# Musical style identification using self-organising maps

Pedro J. Ponce de León and José M. Iñesta
Dept. Lenguajes y Sistemas Informáticos
Universidad de Alicante, Spain
pjleon@mail.ono.es , inesta@dlsi.ua.es

## Abstract

*In this paper the capability of using self-organising neural maps (SOM) as music style classifiers from symbolic specifications of musical fragments is studied. From MIDI file sources, the monophonic melody track is extracted and cut into fragments of equal length. From these sequences, melodic, harmonic, and rhythmic numerical descriptors are computed and presented to the SOM. Their performance is analysed in terms of separability in different music classes from the activations of the map, obtaining different degrees of success for classical and jazz music. This scheme has a number of applications like indexing and selecting musical databases or the evaluation of style-specific automatic composition systems.*

## 1 Introduction

There are a number of applications in computer music to the possibility of melodic fragment comparison. Two main representations of music can be found: sounds (recorded from human or computer interpretation of a music score) and symbols (representation codes independent of the sonic outcome of an interpretation). The automatic machine learning and pattern recognition techniques available, successfully employed in other fields, can be also applied in music analysis. Immediate applications are the classification, indexation and content-based search of digital music libraries, where digitised (MP3), sequenced (MIDI) or structurally represented (XML) music can be found.

One of the tasks that can be posed is the modelization of the music style, providing the computer with the capability of discrimination between musical styles or sub-styles, or even between different composers. Even more, the computer could be trained in a given user musical taste in order to look for that kind of music over large musical databases. Such a model could be used in cooperation with automatic composition procedures to guide this latter process according to some stylistic profile provided by the user.

The aim of this work is to develop a system able to distinguish among a set of musical styles from a symbolic representation of a melody. We have chosen random, jazz and classical melodies for our experiments. Simultaneously we will investigate whether such a representation by itself has enough information to achieve this goal or, on the contrary, there is also timbric information that has to be included for that purpose.

The key point of this work is to test the ability of self-organising maps (SOM) [2], to automatically perform this task. SOM are neural methods able to obtain approximate projections of high-dimensional data distributions in low-dimensional spaces, usually bidimensional. With the map, different clusters in the input data can be located. These clusters can be semantically labelled to characterise the training data and also hopefully future new inputs.

### 1.1 Related work

In a very recent paper, Rauber and Frühwirth [4] pose the problem of organising music digital libraries according to the sound features of each musical theme, in such a way that similar themes are clustered. This would allow the user to locate sections within the library according to stylistic similarities. The authors utilise a system with two SOMs hierarchically organised in order to create a map of the digital library, where similar music themes can be found in zones close to one another in the bidimensional SOM. After finding a given music in the map, others related can be found with an exploration of the surroundings of that point in the map, permitting an intuitive exploration of the library. This is, therefore, a content-based classification of the data (sounds in that case).

Other related work is that of Whitman and Flake [7] in which they present a system named Minnowmatch, based on neuronal nets and support vector machines, able to classify a sound musical fragment into a given source or artist. The system achieves a success rate of 91% with 5 different artists or sources, of 70% with 10 artists and of 46% with 21 artists.

In a similar work to the latter [5], the authors describe a system to recognise music types using an explicit-time modelling (ETM) neural net that codes an abstraction of acoustic events in the hidden layer of the net representing temporal structures of the musical parts. This abstractions are then used to discriminate among different types of music. The experiments show that the system improves the recognition rate of other methods like recurrent neural nets or hidden Markov models.

In [1] the authors present a hierarchical SOM able to analyse time series of musical events. The model can recognise instances of a reference sequence (a fugue by J.S. Bach) in presence of noise, and even discriminate those instances in a different musical context. In this work, the SOM is used as sequence recognisers, using a time integration mechanism in the input layer of two SOM, arranged one on top of the other, to represent the reference monophonic melodic sequence in order to provide the SOM with the ability of processing time sequences.

In the work by Thom [6] pitch histograms (measured in semitones relative to the central pitch of the tonality and independent of the octave) are used to describe blues fragments of the saxophonist Charlie Parker. The pitch frequencies are used to train a SOM.

