

## MIREX 2008

# AUDIO MUSIC CLASSIFICATION USING A COMBINATION OF SPECTRAL, TIMBRAL, RHYTHMIC, TEMPORAL AND SYMBOLIC FEATURES

T. Lidy, A. Rauber

Vienna University of Technology, Austria  
Department of Software Technology  
and Interactive Systems

A. Pertusa, P. J. Ponce de León, J. M. Iñesta

University of Alicante, Spain  
Departamento de Lenguajes y  
Sistemas Informáticos

### ABSTRACT

The novel approach of combining audio and symbolic features for music classification from audio enhanced previous audio-only based results in MIREX 2007. We extended the approach by including temporal audio features, enhancing the polyphonic audio to MIDI transcription system and including an extended set of symbolic features. Recent research in music genre classification hints at a glass ceiling being reached using timbral audio features.

### 1 INTRODUCTION

Classification of music by genre, artist or mood are important tasks for retrieval and organization of music databases. Traditionally the research domain of music classification was divided into the audio and symbolic music analysis and retrieval domains. Our work is aimed at combining approaches from both directions that have proved their reliability in their respective domains. We are combining spectrum-based audio feature extractors, that include aspects such as rhythm, timbre and temporal evolution of signals on various critical frequency bands, with symbolic descriptors, based on note onsets and statistics, using a polyphonic transcription system as an intermediate step. These features are complementary; a score can provide very valuable information, but audio features (e.g., the timbral information) are also very important for classification, e.g. into various genres.

To extract symbolic descriptors from an audio signal it is necessary to first employ a transcription system in order to detect the notes stored in the signal. Transcription systems have been investigated previously but a well-performing solution for polyphonic music and a multitude of genres has not yet been found. Though these systems might not be in a final state for solving the transcription problem, our hypothesis is that they are able to augment the performance of music classification by introducing features on the symbolic level.

The overall scheme of our proposed genre classification system is shown in Figure 1. It processes an audio file in two

ways to predict its genre. While in the first branch, the audio feature extraction methods described in Section 2.1 are applied directly to the audio signal data, there is an intermediate step in the second branch. A polyphonic transcription system, described in Section 2.2.1, converts the audio information into a symbolic notation (i.e. MIDI files). Then, a symbolic feature extractor is applied on the resulting representation, providing a set of symbolic descriptors as output. The audio and symbolic features extracted from the music serve as combined input to a classifier.

The basic system is outlined and described in more detail in [2]. We extended the approach by including temporal audio features, enhancing the polyphonic transcription system and including an extended set of symbolic features, as outlined in the following section.

### 2 SYSTEM DESCRIPTION

#### 2.1 Audio Feature Extraction

All the following descriptors are extracted from a spectral representation of 6 sec. segments in the audio signal. While in full length songs, the number of segments varies and can be controlled using a 'step\_width' parameter, in a 30-second audio clip, usually 5 segments are extracted. Rhythm Patterns and Rhythm Histograms are summarized using the median over the 5 segments, Statistical Spectrum Descriptors are summarized computing the mean. For Temporal Rhythm Histograms and Temporal Statistical Spectrum Descriptors statistics that measure variation over time

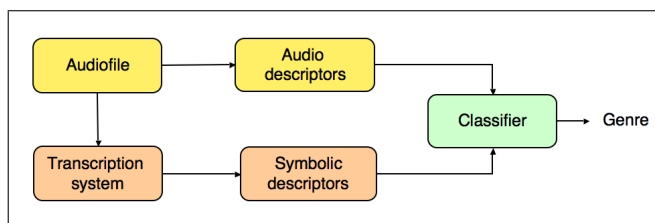


Figure 1. General framework of the system

(i.e. over the 5 segments) are computed. Note that in contrast to MIREX 2007 we did not include Onset features in this submission (due to a change in implementation).

