

# WAVELET SCATTERING FOR AUTOMATIC CHORD ESTIMATION

First Author

Affiliation1

author1@ismir.edu

Second Author

Retain these fake authors in

submission to preserve the formatting

Third Author

Affiliation3

author3@ismir.edu

## ABSTRACT

State-of-the-art automatic chord recognition systems rely on multi-band chroma representations, Gaussian Mixture Model pattern matching, and Viterbi decoding. This paper explores the use of Haar wavelet transforms and scattering in place of multi-band chroma. Wavelets operating across octaves encode sums and differences in chroma bins at different scales. We describe both the Haar wavelet transform and deep wavelet scattering and develop an efficient algorithm for their computation. Potential benefits of wavelet representations, including stability to octave deformations, over multi-band chroma are discussed. Accuracy of wavelet representations used for chord recognition is analyzed over a large vocabulary of chord qualities.

## 1. INTRODUCTION

Along with lyrics and melody, chord sequences provide a succinct description of tonal music. As such, they are often written down under the form of lead sheets, for the use of accompanists and improvisers. Besides its original purpose in music education and transmission, the knowledge of harmonic content has been leveraged in music information research to address higher-level tasks, including cover song identification [3], genre recognition [9], and lyrics-to-audio alignment [7]. We refer to the review of McVicar et al. [8] for a recent state of the art.

As stressed by Humphrey and Bello [4], chord labels are not mutually exclusive categories, but instead follow a relation of partial order. For instance, the tetrad  $A:min7$  is contained in the triad  $A:min$ , which in turn is contained in the power chord  $A:5$ . By setting up equivalence rules for chord labels that belong to a common superset, discrepancies between conflicting annotations can be resolved up to a desired level of specificity. In spite of this variety of settings, all evaluation metrics for automatic chord estimation share the following minimal property: a chord label remains the same if all its components are jointly transposed by one octave, be it upwards or downwards.

In order to comply with this requirement, the vast majority of existing systems rely on the chroma representation, i.e. a 12-dimensional vector derived from the



**Figure 1.** Three possible voicings of the pitch class set  $\{C, E, G, A\}$ , resulting either in the chord  $A:min7$  or  $C:maj6$ . See text for details.

constant-Q spectrum by summing up all frequency bands which share the same pitch class, according to the twelve-tone equal temperament. However, it should be noted that the chroma representation is not only invariant to octave transposition, but also to any permutation of the chord factors – an operation known in music theory as inversion. Yet, although major and minor triads are unchanged by inversion, some rarer chords, such as augmented triads and minor seventh tetrads, are conditional upon the position of the root.

Figure 1 illustrates the importance of disambiguating inversions when transcribing chords. The first two voicings are identical up to octave transposition of all the chord factors, and thus have the same chord label  $A:min7$ . In contrast, the third voicing is labeled as  $C:maj6$  in root position, although its third inversion would correspond to the first voicing. Henceforth, assigning a pure tone to each note would yield the same chroma representations for all three voicings, despite the fact that the underlying synthetic sounds would be labeled differently by a human listener.

With the aim of improving automatic chord estimation under fine-grained evaluation metrics, this article introduces new feature extraction methods that are invariant to octave transposition, yet sensitive to chord inversion. We build upon the software of Cho and Bello [?], which holds state-of-the-art

A system for automatic chord estimation (ACE) typically consists of two stages: feature extraction and acoustic modeling. At the first stage, the audio query is converted into a time series of pitch class profiles, which represent the relative salience of pitch classes according to the twelve-tone equal temperament. At the second stage, each frame in the time series is assigned a chord label among a predefined vocabulary.

**Feature extraction** in the ACE systems of [2] [1] uses a Constant-Q transform to create a frequency-space representation on a perceptual log-scale of ten octaves. That is, a



time-frequency representation  $\mathbf{X}[t, \gamma]$  encodes the energy in a temporal window  $t$  at a frequency band whose center is proportional to  $2^{\gamma/12}$ , where each  $\gamma$  corresponds to a pitch in twelve-tone equal temperament. A  $K$ -multi-band chroma representation is then created by first windowing the frequency scale with  $K$  Gaussians distributed evenly over the support of  $\lambda$ , and then by wrapping each windowed representation  $\mathbf{X}[t, \gamma]$  into a chroma vector. This creates a new representation  $\mathbf{Y}[t, q, k]$  which is invariant to octave transpositions where  $q$  is an integer  $\llbracket 0, 12 \rrbracket$  representing a chroma (or pitch) value, and  $K$  is the windowed band.