All these works attack the same problem that we face here, and most of them use digital sound files as an input. Only the last two ones use symbolic representations for recognising musical parts, not styles. The approach we propose here is to use the symbolic representation of music as the input to the self-organising maps for classification of musical fragments into a initially reduced set of styles. The success of this approach would permit to extend it to other styles and to apply this methodology to the huge amount of symbolic data stored in music databases all over the Internet. We use standard MIDI files as the source of monophonic melodies that will be preprocessed to provide melodic, harmonic and rhythmic descriptors to the SOM.

## 2 Methodology

The monophonic melodies are isolated for the experiments from the rest of the musical content in the MIDI files. This way we have a sequence of musical events that can be either notes or silences. Other kind of MIDI events are filtered out. Each note can take a value from 0 to 127 (the pitch) and the duration is the distance from the event that onsets the sound of a note to the event that finishes it (there is no limit to this in theory). Note that this symbolic representation implies the lack of timbre information. The situation is much like an expert trying to classify tunes from an overview of the scores, rather than hearing an interpretation from an instrument playing the score.

Even dealing with monophonic melodies the search

space is very vast. Nevertheless, we can think that melodies from a same musical genre may share some common features that make possible that a experienced listener is able to assign a musical style to them.

Each melody has a number of events that is a function of some features like the time signature and the number of bars of total duration, among others. Here we will deal only with melodies written in 4/4. In order to have more restricted data, fragments of 8 bars are taken (enough to get a good sense of the melodic phrase in the context of a 4/4 signature). For this, each melody sequence has been cut into fragments of such duration.

We have chosen a vector of musical descriptors of the melodies as the input for the SOM, rather than the explicit representation of the melodies. Thus, a description model is needed. This model is composed of three groups of features: melodic, harmonic and rhythmic properties. From each melody, a record is generated containing all the descriptors needed for training and testing the SOM.

For the experiments we have considered, along with real melodies, other randomly generated melodies in order to test the ability to separate well constructed melodies from other non sense musical constructions. For generating this kind of melodies each bar was divided into $Q$ pulses (quantisation) and the melody was considered as formed by three kinds of events that can appear at each pulse: note onsets, silences and continuation of the previous event.

Onset events can take values in [0,127] (indicating which note to play in the range of MIDI *note on* events), but this possibility generates totally unnatural sequences. The notes have been restricted to a more natural range of two octaves and a half (30 possible pitches) in the range [45,82], heuristically determined after an analysis of a large number of real melodies. In 8 bars we will have $8 \times Q$ events. Each melody was generated with a proportion of notes / silences / continuations among this possibilities:

| 1-1-1 | 1-1-2 | 1-2-1 |
|-------|-------|-------|
| 1-2-3 | 2-1-1 | 2-1-3 |
| 2-3-1 | 3-1-2 | 3-2-1 |

where $N$-$S$-$C$ indicates the probability of generating a note onset ($N$), a silence ($S$) or a continuation event ($C$), according to the expression $X/(N + S + C)$ where $X$ can be $N, S$ or $C$. Therefore a melody generated according to the pattern 2-3-1 will have nearly a 33% of note onset events, a 50% of silence events and a 17% of continuation events.

In the next experiments we will initially consider the next list of descriptors:

- Overall descriptors:

