

Detecting Harmonic Change In Musical Audio

Christopher Harte and Mark Sandler
Centre for Digital Music
Queen Mary, University of London
London, UK
christopher.harte@elec.qmul.ac.uk

Martin Gasser
Austrian Research Institute for Artificial
Intelligence(ÖFAI)
Freyung 6/6, 1010 Vienna, Austria
martin.gasser@ofai.at

ABSTRACT

We propose a novel method for detecting changes in the harmonic content of musical audio signals.

Our method uses a new model for Equal Tempered Pitch Class Space. This model maps 12-bin chroma vectors to the interior space of a 6-D polytope; pitch classes are mapped onto the vertices of this polytope. Close harmonic relations such as fifths and thirds appear as small Euclidian distances.

We calculate the Euclidian distance between analysis frames $n + 1$ and $n - 1$ to develop a harmonic change measure for frame n . A peak in the detection function denotes a transition from one harmonically stable region to another. Initial experiments show that the algorithm can successfully detect harmonic changes such as chord boundaries in polyphonic audio recordings.

Categories and Subject Descriptors

J.0 [Computer Applications]: General

General Terms

Algorithms, Theory

Keywords

Pitch Space, Harmonic, Segmentation, Music, Audio

1. INTRODUCTION

In this paper we introduce a novel process for detecting changes in the harmonic content of audio signals. The Harmonic Change Detection Function (HCDF) combines a new theoretical model for equal tempered pitch space with a DSP front end to extract this information from digital audio recordings. The pitch space model projects collections of pitches as Tonal Centroid points in a 6-D space(section 2).

Event-driven feature analysis has been shown to give more accurate musical feature extraction than more traditional approaches based on frames of equal length [2]. The HCDF

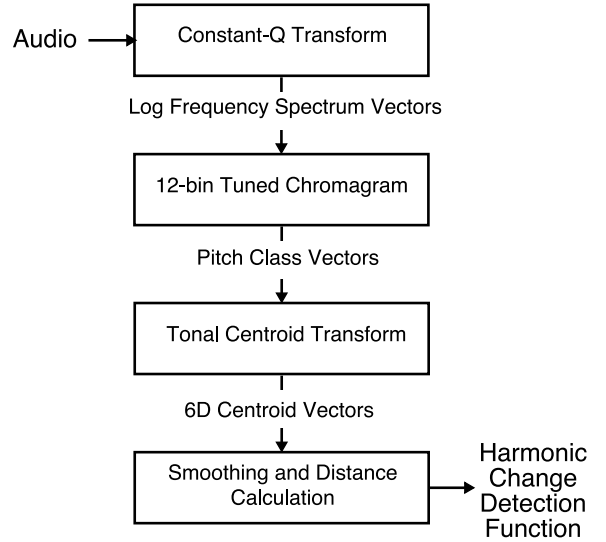


Figure 1: Flow Diagram of HCDF System.

has many potential applications in the segmentation of audio signals, particularly as a preprocessing stage for further harmonic recognition and classification algorithms. The primary motivation for this work is for use in chord recognition from audio. Used as a segmentation algorithm, this approach provides a good foundation for solving the general chord recognition problem of needing to know the positions of chord boundaries in the audio data before being able to successfully identify possible chord symbols as discussed in [16].

The HCDF system comprises several distinct elements (Figure 1). At the lowest level there is a Constant-Q spectral analysis followed by a 12-semitone Chromagram decomposition. A harmonic Centroid transform is then applied to the Chroma vectors which is then smoothed with a Gaussian filter before the distance measure is calculated (sections 3.1 and 3.2).

The results of our preliminary experiments are very promising with an f-measure value for overall chord change detection of 64.9% (section 4).

2. MODELS FOR TONAL SPACE

The Harmonic Network or Tonnetz shown in figure 2 is a well known planar representation of pitch relations first attributed to Euler [7], later used extensively by 19th century music theorists such as Riemann and Oettingen and in re-

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.

AMCMM '06 Santa Barbara, CA USA

Copyright 200X ACM X-XXXXX-XX-X/XX/XX ...\$5.00.

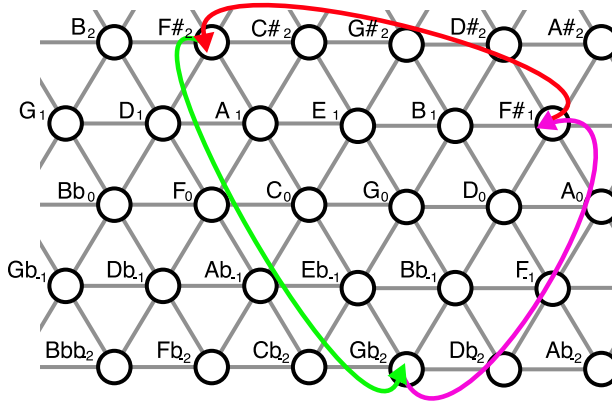


Figure 2: The **Harmonic Network** or **Tonnetz**. **Arrows show the three circularities inherent in the network** if enharmonic and octave equivalence are assumed.