**Pattern matching** is done by creating a chord model for all 157 chords with Gaussian Mixture Models (GMMs). This is a supervised learning approach where features are extracted for a large training set that contains examples of all 157 chords (with annotations) and are then used to train GMM templates for each chord. Data augmentation is employed by rotating each chord through the chroma space so that one example of C minor can also train a model for C# minor, D minor, etc. For any temporal frame  $t$ , each of the  $K$  chroma representations  $\mathbf{Y}_{t,q,k}$  are matched using GMMs and then fused together via geometric mean. Finally, **decoding** uses Hidden Markov Model (HMM) systems to post-filter the pattern matched data.

Here, we use the Viterbi algorithm, which describes the transition probability from one chord  $C_i$  state to another –  $\mathbb{P}(C_i \rightarrow C_j)$ , therefore filtering out chord transitions that are highly unlikely. The matrix of transition probabilities for Viterbi decoding is generated from the labeled training data.

## 2. EXPERIMENTAL SETUP AND EVALUATION

Evaluation of an ACE system was determined using some of the tools found in the `mir_eval` library [10] and written into a simple python script which compared the estimated chords coming out of the ACE system to the ground truth annotations for all songs in the testing set. All experiments were run on the High Performance Computing (HPC) environment at NYU’s Courant Institute.

For each experiment, a chord model and Viterbi transition probability matrix are generated from a training set of 451 songs. The ‘band’ ( $K$ ) of any experiment determines the amount of chroma vectors at any given temporal window – equivalent to the number of bands in the multiband chroma representation, and the maximum wavelet scale in the wavelet and scattering representations (i.e. the number of wavelet coefficients). Chord recognition is then carried out on the testing set of 65 songs. Both the training set and testing set of songs are kept constant across all experiments.

After generating estimated chord labels for each song in the test set, a Python script evaluates the results. As per [10], there is “no single right way to compare two sequences of chord labels,” and the `mir_eval` package computes ACE accuracy based along metrics such as root, triads, maj/min, sevenths, inversions, etc. For simplicity, we evaluate all experiments using the `mirex` measure, which

“considers a chord correct if it shares at least three pitch classes in common” [10].

State of the art systems such as those in [2] use a multi-band chroma feature extraction with  $K = 4$  windowed bands. To begin, some baseline trials were run on a training set consisting of 108 songs from the Beatles discography, 99 RWC pop songs, 224 songs from the Billboard dataset, and 20 Queen songs, for a total of 451 songs. Our testing dataset comprised of 65 songs from the Beatles and us-pop datasets that were not part of the training set and that contained a sufficient number of examples of each chord quality.

The constant-Q spectrogram  $\mathbf{X}[t, \gamma]$  is reshaped into a time-chroma-octave representation  $\mathbf{X}[t, q, u]$  where the chroma index  $q \in \llbracket 0; Q \rrbracket$  and octave index  $u \in \llbracket 0; K \rrbracket$  are derived from the log-frequency index  $\gamma$  by Euclidean division:  $\gamma = Q \times u + q$ . All subsequent operations apply to the octave variable  $u$ , and are vectorized in terms of time  $t$  and chroma  $u$ . To alleviate notations, we replace the three-way tensor  $\mathbf{X}[t, q, u]$  by a vector  $\mathbf{x}[u]$ , thus leaving the indices  $t$  and  $q$  implicit.

## 3. HAAR WAVELET TRANSFORM

Dating back to 1909, the Haar wavelet  $\psi$  is a piecewise constant, real function of compact support, consisting of two steps of equal length and opposite values. Within a discrete framework, it is defined by the following formula:

$$\forall u \in \mathbb{Z}, \psi[u] = \begin{cases} \frac{-1}{\sqrt{2}} & \text{if } u = 0 \\ \frac{1}{\sqrt{2}} & \text{if } u = 1 \\ 0 & \text{otherwise} \end{cases} \quad (1)$$