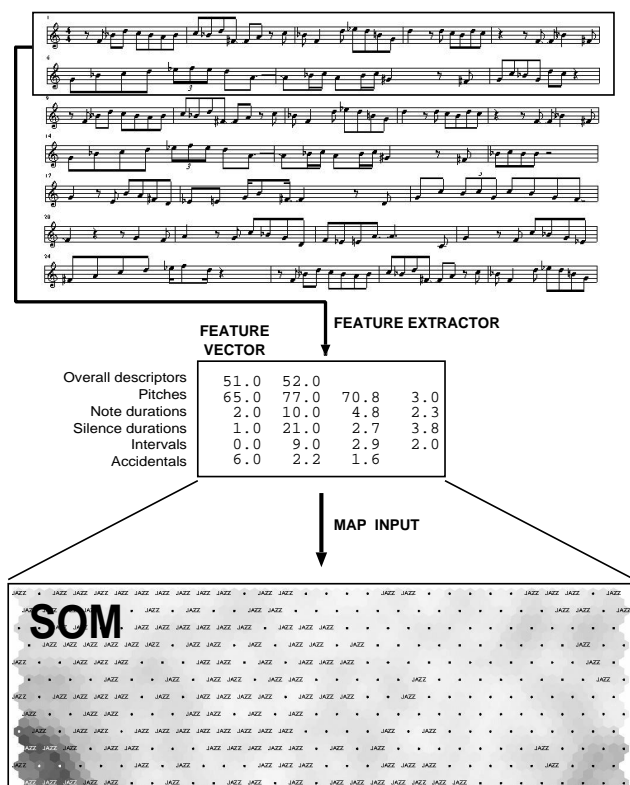  - Number of notes in the windowed melody.
  - Number of silences.

**Figure 1. Structure of the system: musical descriptors are computed from a window 8-bar wide and provided to the SOM.**

- Pitch descriptors:

  - Lowest, highest (this with the previous one provide information about the pitch range of the melody), average, standard deviation (this and the former one provide information about how the notes are distributed in the score).

- Note duration descriptors (these descriptors are measured in pulses):

  - Minimum, maximum, average, and standard deviation.

- Silence duration descriptors (in pulses):

  - Minimum, maximum, average, and standard deviation.

- Interval descriptors (distance in pitch between two consecutive notes, measured in semitones):

  - Minimum, maximum, average, and standard deviation.

These 18 features are melodic descriptors. We have considered durations as melodic rather than rhythmic descriptors. The number of features will be increased in further experiments in order to better evaluate harmonic and rhythmic aspects of the melodies. In Fig. 1 an example of the system is presented. This example is based on the jazz piece "Dexterity". The features are computed using a quantisation $Q = 48$ pulses per bar.

## 3  Experiments and results

The experiments are divided into two phases: first, a set of random melodies with different proportions of notes, silences and continuation events are generated, and a set of real melodies, extracted from jazz standards, is built. We put the capability of SOM for this task to a test with an, a priori, easy task: to separate random musical sequences from melodies with real musical feeling. The jazz samples were taken from a book of jazz standards and the melodies were sequenced in real time. Authors included Charly Parker, Wayne Shorter, John Coltrane, George Shearing, Miles Davies, Duke Ellington, Dexter Gordon, Herbie Hancock, Thelonio Monk, Dizzie Gilespie, Bill Evans, Antonio Carlos Jovim, Vernon Duke, Oliver Nelson, Richard Rogers and Lorenz Hart, amnog others. These tunes are from different jazz styles like be-bop, hard-bop, big-band swing, etc.

The second phase consists of substituting music of other type different from jazz for the random melodies in order to test the ability for style discrimination. Classical music was chosen and melodic samples were taken from works by Mozart, Bach, Schubert, Chopin, Grieg, Vivaldi, Schumann, Brahms, Beethoven, Dvorak, Haendel, Paganini and Mendhelson. All of them were downloaded from the Internet and selected using a criterion of monophony for the melodic track. Styles included baroque, romantic, renaissance, impressionism, etc.

For SOM implementation and graphic representations the SOMPAK software [3] has been used. For the experiments a hexagonal geometry for unit connections and a bubble neighbourhood for training have been selected. The value for this neighbourhood is constant for all the units in it and decreases as a function of time.

