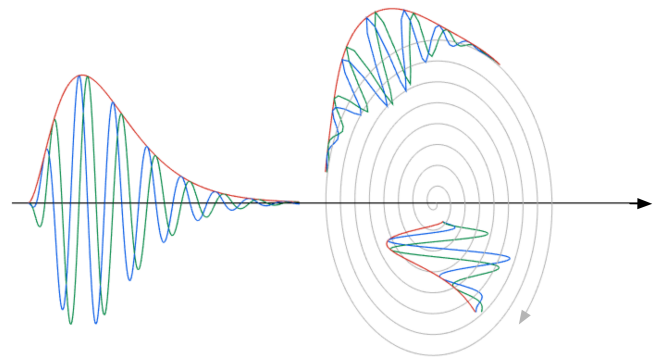


Figure 1: The spiral wavelet is a product of wavelets along time t , log-frequency $\log \lambda_1$ (angular variable), and octave index j (radial variable). Blue and green oscillations represent the real and imaginary parts. The red envelope represents the complex modulus. Partial of an hypothetical harmonic sound are marked as thick dots.

pretation of the spiral wavelet coefficients in a music signal with extended instrumental techniques.

2. FROM TIME SCATTERING TO SPIRAL SCATTERING

This section builds the spiral scattering transform progressively as a cascade of wavelet transforms along time, log-frequency, and octave index.

2.1. Time scattering

An analytic "mother" wavelet is a complex filter $\psi(t)$ whose Fourier transform $\hat{\psi}(\omega)$ is concentrated over the dimensionless frequency interval $[1 - 2^{1/2Q}, 1 + 2^{1/2Q}]$, where the quality factor Q is in the typical range 12–24. Dilations of this wavelet define a family of bandpass filters centered at frequencies $\lambda_1 = 2^{j_1 + \frac{\chi_1}{Q}}$, where the indices $j_1 \in \mathbb{Z}$ and $\chi_1 \in \{1 \dots Q\}$ respectively denote octave and chroma:

$$\hat{\psi}_{\lambda_1}(\omega) = \hat{\psi}(\lambda_1^{-1}\omega) \quad \text{i.e.} \quad \psi_{\lambda_1}(t) = \lambda_1 \psi(\lambda_1 t). \quad (1)$$

The wavelet transform convolves an input signal $x(t)$ with the filter bank of ψ_{λ_1} 's. Its modulus is the wavelet "scalogram"

* This work is supported by the ERC InvariantClass 320959.

$$x_1(t, \log \lambda_1) = |x * \psi_{\lambda_1}| \quad \text{for all } \lambda_1 > 0, \quad (2)$$

whose frequential axis is uniformly sampled by the binary logarithm $\log \lambda_1$. The scalogram x_1 localizes the energy of $x(t)$ around the frequencies λ_1 over durations $2Q\lambda_1^{-1}$, trading frequency resolution for time resolution.

The constant-Q transform (CQT) S_1x corresponds to a low-pass filtering of x_1 with a window $\phi(t)$ of size T :

$$S_1x(t, \log \lambda_1) = x_1 * \phi_T = |x * \psi_{\lambda_1}| * \phi_T. \quad (3)$$

To recover the amplitude modulations lost when averaging by ϕ_T in Equation (3), the time scattering transform also convolves x_1 with a second filterbank of wavelets ψ_{λ_2} and applies complex modulus to get

$$x_2(t, \log \lambda_1, \log \lambda_2) = |x * \psi_{\lambda_2}| = ||x * \psi_{\lambda_1}| * \psi_{\lambda_2}|. \quad (4)$$

The wavelets $\psi_{\lambda_2}(t)$ have a quality factor in the range 1–2, though we choose to keep the same notation ψ for simplicity. Like in Equation (3), averaging in time creates invariance to translation in time up to T , yielding

$$S_2x(t, \log \lambda_1, \log \lambda_2) = x_2 * \phi = ||x * \psi_{\lambda_1}| * \psi_{\lambda_2}| * \phi_T. \quad (5)$$

Due to the constant-Q property, S_1x and S_2x are stable to small time warps of $x(t)$ as long as they do not exceed Q^{-1} , i.e. one semitone. This implies that small modulations, such as tremolo and vibrato, are accurately linearized [1].