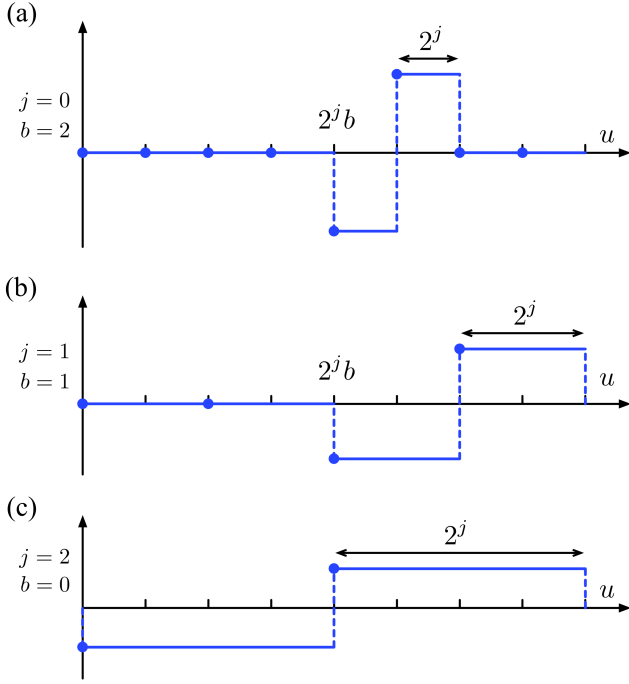
The “mother” wavelet  $\psi[u]$  is translated and dilated by powers of two, so as to produce a family of discrete sequences  $\psi_{j,b}[u] = 2^{j/2} \psi[2^j(u - 2b)]$  indexed by the scale parameter  $j \in \mathbb{N}$  and the translation parameter  $b \in \mathbb{Z}$ . After endowing them with the Euclidean inner product

$$\langle \psi_{j,b} | \psi_{j',b'} \rangle = \sum_{u=-\infty}^{+\infty} \psi_{j,b}[u] \psi_{j',b'}[u], \quad (2)$$

the wavelets  $\{\psi_{j,b}\}_{j,b}$  form an orthonormal basis of finite-energy real sequences. Moreover, the Haar wavelet is the shortest function of compact support such that the family  $\{\psi_{j,b}\}_{j,b}$  satisfies this orthonormality property. On the flip side, it has a poor localization in the Fourier domain, owing to its sharp discontinuities. Elements of the Haar wavelet basis are shown on Figure 3 for various values of  $j$  and  $b$ .

The wavelet transform coefficients of some finite-energy sequence  $\mathbf{x} \in \ell^2(\mathbb{Z})$  are defined by  $W\mathbf{x}[j, b] = \langle \mathbf{x} | \psi_{j,b} \rangle$ . Since  $\mathbf{x}[u]$  has a finite length  $K = 2^J$ , the above decomposition is informative only for indices  $(j, b)$  such that  $j < J$  and  $2^j b < K$ , that is  $b < 2^{J-j}$ . The number of coefficients in the Haar wavelet transform of  $\mathbf{x}[u]$  is thus equal to  $\sum_{j < J} 2^{J-j} = 2^J - 1$ . For the wavelet representation to preserve energy and allow signal reconstruction, a residual term

$$\mathbf{A}_J \mathbf{x} = \mathbf{x}[0] - \sum_{j,b} \langle \mathbf{x} | \psi_{j,b} \rangle \psi_{j,b}[0] = \sum_{u < K} \mathbf{x}[u] \quad (3)$$



**Figure 2.** Three elements of the Haar wavelet basis  $\{\psi_{j,b}\}$  for various values of the scale index  $j$  and the translation index  $b$ . See text for details.

must be appended to the wavelet coefficients. Observe that  $A_J x$  computes a delocalized average of all signal coefficients, which can equivalently be formulated as an inner product with the constant function  $\phi[u] = 2^{-J/2}$  over the support  $\llbracket 0; K \rrbracket$ . Henceforth, it corresponds to the traditional chroma representation, where spectrogram bands of the same pitch class  $q$  are summed across all  $K$  octaves.

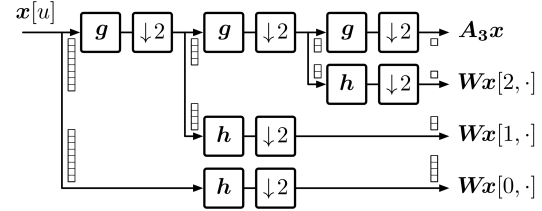
Since the wavelet representation amounts to  $K$  inner products in  $\mathbb{R}^K$ , its computational complexity is  $\Theta(K^2)$  if implemented as a matrix-vector product. Implementing these inner products as convolutions would bring the complexity to  $\Theta(K(\log_2 K)^2)$  by using the Fast Fourier Transform (FFT) algorithm. To improve this, Mallat has developed a recursive scheme, called *multiresolution pyramid* [6], which operates as a cascade of convolutions with some pair of quadrature mirror filters ( $g, h$ ) and progressive downsamplings by a factor of two.

$$Wx[j, b] = (h_{\downarrow 2} \circ g_{\downarrow 2} \circ \dots \circ g_{\downarrow 2} x)[b] \quad (4)$$

$$A_J x = g_{\downarrow 2} \circ \dots \circ g_{\downarrow 2} x \quad (5)$$

The associated Z-transforms  $G$  and  $H$  of  $g$  and  $h$  satisfy the quadrature mirror property  $G(z) = H(-z)$  and energy conservation:  $|G(z)|^2 + |H(z)|^2 = 1$ . The interested reader can refer to chapter 7 of Mallat’s textbook [5], of which we adopt the mathematical notations.