In this paper, two main kinds of map representations are shown: the Sammon projection, as a way to display in 2D the organisation of the weight vectors in the weight space, and the U-map representation, where the units are represented by hexagons with a dot or label in their centre. The grey level of unlabelled hexagons represents the distance between neighbour units (the clearer the closer they are). The grey level of labelled units is an average of those distances. This way, clear zones are clusters of units in the SOM, sharing similar weight vectors. The labels are a re-
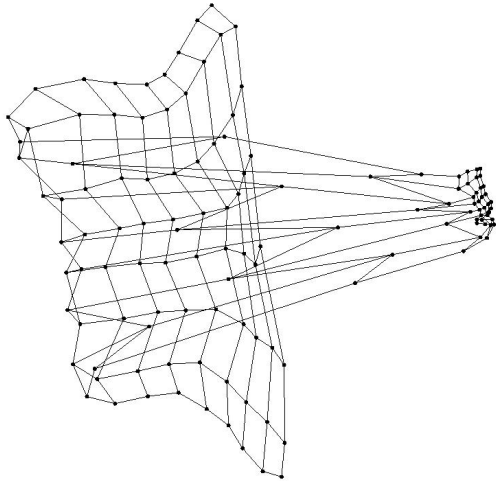
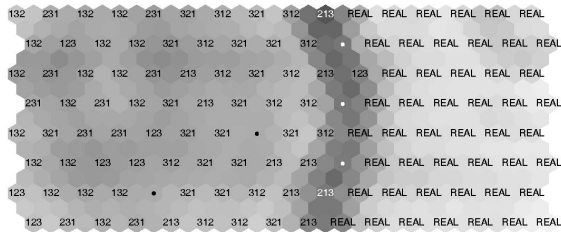**Figure 2. Sammon projection of the** $16 \times 8$ **map of figure 3: random versus real melodies.**



**Figure 3. SOM map for the same weights as in figure 2.**

sult of calibrating the map with a series of test samples and indicate the class of the sample that activates that unit more times.

### 3.1 Random versus jazz melodies

400 random samples have been generated and 430 jazz samples have been extracted from 54 MIDI sequences of jazz standards, all of them made up of 8 bars with a quantisation of $Q = 8$ pulses per bar (64 events per melody). From them, the 20 descriptors listed above were computed. Using these sets a SOM of 16 neurons for the $OX$ axis and 8 for the $OY$ axis was trained. The training consisted of two stages: a coarse one of 1000 iterations with wide neighbourhoods (12 units) and a high learning rate (0.1) and then a fine one of 10,000 iterations with smaller neighbourhood ratio (4 units) and learning rate (0.05). These training parameters can be applied to the rest of experiments with little variations.

In figures 2 and 3 the Sammon projection and the SOM map are displayed after training for that experiment. Note
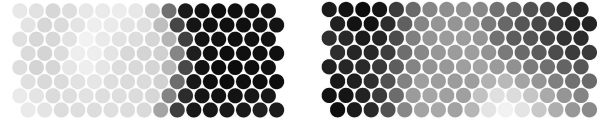


**Figure 4. Contributions to discrimination: (left) plane of the weight space for average intervals correlate with map areas; (right) plane for minimum silence durations does not.**

that there exists a clear gap between two zones in the map. The small cluster on the right in the Sammon projection corresponds to the real melodies and that of the left to the random melodies. In the map the same can be observed for the two areas clearly separated: random samples on the right and real samples in the left. The dark strip represents the separation between both zones. The SOM has been labelled using the training samples. The "REAL" cluster has less extension than that of random samples (labelled according to the event proportions), because the latter have more variability. There was an almost total lack of overlapping (units labelled with both styles) between the zones.

It is clear that the distinction between both zones in the map corresponds to real differences between the random and jazz melodies, and the SOM has been able to capture those differences.

The descriptors that contribute more to that separation are those having a higher correlation among the samples in each of the zones. The planes in the weight space (see Fig. 4) corresponding to each descriptor provide information about this. An analysis of the most contributive features is an indication about how the discrimination has been carried out:

- *Maximum and minimum pitches*. In the random melodies, extremal note pitches appear more often than in real melodies.

- *Standard deviation of the pitches*. Is clearly larger for random melodies, due to the lack of a 'tonal centre' that acts like an attractor for the melodic line.

- *Maximum and average interval*. Much higher in random samples. Intervals higher than an octave are seldom present in real melodies. The average interval is usually between 2 and 3 for real melodies and higher than 10 for random ones.

### 3.2 Jazz versus classical music

We have seen that, when the melodies are clearly different, the use of adequate descriptors makes the discrimination problem an easy task using SOM (and probably with

**Figure 5. SOM map trained with jazz and classical melodies.**

other statistical classifiers). Once tested the ability of SOM we will substitute other real melodies for the random samples. The new set is composed of monophonic fragments of classical music. In order to have more resolution in the sequences a higher quantisation of 48 pulses by bar has been used.