### 2.1.1 Rhythm Pattern (RP)

The feature extraction process for a Rhythm Pattern [4, 1] is composed of two stages. First, the specific loudness sensation on 24 critical frequency bands is computed, by using a Short Time FFT, grouping the resulting frequency bands to the Bark scale, applying spreading functions to account for masking effects and successive transformation into the Decibel, Phon and Sone scales. This results in a psycho-acoustically modified Sonogram representation that reflects human loudness sensation. In the second step, a discrete Fourier transform is applied to this Sonogram, resulting in a (time-invariant) spectrum of loudness amplitude modulation per modulation frequency for each individual critical band. After additional weighting and smoothing steps, a Rhythm Pattern exhibits magnitude of modulation for 60 modulation frequencies (between 0.17 and 10 Hz) on the 24 critical bands. Note that when using 22kHz audio, the number of critical bands is reduced to 20 and the final Rhythm Pattern has 1200 dimensions. For details refer to [4, 1].

### 2.1.2 Rhythm Histogram (RH)

A Rhythm Histogram (RH) aggregates the modulation amplitude values of the individual critical bands computed in a Rhythm Pattern and is thus a lower-dimensional descriptor for general rhythmic characteristics in a piece of audio [1]. A modulation amplitude spectrum for critical bands according to the Bark scale is calculated, as for Rhythm Patterns. Subsequently, the magnitudes of each modulation frequency bin of all critical bands are summed up to a histogram, exhibiting the magnitude of modulation for 60 modulation frequencies between 0.17 and 10 Hz.

### 2.1.3 Temporal Rhythm Histogram (TRH)

Statistical measures (mean, median, variance, skewness, kurtosis, min and max) are computed over the individual Rhythm Histograms extracted from various segments in a piece of audio. Thus, change and variation of rhythmic aspects in time are captured by this descriptor.

### 2.1.4 Statistical Spectrum Descriptor (SSD)

In the first part of the algorithm for computation of a Statistical Spectrum Descriptor (SSD) the specific loudness sensation is computed on 24 Bark-scale bands, equally as for a Rhythm Pattern. Subsequently, the mean, median, variance, skewness, kurtosis, min- and max-value are calculated for each individual critical band. These features computed

for the 24 bands constitute a Statistical Spectrum Descriptor. SSDs describe fluctuations on the critical bands and are able to capture additional timbral information compared to a Rhythm Pattern, yet at a much lower dimension of the feature space, as shown in the evaluation in [1].

### 2.1.5 Temporal Statistical Spectrum Descriptor (TSSD)

Statistical measures (mean, median, variance, skewness, kurtosis, min and max) are computed over the individual Statistical Spectrum Descriptors extracted from the various segments of a piece of audio. This captures timbral variations and changes over time in the spectrum on the individual critical frequency bands.

### 2.1.6 Modulation Frequency Variance Descriptor (MVD)

This descriptor measures variations over the critical frequency bands for a specific modulation frequency (derived from a Rhythm Pattern). Consider a Rhythm Pattern, i.e. a matrix representing the amplitudes of 60 modulation frequencies on 24 critical bands: The MVD vector is computed by taking statistics (mean, median, variance, skewness, kurtosis, min and max) for one modulation frequency over the 24 (resp. 20) bands. A vector is computed for each of the 60 modulation frequencies. The MVD descriptor for an audio file is computed from the mean over the multiple MVDs of its segments.

## 2.2 Symbolic Feature Extraction

### 2.2.1 Transcription System

To complement the audio features with symbolic features we developed a polyphonic transcription system to extract the notes. This system converts the audio signal into a MIDI file that will later be analyzed to extract the symbolic descriptors. It does not consider rhythm, only pitches and note durations are extracted. Therefore, the transcription system converts a mono audio file sampled at 22 kHz into a sequence of notes. The transcription system is described in detail in [2], the latest version is presented at the MIREX 2008 Multiple Fundamental Frequency Estimation & Tracking task with an abstract describing the system.