In the case of Haar wavelets, the low-pass filtering  $(x * g)$  consists of the sum between adjacent coefficients, whereas the high-pass filtering  $(x * h)$  is the corresponding difference, up to a renormalization constant:



**Figure 3.** Discrete wavelet transform of a signal of length 8, as implemented with a convolutional pyramid scheme. See text for details.

	operations	memory
Matrix-vector product	$\Theta(K^2)$	$\Theta(K)$
Fast Fourier transform	$\Theta(K(\log K)^2)$	$\Theta(K)$
Multiresolution pyramid	$\Theta(K \log K)$	$\Theta(K)$

**Table 1.** Computational complexity and memory usage of various implementations of the Haar wavelet transform, for a one-dimensional signal of length  $K$ . See text for details.

$$(x * g)[2b] = \frac{x[2b+1] + x[2b]}{\sqrt{2}} = x[2b] - \langle x | \psi_{0,b} \rangle \quad (6)$$

$$(x * h)[2b] = \frac{x[2b+1] - x[2b]}{\sqrt{2}} = \langle x | \psi_{0,b} \rangle. \quad (7)$$

When applied to the chroma representation  $\mathbf{X}(t, q, z)$ , where  $t$  is a temporal window,  $q$  a chroma pitch class (an integer between 1 and 12), and  $z$  an octave, the Haar wavelet transform encodes the sums and differences in chroma bins across octaves. The number of output “bands” from the Haar wavelet transform, given by  $k$ , determines the number of scales  $J$  for which the basis Haar wavelet is dilated:  $2^J = k$ .

The use of the Haar wavelet transform in the feature extraction stage for our ACE system is motivated by the fact that it is stable to octave deformations – something that the multiband representation is not. By coding for chroma relationships across varying octave scales, we are able to capture representations at both low scales (less invariance, more stability) and higher scales (higher invariance, less stability), thereby reducing the variance of the GMMs.

Upon receiving a CQT-derived pitch representation of 8 octaves  $\mathbf{X}[t, \lambda]$ , the lowest and highest octaves are thrown out of a total of 8 octaves (this follows suit from the multiband chroma representation).  $\mathbf{X}$  is therefore a matrix of size  $T \times \Lambda$ , i.e. the number of frames in the song  $T$  times the perceptual frequency bins  $\Lambda = (12 \times 8) = 96$ . A chroma representation is then created by reshaping  $\mathbf{X}[t, \lambda]$  into a tensor of size  $T \times Q \times Z$  where  $Q = 12$  (number of chroma pitch classes) and  $Z = k = 8$  (number of octaves as well as number of wavelet coefficients). In the sequel, we keep the notation  $\mathbf{X}[t, q, z]$  for simplicity.

The Haar wavelet transform is then computed iteratively over the scale  $J = \log_2(k)$  and stored into a cell of size

	operations	memory
Matrix-vector product	$\Theta(K^3)$	$\Theta(K^2)$
Fast Fourier transform	$\Theta(K^2(\log K)^2)$	$\Theta(K^2)$
Multiresolution pyramid	$\Theta(K^2 \log K)$	$\Theta(K^2)$

**Table 2.** Computational complexity and memory usage of various implementations of the deep Haar scattering transform, for a one-dimensional signal of length  $K$ . See text for details.

$k$ , with each row in the cell containing a chroma matrix of size  $T \times Q$  (same as in the multiband case, but with  $k = 8$  cell rows instead of 4). At each scale  $j \in \llbracket 1, J \rrbracket$ , the chroma tensor  $\mathbf{X}[t, q, z]$  is reshaped to a four dimensional tensor  $g$  from which the sum  $g$  and difference  $h$  coefficients are calculated:

$$\mathbf{X}[t, q, z] \rightarrow g[t, q, w, v] \quad (8)$$

where

$$\text{size}(\mathbf{X}) = (T \times Q \times Z) \rightarrow \text{size}(g) = (T \times Q \times 2 \times \frac{k}{2^j}) \quad (9)$$

$h$  takes the difference of  $g$  along the third dimension ( $\partial g / \partial w = h$ ). For each  $v \in [1, \frac{k}{2^j}]$ , the modulus  $|h_{t,q,w,v}|$  is stored as a row in the output Haar wavelet representation cell.  $g$  is then summed along the third dimension ( $w$ ) and squeezed down into a three dimensional tensor (singleton dimensions due to sum removed). After iterating at all scales, the resulting  $g$  is stored as the last ( $k$ 'th) row in the output Haar cell array.