2.2. Joint time-frequency scattering

The time scattering transform defined in Equation (4) decomposes each frequency band separately, and so cannot properly capture the coherence of time-frequency structures, such as those induced by pitch contour. To remedy this, Andén [2] has redefined the wavelets ψ_{λ_2} 's as functions of both time and log-frequency, indexed by pairs $\lambda_2 = (\alpha, \beta)$, where α is measured in Hertz and β is measured in cycles per octaves. The joint wavelets $\psi_{\lambda_2}(t, \log \lambda_1)$ factorize as

$$\psi_{\lambda_2}(t, \log \lambda_1) = \psi_{\alpha}(t) \times \psi_{\beta}(\log \lambda_1). \quad (6)$$

We write $\overset{\chi_1}{*}$ to denote convolutions along the log-frequency axis, i.e. along chromas. Wavelet scattering is extended to two-dimensional convolutions by plugging Equation (6) into the definition of x_2 in Equation (4).

$$x_2(t, \log \lambda_1, \log \lambda_2) = |x_1 \overset{t, \chi_1}{*} \psi_{\lambda_2}| = |x_1 * \psi_{\alpha} \overset{\chi_1}{*} \psi_{\beta}|. \quad (7)$$

The joint time-frequency scattering transform corresponds to the "cortical transform" introduced by Shamma to formalize his findings in auditory neuroscience.

2.3. Spiral scattering

The time-frequency scattering transform defined in Equation (??) provides template-free features for pitch variability along time. However, it is unaware of the harmonic structure of voiced sounds, such as vowels or musical notes. The temporal evolution of this

structure yields relevant information about attack transients and formantic changes, almost independently from the pitch contour.

In order to disentangle variabilities in pitch and spectral envelope, we extend the joint time-frequency scattering transform to encompass regularity in time across octaves, in conjunction with regularity along neighboring constant-Q bands.

We roll up the log-frequency variable $\log \lambda_1$ into a pitch spiral making one full turn at each octave (see Fig. 1). Since a frequency interval of one octave corresponds to one unit in binary logarithms $\log \lambda_1$, pitch chroma and pitch height in the spiral correspond to integer part $\lfloor \log \lambda_1 \rfloor$ and fractional part $\{\log \lambda_1\}$:

$$\log \lambda_1 = \lfloor \log \lambda_1 \rfloor + \{\log \lambda_1\}. \quad (8)$$

In this setting, the fundamental frequency f_0 is aligned with its power-of-two harmonics $2f_0, 4f_0, 8f_0$ and so forth. Likewise, the perfect fifth $3f_0$ is aligned with $6f_0$. As the number of harmonics per octave increase exponentially, the alignment of upper harmonics — $5f_0, 7f_0$, and so forth — in the spiral is less crucial, because it can also be recovered with chroma scattering for β^{-1} of the order of a few semitones.

We cascade three one-dimensional wavelet transforms in time, log-frequency, and octave index, to build a so-called spiral scattering transform:

$$\psi_{\lambda_2}(t, \log \lambda_1) = \psi_{\alpha}(t) \times \psi_{\beta}(\log_2 \lambda_1) \times \psi_{\gamma}(\lfloor \log \lambda_1 \rfloor) \quad (9)$$

To ensure invertibility and energy conservation, the quefren- cies β and γ must take negative values, including zero. We adopt the shorthand notation

$$\log \lambda = (\log \alpha, \log |\beta|, \text{sign } \beta, \log |\gamma|, \text{sign } \gamma) \quad (10)$$

In the special case $\beta = 0$, ψ_{β} is no longer a wavelet but a low-pass filter whose support covers one octave. By convention, the corresponding log-frequency index is $\log |\beta| = -\infty$. The same remark applies to ψ_{γ} for $\gamma = 0$, which covers six octaves.

Since its Fourier transform $\widehat{\psi_{\lambda_2}}$ is centered at (α, β, γ) , the spiral wavelet ψ_{λ_2} has a pitch chroma velocity of α/β and a pitch height velocity of α/γ . Both velocities are measures in octaves per second.

We write $\overset{j}{*}$ to denote convolutions across neighboring octaves. The definition for x_2 is the same as Equation 7:

$$\begin{aligned} x_2(t, \log_2 \lambda_1, \log_2 \lambda_2) &= |x_1 \overset{t, \chi_1, j}{*} \psi_{\lambda_2}| \\ &= |x_1 * \psi_{\alpha} \overset{\chi_1}{*} \psi_{\beta} \overset{j}{*} \psi_{\gamma}|. \end{aligned} \quad (11)$$

Rolling up pitches into a spiral is a well-established idea in music, if only because of circularity of musical pitch classes. It has been studied by Shepard [3], Risset [4], and Deutsch [5] to build paradoxes in perception of pitch, and is corroborated by functional imaging of the auditory cortex [6].