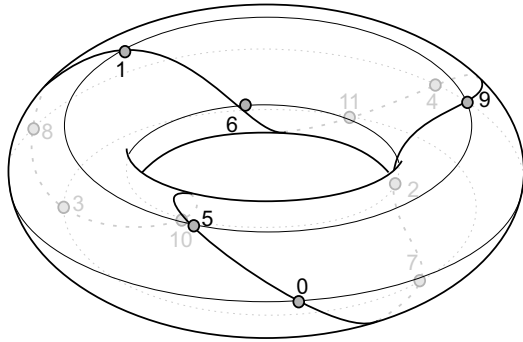


Figure 3: A projection showing how the **Tonnetz** wraps around the surface of a **Hypertorus** with the pitch classes following the spiral of fifths when enharmonic and octave equivalence are assumed.

cent years by Neo-Riemannian Music Theorists [10, 7, 12]. Close harmonic relations are modelled by small distances on the plane. **Lines of fifths** travel from left to right, **lines of major thirds** travel from bottom left to top right and **lines of minor thirds** travel from top left to bottom right.

In Just Intonation, the Tonnetz is an infinite plane [13]. If it is assumed that a particular note spelling on one row is equivalent to the same note spelling on the next row (i.e. $F\sharp_1 \equiv F\sharp_2$ etc. in fig. 2), the plane wraps up and forms a tube with the line of fifths becoming a helix on its surface. In the case where the helix is wrapped so that **major third intervals are directly above each other on the surface** of the tube this is Chew's Spiral Array [5]. Chew's model allows chords and keys to be projected as objects in a 3-D space on the interior of the tube and has been applied successfully to problems such as **key finding** and **pitch spelling** from symbolic data [6].

In the case of data derived from audio, it is very difficult to directly extract the correct spelling of pitches. This is partly due to the fact that high resolution frequency analysis would be needed to resolve the small differences between them. Equally, on a more practical level, it is because the majority of keyboard instruments are now tuned to twelve-tone equal

temperament so the differences would not be present.

If enharmonic equivalence is assumed then instead of dealing with a theoretically infinite number of pitch names, there are now just the **twelve different pitch classes** (here we reference C as pitch class 0). In the Spiral Array model, **this has the effect of joining the two ends of the tube together and the result is a hypertorus with the circle of fifths wrapping around its surface three times** (see Figure 3). A form of this Hypertorus appears in many different areas of music research [10, 7, 11, 14].

We now propose a **6-dimensional interior space** contained by the surface of the Hypertorus. This allows us to apply the same technique that Chew uses to develop the Centre of Effect in the Spiral Array to this equal tempered model for pitch space.

Since it is not possible to directly visualise 6-D space, it is helpful to imagine it as a projection onto the three circularities in the equal tempered Tonnetz: **the circle of fifths, the circle of minor thirds and the circle of major thirds** (figure 4). Here, the six dimensions are viewed as three co-ordinate pairs x_1, y_1, x_2, y_2 and x_3, y_3 . A collection of pitches (i.e. a chord) can be described as a single centroid point in the space. Chords with a **tonal centre** (such as the A major shown as point A in figure 4) can be clearly assigned to a point in the circle of fifths. However, **there are chords without defined tonal centres** (e.g. diminished 7th and augmented chords). The centroid of each of these chords lies in the centre of the circle of fifths. On the circle of **minor thirds**, however, augmented chords can be unambiguously identified, while the circle of **major thirds** can uniquely depict diminished 7th chords.

3. ALGORITHM

The first stage of the system is the **Constant-Q spectral analysis**. This is a logarithmic frequency analysis based on the efficient algorithm described in [3]. We calculate a **36 bins-per-octave transform** across five octaves between $f_{min} = 110\text{Hz}$ (A2) and $f_{max} = 3520\text{Hz}$ (A7) from a 11025Hz mono audio signal. To obtain this resolution at the lowest analysed frequencies requires a 743ms window length. This is a long analysis window in terms of musical signals so to improve time resolution we overlap analysis frames by $\frac{1}{8}$ th of a window length giving an effective frame length of 93ms. A **12-bin tuned Chromagram** is then calculated from the Constant-Q spectra using the method described in [9] giving a **12-dimensional chroma vector \mathbf{c}** for every frame.