522 classical music melody fragments of eight bars of length were extracted from MIDI files for the training set along with the previous 430 jazz samples, now quantised to $Q = 48$ pulses by bar.

After training and labelling the SOM it can be observed in Fig. 5 that the "JAZZ" labels are denser on the left side of the map and becomes sparser on the right, and the contrary happens with the labels of the classical composers. Note that some units are labelled for both music styles because they were activated by samples from both styles. In these cases there is always a winner label (we call it the *first label* according to the number of activations and a loser (*second*) label. The proportion of units with both labels is the overlapping degree, that in this experiment is rather high, a 39.0% of the units were labelled for both styles. This fact suggests that differences were detected between both styles but maybe there is a lack of information to take decisions. For this, harmonic features were added to the set of 18 descriptors already considered.

### 3.2.1 Addition of harmonic descriptors

Most of western music is based in a number of scales (sets of notes ordered by pitch) and melodies can be formed taking notes from those sets. A *diatonic* melody is made up with the natural notes of a scale, without sharp or flat notes (named *accidentals*). In western music most of the melodies belong to one of two main scale types: major or minor

scales, and people usually know how they sound in an intuitive way. The first note of a scale determines its 'tonality' or 'key' and in any melody diatonic and accidental notes can appear.

If the overall key and kind of scale (major or minor) of a melody are known, the set of diatonic pitches is also known and any note event can be classified into diatonic or accidental, and some harmonic information can be evaluated, like the proportion of diatonic notes with respect to the total. If the proportion is high then it is an indication of small key changes or modulations, if any. On the other hand, a low proportion indicates that there are a lot of key changes.

The detection of the key and the definition of the diatonic scale utilised is based on musicological criteria and their description is outside the scope of this paper.

We number the accidental notes of a given scale from 1 to 5 according to their distance in pitch from the key note of the scale. We will call this the *accidental degree*. According to this criterion, three harmonic descriptors are defined:

- *Number of accidental notes.* An indication of frequent excursions outside tonality or modulations.

- *Average degree of accidental notes.* Describes the kind of excursions.

- *Standard deviation of degrees of accidental notes.* Indicates a higher variety in the modulations. If this value is close to 0 the accidentals are probably caused by chromatic approximations or adorns rather than real harmonic modulations or key changes.

From the MIDI file the key is extracted and the diatonic notes determined. Then, the harmonic descriptors are computed.

A new experiment is designed using the 21 descriptors already defined. 522 samples of 8 bars of classical music have been utilised, 430 for training and the rest for test, and also 428 samples of the same length as jazz.

The size of the map has been also increased according to the higher dimensionality of the input vectors. The number of neurons is now $30 \times 12$. The neighbourhood radius has also been adapted to the new dimensions, having a radius of 20 units for the first training phase (coarse) and 6 units for the second (fine). Also, the duration of these phases is now ten times longer than where the SOM where smaller. After training and labelling, the maps in figure 6 have been obtained. It is observed how the labelling process has located the "JAZZ" labels mainly in the right and upper zone, and those corresponding to classical composers mainly in the lower left zone. The percentage of overlapping in this experiment was in this case very low: 11.1%. Now a clear distinction of styles has been achieved.

**Figure 6. SOM map after being labeled with jazz (top) and classical (down) melodies.**



**Figure 7. SOM after training using also rhythmic descriptors: (top) units labelled with "JAZZ" and (bottom) units labelled with "CLAS".**

### 3.2.2 Addition of rhythmic descriptors

One of the features of certain styles (like jazz) is the abundance of syncopation: notes not beginning in the rhythm beats but in some places between them (usually in the middle). Beats are the main pulses of a bar and the syncopation provides a very particular feeling to music. On the other hand, most of classical styles do not make an extensive use of this feature. Therefore the quantification of these features are a priori interesting for melody description.