3. DEFORMATIONS OF THE SOURCE-FILTER MODEL

A classical model for voiced speech production consists in the convolution of a harmonic glottal source $e(t)$ with a vocal tract filter $h(t)$. Introducing independent deformations to both components brings realistic variability to pitch and spectral envelope. This section studies the decomposition of the deformed source-filter model in the spiral scattering transform.

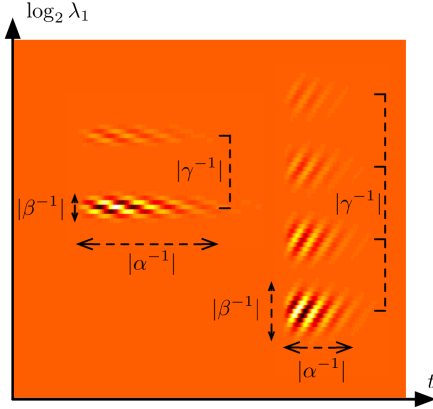


Figure 2: Two spiral wavelets $\psi_{\lambda_2}(t, \log_2 \lambda_1, \lfloor \log_2 \lambda_1 \rfloor)$ in the time-frequency plane, with different values of $\lambda_2 = (\alpha, \beta, \gamma)$. Left: $\alpha^{-1} = 120\text{ms}$, $\beta^{-1} = -0.25$ octave, $\gamma^{-1} = +2$ octaves. Right: $\alpha^{-1} = 60\text{ ms}$, $\beta^{-1} = +0.5$ octave, $\gamma^{-1} = -4$ octaves. Darker color levels corresponds to greater values of the real part.

3.1. Overview

Let $e(t) = \sum_n \delta(t - 2\pi f_0^{-1}n)$ be a harmonic signal and $t \mapsto \theta(t)$ a time warp function. We define a warped source as $e_\theta(t) = (e \circ \theta)(t)$. Similarly, we compose a filter $h(t)$ and a warp $t \mapsto \nu(t)$ to define $h_\nu(t) = (h \circ \nu)(t)$. The warped source-filter model is

$$x(t) = [e_\theta * h_\nu](t). \quad (12)$$

Observe that $\dot{\theta}(t)$ induces a change of fundamental frequency, whereas $\dot{\nu}(t)$ accounts for a local dilation of the spectral envelope $|\hat{h}|(\omega)$. We show in this section that, for $\dot{\theta}(t)$ and $\dot{\nu}(t)$ reasonably regular over the support of first-order wavelets, the local maxima of x_2 are clustered on a plane in the (α, β, γ) space of scattering coefficients. This plane satisfies the Cartesian equation

$$\alpha + \frac{\ddot{\theta}(t)}{\dot{\theta}(t)}\beta + \frac{\ddot{\nu}(t)}{\dot{\nu}(t)}\gamma = 0. \quad (13)$$

This result is likely to help automated transcription of polyphonic music, since notes overlapping both in time and frequency could be disentangled according to their respective source-filter velocities. Let $e_{\theta,1}(t, \log_2 \lambda_1)$ and $h_{\nu,1}(t, \log_2 \lambda_1)$ be the respective scalograms of $e_\theta(t)$ and $h_\nu(t)$.

3.2. Source-filter factorization in the scalogram

Given λ_1 near $pf_0\dot{\theta}(t)$ where $p \in \mathbb{N}$, the first step is to linearize $\theta(t)$ and $\nu(t)$ over the support of the first-order wavelet $\psi_{\lambda_1}(t)$. We work with the following assumptions:

- (a) Q large enough to discriminate the p^{th} partial: $Q > 2p$,
- (b) slowly varying source: $\|\ddot{\theta}/\dot{\theta}\|_\infty \ll \lambda_1/Q$, and
- (c) slowly varying filter: $\|\ddot{\nu}/\dot{\nu}\|_\infty \ll \lambda_1/Q$.