3.1 Tonal Centroid Calculation

The **six dimensional tonal centroid vector, ζ** , for time frame n is given by the multiplication of the **chroma vector, \mathbf{c}** , and a **transformation matrix Φ** . Dividing by the L_1 norm of \mathbf{c} prevents numerical instability and ensures that the tonal centroid always lies within the 6-D polytope (equation 1).

$$\zeta_n(d) = \frac{1}{\|\mathbf{c}_n\|_1} \sum_{l=0}^{11} \Phi(d, l) \mathbf{c}_n(l) \quad \begin{matrix} 0 \leq d \leq 5 \\ 0 \leq l \leq 11 \end{matrix} \quad (1)$$

where l is the **chroma vector pitch class index** and d denotes which of the six dimensions of ζ_n is being evaluated. The transformation matrix Φ represents the basis of the 6-D space described in section 2 and is given as:

$$\Phi = [\phi_0, \phi_1 \dots \phi_{11}] \quad (2)$$

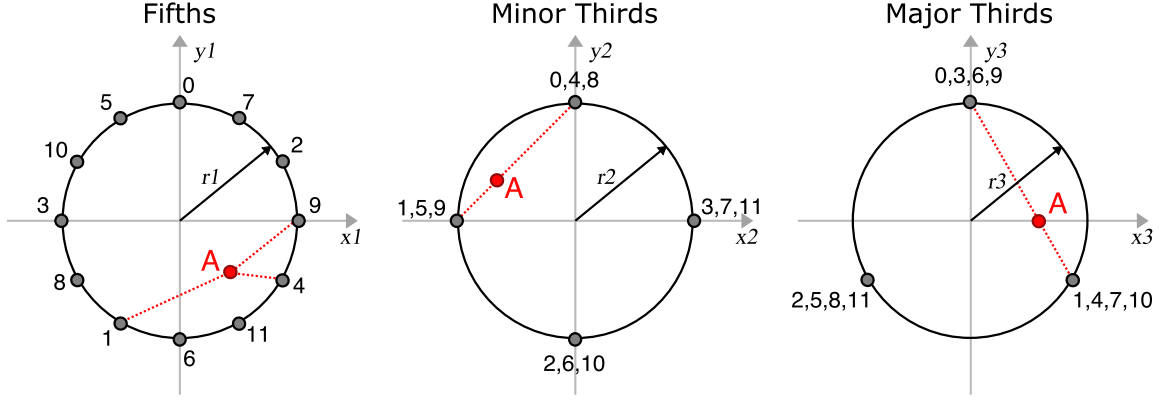


Figure 4: Visualising the the 6-D Tonal Space as three circles. Circles left to right: **Fifths**, **Minor Thirds** and **Major Thirds**. The Tonal Centroid for **chord A Major** (pitch classes 9,1 and 4) is shown at point A

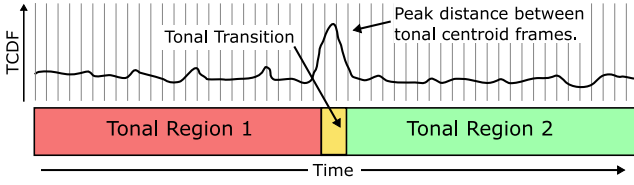


Figure 5: Harmonic Change Detection.

where

$$\phi_l = \begin{bmatrix} \Phi(0, l) \\ \Phi(1, l) \\ \Phi(2, l) \\ \Phi(3, l) \\ \Phi(4, l) \\ \Phi(5, l) \end{bmatrix} = \begin{bmatrix} r_1 \sin l \frac{7\pi}{6} \\ r_1 \cos l \frac{7\pi}{6} \\ r_2 \sin l \frac{3\pi}{2} \\ r_2 \cos l \frac{3\pi}{2} \\ r_3 \sin l \frac{2\pi}{3} \\ r_3 \cos l \frac{2\pi}{3} \end{bmatrix} \quad 0 \leq l \leq 11 \quad (3)$$

The values r_1 , r_2 and r_3 are the radii of the three circles in Figure 4. To ensure that the distances between pitch classes in the 6-D space correspond to our perception of harmonic relations between pitches (i.e. that the fifth is the closest relation followed by the major third then the minor third and so on) we set the r_1 , r_2 and r_3 to 1, 1 and 0.5 respectively. These values have been derived using the same approach that Chew uses to define the ratio of height to diameter for the Spiral Array [5].