#### 4. DEEP HAAR SCATTERING

Two neighboring samples  $x[2b]$  and  $x[2b+1]$  can be retrieved, up to a permutation, from the absolute values of the subsampled QMF decomposition ( $g_{\downarrow 2}x, h_{\downarrow 2}x$ ):

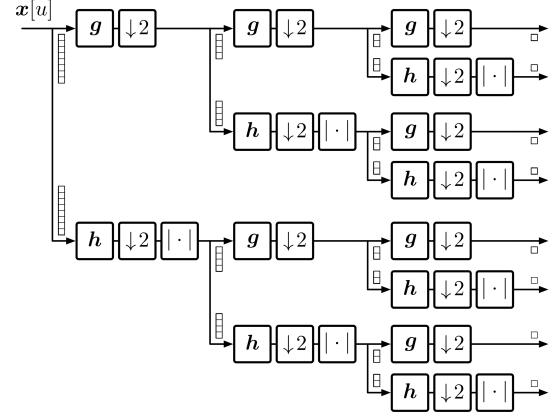
$$\max_{u \in \{2b, 2b+1\}} x[u] = \frac{|g_{\downarrow 2}x[b]| + |h_{\downarrow 2}x[b]|}{\sqrt{2}} \quad (10)$$

$$\min_{u \in \{2b, 2b+1\}} x[u] = \frac{|g_{\downarrow 2}x[b]| - |h_{\downarrow 2}x[b]|}{\sqrt{2}} \quad (11)$$

As opposed to the Haar wavelet transform representation, a Haar scattering representation concatenates sum and difference terms at different scales and then iterates. From the chroma representation  $\chi_{t,q,z}$ , we reshape to a tensor of dimension  $J+2$  called  $\phi$ :

$$\chi_{t,q,z} \rightarrow \phi_{t,q,\sigma_1,\sigma_2,\dots,\sigma_J} \quad (12)$$

where the size of  $\phi$  is  $T \times Q \times 2_1 \times 2_2 \dots \times 2_J$ . At each scale  $j \in [1, J]$ , we first create a copy of  $\phi$  called  $\hat{\phi}$ , and then concatenate the sum of  $\hat{\phi}$  along scale  $\sigma_j$  and the difference of  $\hat{\phi}$  along  $\sigma_j$ , setting the result as  $\phi$ . We then take the modulus  $|\phi|$  and iterate.



**Figure 4.** Deep scattering transform of a signal of length 8, as implemented with a convolutional pyramid scheme. See text for details.

#### 5. RESULTS

Results and discussion forthcoming.

#### 6. DISCUSSION

Results and discussion forthcoming.

#### 7. REFERENCES

- [1] Taemin Cho. *Improved Techniques For Automatic Chord Recognition From Music Audio Signals*. PhD thesis, New York University, 2013.
- [2] Taemin Cho and Juan P. Bello. On the relative importance of individual components of chord recognition systems. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.*, 22(2):477–492, February 2014.
- [3] Daniel P. W. Ellis and Graham E Poliner. Identifying “cover songs” with chroma features and dynamic programming beat tracking. In *Proc. ICASSP*, 2007.
- [4] Eric J. Humphrey and Juan P. Bello. Four timely insights on automatic chord estimation. In *ISMIR*, 2015.
- [5] Stéphane Mallat. *A Wavelet Tour of Signal Processing, 3rd edition: The Sparse Way*. Academic press, 2008.
- [6] Stephane G Mallat. A theory for multiresolution signal decomposition: the wavelet representation. *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, 11(7):674–693, 1989.
- [7] M. Mauch, H. Fujihara, and M. Goto. Integrating additional chord information into hmm-based lyrics-to-audio alignment. *IEEE TASLP*, 20(1):200–210, Jan 2012.

- [8] M. McVicar, R. Santos-Rodríguez, Y. Ni, and T. D. Bie. Automatic chord estimation from audio: A review of the state of the art. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 22(2):556–575, Feb 2014.
- [9] Carlos Pérez-Sancho, David Rizo, and José M Inesta. Genre classification using chords and stochastic language models. *Connection science*, 21(2-3):145–159, 2009.
- [10] Colin Raffel, Brian McFee, Eric J. Humphrey, Justin Salamon, Oriol Nieto, Dawen Liang, and Daniel P. W. Ellis. mir\_eval: A transparent implementation of common mir metrics. In *ISMIR*, 2014.