According to (a), partials $p' \neq p$ have a negligible contribution to e_1 at the log-frequency $\log_2 \lambda_1$. For lack of any interference, e_1 is constant through time, and we may drop the dependency in t :

$$e_1(\log_2 \lambda_1) = |\widehat{\psi_{\lambda_1}}(pf_0)| \quad (14)$$

According to (b), the scalogram of the deformed source can be replaced by the scalogram of the original source translated along the log-frequency at the velocity $\log_2 \dot{\theta}(t)$:

$$e_{\theta,1}(t, \log_2 \lambda_1) = e_1(\log_2 \lambda_1 - \log_2 \dot{\theta}(t)). \quad (15)$$

Similarly, we leverage (c) to linearize $\nu(t)$ over the support of $\psi_{\lambda_1}(t)$. The spectral smoothness assumption allows to approximate $\hat{h}(\omega)$ by a constant over the frequential support of the wavelet, hence to factorize the filtering as a product:

$$[h_\nu * \psi_{\lambda_1}] = h_1(\log_2 \lambda_1 - \log_2 \dot{\nu}(t))\psi_{\lambda_1}\left(\frac{\nu(t)}{\dot{\nu}(t)}\right). \quad (16)$$

By plugging Equation 15 into Equation 16, x_1 appears as a separable product between e_1 and h_1 , moving in log-frequency at respective velocities $\log_2 \dot{\theta}(t)$ and $\log_2 \dot{\nu}(t)$:

$$x_1(t, \log_2 \lambda_1) = e_1(\log_2 \lambda_1 - \log_2 \dot{\theta}(t))h_1(\log_2 \lambda_1 - \log_2 \dot{\nu}(t)). \quad (17)$$

3.3. Harmonicity and spectral smoothness properties

The second step in the proof consists in showing that the convolution along chromas with ψ_β only applies to $e_{1,\theta}$, whereas the convolution across octaves with ψ_γ only applies to $h_{1,\nu}$. Indeed, all wavelets are designed to carry a negligible mean value, i.e. convolving them with a constant yields zero. Therefore, the harmonicity and spectral smoothness properties rewrite as

$$e_{\theta,1} \overset{j}{*} \psi_\gamma \approx 0 \quad \text{and} \quad h_{\nu,1} \overset{\chi_1}{*} \psi_\beta \approx 0. \quad (18)$$

Gathering Equations 17 and 18 into the definition of spiral scattering yields

$$x_1 \overset{t, \chi_1, j}{*} \psi_{\lambda_2} = \left[\left(e_{1,\theta} \overset{\chi_1}{*} \psi_\beta \right) \times \left(h_{1,\nu} \overset{j}{*} \psi_\gamma \right) \right] \overset{t}{*} \psi_\alpha, \quad (19)$$

where the superscripts t , χ_1 and j denote convolutions along time, chromas and octaves respectively.

3.4. Extraction of instantaneous frequencies

As a final step, we state that the phase of $[e_{\theta,1} \overset{\chi}{*} \psi_\beta]$ is $\beta \times (\log_2 \lambda_1 - \log_2 p\dot{\theta}(t))$. By differentiating this quantity along t for fixed $\log_2 \lambda_1$, we obtain an instantaneous frequency of $-\beta\ddot{\theta}(t)/\dot{\theta}(t)$. Similarly, the instantaneous frequency of the convolution $[h_{\nu,1} \overset{j}{*} \psi_\gamma]$ is $-\gamma\ddot{\nu}(t)/\dot{\nu}(t)$. As long as

$$\alpha \geq \left| \frac{\ddot{\theta}(t)}{\dot{\theta}(t)}\beta \right| \quad \text{and} \quad \alpha \geq \left| \frac{\ddot{\nu}(t)}{\dot{\nu}(t)}\gamma \right|, \quad (20)$$

the envelopes of these two convolutions are almost constant over the support of $\psi_\alpha(t)$. We conclude with the following approximate closed-form expression for the spiral scattering coefficients of the deformed source-filter model:

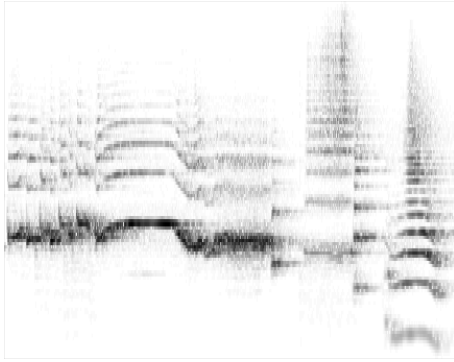


Figure 3

$$x_2(t, \log_2 \lambda_1, \log_2 \lambda_2) = |e_{\theta,1} \overset{x}{*} \psi_\beta| \times |h_{\nu,1} \overset{j}{*} \psi_\gamma| \times \left| \widehat{\psi}_\alpha \left(-\frac{\ddot{\theta}(t)}{\dot{\theta}(t)}\beta - \frac{\ddot{\nu}(t)}{\dot{\nu}(t)}\gamma \right) \right|. \quad (21)$$