3.2 Harmonic Change Detection Function

To reduce the effects of transient frames, the sequence of tonal centroid vectors is convolved with a Gaussian with σ value of 8 in a row-by-row fashion (i.e., the individual dimensions are smoothed over time). We define the HCDF, ξ , as the overall rate of change of the smoothed tonal centroid signal. ξ_n is the euclidian distance between the tonal centroid vectors ζ_{n-1} and ζ_{n+1} (equation 4). Peaks in this signal indicate transitions between regions that are harmonically stable (see figure 5); an approach inspired by Chew's key modulation finding algorithm described in [6].

$$\xi_n = \sqrt{\sum_{d=0}^5 [\zeta_{n+1}(d) - \zeta_{n-1}(d)]^2} \quad (4)$$

4. RESULTS

The HCDF was implemented in Matlab for experiments and also in C++ as a visualisation plugin for Sonic Visualiser [4]. To test how the HCDF performed as a chord segmentation algorithm we analysed a set of sixteen Beatles songs (two from each of the first eight albums picked) for which we have chord transcription files. We apply peak picking to the HCDF in order to identify harmonic transition times and these are compared against the times of chord changes in the transcriptions.

The results of this experiment are shown in Table 1. We also give results for a harmonic onset detection algorithm by Hainsworth and Macleod [8] for comparison. We defined a hit as a match within ± 3 frames (278ms). The three performance measures used here are *Precision* (P), the ratio of Hits to Detected Changes and *Recall* (R), the ratio of hits to transcribed changes and the *f-measure* (F) which combines the two (see equation 5) [1].

$$F = \frac{2RP}{R+P} \quad (5)$$

The f-measure scores show the HCDF to have the better overall performance of the two algorithms for chord boundary detection with a score of 64.9% compared to 45.8% for the Hainsworth algorithm. The Recall scores for both algorithms are fairly high with an average of 88% for Hainsworth's approach and 84% for the HCDF. However, the Precision scores for the two algorithms were much lower with averages of 31% for the Hainsworth algorithm and 53% for the HCDF. The significantly better Precision scores for the HCDF suggest that the new algorithm is better at discriminating important harmonic changes in the signal. This can be seen from the number of detection function peaks for each algorithm where the Hainsworth algorithm has an average of almost twice as many detections for each song as the HCDF. The low Precision scores for both algorithms can be explained by the fact that the transcription files only label chord changes. Both the detection functions here, however, pick up not only chord changes but also changes in harmonic content caused by strong melody or bass lines that include non-chord tones. Because of this a high number of false positives is to be expected for this experiment. Most misses in the HCDF are caused by transient signals introducing false peaks that mask the smaller peaks associated with the

Table 1: Experimental results for HCDF peaks compared with hand labelled chord changes for 16 Beatles Songs (Songs arranged in chronological order of release date).

Song Title	TC	Mn	Mp	Mr	Mf	Hn	Hp	Hr	Hf
Please Please Me	78	267	27%	92%	41%	128	53%	87%	65.9%
Do You Want To Know A Secret	113	214	51%	97%	66.9%	126	72%	80%	75.8%
All My Loving	74	230	26%	81%	39.4%	129	50%	86%	63.2%
Till There Was You	93	265	33%	88%	48%	142	59%	90%	71.3%
A Hard Day's Night	101	343	28%	94%	43.1%	158	45%	70%	54.8%
If I Fell	81	237	32%	95%	47.8%	133	53%	87%	65.8%
Eight Days A Week	101	316	25%	80%	38%	169	49%	82%	61.3%
Every Little Thing	96	247	34%	84%	48.4%	133	49%	67%	56.6%
Help!	59	269	20%	84%	32.3%	152	29%	74%	41.7%
Yesterday	97	186	47%	83%	60%	132	64%	86%	73.4%
Drive My Car	84	328	24%	94%	38.2%	148	49%	86%	62.4%
Michelle	94	272	31%	90%	46.1%	160	53%	90%	66.7%
Eleanor Rigby	54	234	22%	96%	35.8%	128	34%	81%	47.9%
Here There And Everywhere	98	240	32%	73%	44.5%	127	65%	83%	72.9%
Lucy In The Sky With Diamonds	120	411	28%	97%	43.5%	213	50%	88%	63.8%
Being For The Benefit Of Mr Kite	113	255	38%	80%	51.5%	160	68%	95%	79.5%
Average over sixteen songs	91	270	31%	88%	45.8%	146	53%	84%	64.9%

Key to abbreviations

M	Denotes result for Hainsworth & Macleod's algorithm	p	Precision
H	Denotes result for HCDF algorithm	r	Recall
TC	Number of hand transcribed chord changes	f	f-measure
n	Number of detection function peaks		

desired harmonic changes.