### 2.2.2 Symbolic Features

A set of 53 symbolic descriptors was extracted from the transcribed notes. This set is based on the features described in [3], that yielded good results for monophonic classical/jazz classification, and on the symbolic features described in [6], used for melody track selection in MIDI files. The number of notes, number of significant silences, and the number of non-significant silences were computed. The occupation rate (sounding notes periods with respect

to song length) and polyphony rate (proportion of sounding periods with more than one note active simultaneously) were also computed. Note pitches, pitch intervals, note durations, silence durations, Inter Onset Intervals (IOI) and non-diatonic notes were also analyzed (for the latter, the song key is guessed using the algorithm described in [5]). Each one of this properties is described by their highest and lowest values, their range, average, relative average, standard deviation, and a normality estimation. The total number of IOI was also taken into account, as the number of distinct pitch intervals, the most repeated pitch interval, the sum of all note durations and an estimation of the number of syncopations in the song, completing the symbolic feature set.

### 3 CLASSIFICATION

#### 3.1 Classification Setup

With the availability of multiple feature sets as a source of music description, and potentially also multiple classifiers, there are several alternatives of how to design a music classification system. The option we chose is to merge the different feature sets and provide the combined set as input to a common classifier that receives an extended set of feature attributes on which it bases its classification decision (c.f. Figure 1). For our experiments we chose linear Support Vector Machines. We used the SMO implementation of the Weka machine learning software [7] with pairwise classification and the default Weka parameters (complexity parameter  $C = 1.0$ ). We investigated the performance of the feature sets individually in advance and then decided which feature sets to combine.

#### 3.2 Feature Combinations

Various combinations of different feature sets achieve divergent results on different music test collections. Therefore, we included 4 algorithm variants in our submission, in order to observe the differing performances of different feature set combinations. Moreover, comparing the variants is important to gain insight in the usefulness of symbolic features for audio music classification. The 4 variants are based on the following sets of features (refer to Section 2):

1. RH+SSD (audio only)
2. RP+MVD+SSD (audio only)
3. RP+SSD+TRH+Symbolic
4. RP+RH+TRH+SSD+TSSD+MVD+Symbolic

### 4 ACKNOWLEDGMENTS

This work is supported by the Austrian Academic Exchange Service (ÖAD) through the IMPACT project and the Spanish PROSEMUS project with code TIN2006-14932-C02.

### 5 REFERENCES

- [1] T. Lidy and A. Rauber. Evaluation of feature extractors and psycho-acoustic transformations for music genre classification. In *Proc. ISMIR*, pages 34–41, London, UK, September 11-15 2005.
- [2] T. Lidy, A. Rauber, A. Pertusa, and J.M. Iñesta. Improving genre classification by combination of audio and symbolic descriptors using a transcription system. In *Proc. ISMIR*, Vienna, Austria, 2007.
- [3] P. J. Ponce de León and J. M. Iñesta. A pattern recognition approach for music style identification using shallow statistical descriptors. *IEEE Trans. on Systems Man and Cybernetics C*, 37(2):248–257, 2007.
- [4] A. Rauber, E. Pampalk, and D. Merkl. The SOM-enhanced JukeBox: Organization and visualization of music collections based on perceptual models. *Journal of New Music Research*, 32(2):193–210, June 2003.
- [5] D. Rizo, J.M. Iñesta, and P.J. Ponce de León. Tree model of symbolic music for tonality guessing. In *Proc. of the IASTED Int. Conf. on Artificial Intelligence and Applications, AIA 2006*, pages 299–304, Innsbruck, Austria, 2006. IASTED, Acta Press. ISBN 0-88986-404-7.
- [6] D. Rizo, P.J. Ponce de León, C. Pérez-Sancho, A. Pertusa, and J.M. Iñesta. A pattern recognition approach for melody track selection in midi files. In *Proc. ISMIR*, pages 61–66, Victoria, Canada, 2006.
- [7] I.H. Witten and E. Frank. *Data Mining: Practical machine learning tools and techniques*. Morgan Kaufmann, San Francisco, 2nd edition, 2005.