The Fourier spectrum $|\widehat{\psi}_\alpha(\omega)|$ of $\psi_\alpha(t)$ is a bump centered at the frequency α . Equation 13 follows immediately from the above formula. The same result holds for the averaged coefficients $S_2x = x_2 * \phi_T(t)$ if the velocities $\ddot{\theta}(t)/\dot{\theta}(t)$ or $\ddot{\nu}(t)/\dot{\nu}(t)$ have small relative variations:

$$\left| \frac{\ddot{\theta}(t)}{\dot{\theta}(t)} - \frac{\ddot{\theta}(t)}{\dot{\theta}(t)} \right| \ll T^{-1} \quad \text{and} \quad \left| \frac{\ddot{\nu}(t)}{\dot{\nu}(t)} - \frac{\ddot{\nu}(t)}{\dot{\nu}(t)} \right| \ll T^{-1}. \quad (22)$$

An important caveat is that the inequalities above do not hold at inflexion points of the diffeomorphisms $\theta(t)$ and $\nu(t)$, i.e. where the velocities $\dot{\theta}(t)/\dot{\theta}(t)$ or $\dot{\nu}(t)/\dot{\nu}(t)$ cross zero.

4. CONCLUSIONS

The spiral model we have presented is well-known in music theory and experimental psychology. However, existing methods in audio signal processing do not fully take advantage from its richness, as they either picture pitch on a line (e.g. MFCC) or on a circle (e.g. chroma features).

Future work will be devoted to evaluating the discriminative power of Shepard spiral scattering coefficients over a variety of classification pipelines. Our representation also encompass automatic music transcription, perceptual similarity learning, and new audio transformations as potential applications.

5. REFERENCES

- [1] Joakim Andén and Stéphane Mallat, “Scattering Representation of Modulated Sounds,” *Proc. of the 15th Int. Conference on Digital Audio Effects (DAFx-12)*, , no. 3, pp. 15–18, 2012.
- [2] Joakim Andén, *Time and Frequency Scattering for Audio Classification*, Ph.D. thesis, École Polytechnique, 2014.
- [3] R. Shepard, “Circularity in Judgments of Relative Pitch,” *J. Acoust. Soc. Am.*, vol. 36, no. 12, pp. 2346, 1964.

- [4] Jean-Claude Risset, “Pitch control and pitch paradoxes demonstrated with computer-synthesized sounds,” *J. Acoust. Soc. Am.*, vol. 46, no. 1A, pp. 88–88, 1969.
- [5] Diana Deutsch, Kevin Dooley, and Trevor Henthorn, “Pitch circularity from tones comprising full harmonic series,” *The Journal of the Acoustical Society of America*, vol. 124, no. 1, pp. 589–597, 2008.
- [6] Jason D. Warren, Stefan Uppenkamp, Roy D. Patterson, and Timothy D. Griffiths, “Analyzing Pitch Chroma and Pitch Height in the Human Brain,” *Annals of the New York Academy of Sciences*, vol. 999, no. 17, pp. 212–214, 2003.
- [7] S. Abdallah and M. Plumbley, “Polyphonic Music Transcription By Non-negative Sparse Coding of Power Spectra,” in *Proc. ISMIR*, 2004, vol. 510, pp. 10–14.
- [8] P. Smaragdis and J.C. Brown, “Non-negative Matrix Factorization For Polyphonic Music Transcription,” in *2003 IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (IEEE Cat. No.03TH8684)*, 2003.
- [9] J. Andén and S. Mallat, “Deep Scattering Spectrum,” *IEEE Transactions on Signal Processing*, vol. 62, no. 16, pp. 4114–4128, 2014.
- [10] R. Hennequin, R. Badeau, and B. David, “NMF with Time-frequency Activations to Model Nonstationary Audio Events,” *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 19, no. 4, pp. 744–753, 2011.
- [11] B. Fuentes, R. Badeau, and G. Richard, “Harmonic-adaptive Latent Component Analysis of Audio and Application to Music Transcription,” *IEEE Transactions on Audio, Speech and Language Processing*, vol. 21, no. 9, pp. 1854–1866, 2013.