The songs that the HCDF algorithm performed best on were ones with fast harmonic rhythm such as *For The Benefit Of Mr Kite!*. The fast chord changes reduce the number of false positives and a strong organ part outlining the chords makes boundary detection easier. In contrast, songs such as *Every Little Thing* and *Help!* with slow harmonic rhythm and strong bass and melody lines that produce false chord boundary detections score less highly.

5. CONCLUSIONS AND FURTHER WORK

We have presented a novel feature detection function for audio data. A new model for equal tempered tonal space has been introduced on which the algorithm is based and the results of our preliminary experiments are encouraging. The tests show that the algorithm can successfully detect chord changes. Other changes in harmonic content such as strong melody or bass line movement will also be detected. Applying adaptive thresholding may improve the detection of more important harmonic changes. Strong transient signals can cause true peaks to be masked. Applying a Transient/Steady-State separation to the audio may rectify this problem.

The HCDF has many potential applications in the segmentation of audio signals. It will be particularly useful as a preprocessing stage for many chord recognition and harmonic classification algorithms, especially those based on Hidden Markov Model techniques such as [2] and [15].

6. ACKNOWLEDGMENTS

This research was funded by EU-FP6-IST-507142 project

SIMAC¹(acronym for Semantic Interaction with Music Audio Contents) and the WWTF project Interfaces to Music (I2M)².

7. REFERENCES

- [1] J. P. Bello, L. Daudet, S. Abdallah, C. Duxbury, M. Davies, and M. Sandler. A tutorial on onset detection in music signals. *IEEE Transactions on Speech and Audio Processing*, 13(5, part 2):1035–1047, 2005.
- [2] J. P. Bello and J. Pickens. A Robust Mid-level Representation for Harmonic Content in Musical Signals. In *Proceedings of the 6th International Conference on Music Information Retrieval, London*, 2005.
- [3] J. C. Brown and M. S. Puckette. An Efficient Algorithm for the Calculation of a Constant Q Transform. *Journal of the Acoustical Society of America*, 92(5):2698–2701, 1992.
- [4] C. Cannam, C. Landone, M. Sandler, and J. Bello. The Sonic Visualiser: A Visualisation Platform For Semantic Descriptors From Musical Signals. In *Submitted to ISMIR 2006*, Victoria, Canada, 2006.
- [5] E. Chew. *Towards a Mathematical Model of Tonality*. PhD thesis, Operations Research Center, MIT. Cambridge, MA, 2000.
- [6] E. Chew. The Spiral Array: An Algorithm For Determining Key Boundaries. *Proceedings of the Second International Conference, ICMAI 2002*, pages 18–31, 2002.

¹<http://www.semanticaudio.org>

²<http://www.ofai.at/research/impml/projects/i2mproject.html>

- [7] R. Cohn. Introduction to Neo-Riemannian Theory: A Survey and a Historical Perspective. *The Journal of Music Theory*, 42(2), 1998.
- [8] S. Hainsworth and M. Macleod. Onset Detection in Musical Audio Signals. In *Proceedings of ICMC*, Singapore, 2003.
- [9] C. Harte and M. Sandler. Automatic Chord Identification Using a Quantised Chromagram. In *Proceedings of AES 118th Convention, Barcelona*, 2005.
- [10] B. Hyer. Re-Imagining Riemann. *Journal of Music Theory*, 39(1):101–138, 1995.
- [11] C. L. Krumhansl. *Cognitive Foundations of Musical Pitch*. Oxford University Press, New York, 1990.
- [12] D. Lewin. *Generalized Musical Intervals and Transformations*. Yale University Press, New Haven, 1987.
- [13] H. Longuet-Higgins. Second Letter to a Musical Friend. *The Music Review*, 23, 1962.
- [14] H. Purwins. *Profiles of Pitch Classes Circularity of Relative Pitch and Key Experiments, Models, Computational Music Analysis, and Perspectives*. PhD thesis, Elektrotechnik und Informatik der Technischen Universität Berlin, Berlin, 2005.
- [15] A. Sheh and D. P. Ellis. Chord Segmentation and Recognition using EM-Trained Hidden Markov Models. *Proceedings of the ICMC 2003*, 2003.
- [16] T. Yoshioka, T. Kitahara, K. Komatani, T. Ogata, and H. G. Okuno. Automatic Chord Transcription with Concurrent Recognition of Chord Symbols and Boundaries. In *Proceedings of ISMIR*, 2004.