In this experiment the number of syncopations is added as a rhythmic descriptor to those already defined, having therefore 22 descriptors. Syncopation detection is not a trivial task if quantisations over $Q = 16$ are utilised. For that, a syncopation range $q_{sync}$, measured in pulses and dependent on $Q$, is defined. If a note starts within that range around the middle of a beat, then it is considered as a syncopation. This descriptor is a measurement of the *degree of syncopation* rather than a re-count of the number of such events.

Again, the same map is trained with the same training set. And similar maps and label distributions are obtained, as seen in Fig. 7.

When analysing the labelled map in Fig. 7 it is observed that the left side is more "jazzy" and the right one is more classical, corresponding to both sides with respect to the knot, but there are a lot of mixed labels, making it harder to determine zones that could be assigned clearly to each genre. This is probably due to the fact that there are some jazz melodies that seem classical music at times and vice versa. In spite of that, less overlapping has been found (7,22%).

In the Sammon projection of figure 8 a knot separates two zones in the map. The zone at the left of the knot has a majority presence of units labelled with the jazz label and the zone at the right is mainly classical. With our data, con-
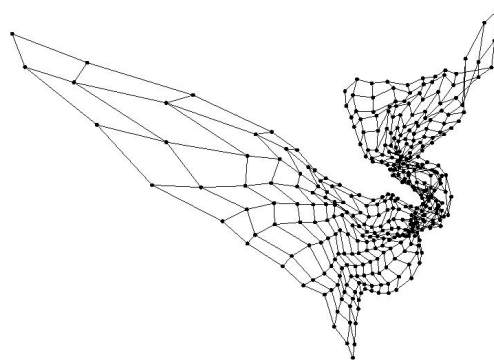


**Figure 8. Sammon projection of the SOM map in figure 7.**

sistently, a higher number of classical units appear in the "jazz zone" than jazz units among the classical ones.

### 3.3 Classification

The results of classifying new melodic fragments, not contained in the training sets, using the different SOM described above are presented in table 1. We have focused in the classification achieved using melodic and harmonic descriptors. The results in the table are those obtained in the next experiments:

- Classification results with the small ($16 \times 8$) SOM using melodic and harmonic descriptors.

- Classification results with the large SOM ($30 \times 12$) using melodic and harmonic descriptors.

**Table 1. Classification results (percentages) using melodic and harmonic descriptors**

|                | JAZZ | CLASSICAL |
|----------------|------|-----------|
| Map dimensions = $16 \times 8$ | | |
| Class success  | 76.9 | 77.5 |
| Class error    | 23.1 | 22.3 |
| Unclassified   | 0.0  | 0.2  |
| One label      | 24.9 | 52.0 |
| First label    | 52.0 | 25.5 |
| Second label   | 14.1 | 18.6 |
| Error          | 9.0  | 3.7  |
| Map dimensions = $30 \times 12$ | | |
| Class success  | 69.8 | 69.4 |
| Class error    | 23.2 | 18.3 |
| Unclassified   | 6.9  | 12.3 |
| One label      | 51.2 | 61.4 |
| First label    | 18.6 | 8.0  |
| Second label   | 6.3  | 12.5 |
| Error          | 16.9 | 5.8  |

For obtaining reliable results a scheme based on *leave-k-out* has been carried out. In our case $k = 10\%$ of the size of the whole database. This way, 10 sub-experiments were performed for each experiment and the results have been averaged. In each experiment the training set was made of a different 90% of the total database and the other 10% was kept for testing.

The data presented in the table are as follows: *One label success* indicates the proportion of melodies to which the map has assigned a unit containing just one label and it was the right one. *First label* success percentage is related to melodies assigned to units with two labels but being the first the right one. Therefore, we are considering as favourable decisions these two criteria: just one correct label or two labels but the first is the correct one (the sum of these two quantities gives the *class success* row).

*Second label* indicates the proportion of times that a melodic fragment has been asigned to a unit that had the correct label as the second one. This fact is considered as a misclassification. *Error* means that the map has assigned a unit containing just one but wrong label. These two answers define the *class error*. Finally, *Unclassified* melodies were those assigned to a unit not containing any label.

The best overall performance was obtained with the smaller map, with a success classification rate of 76.9% for jazz melodies and of 77.5% for classical melodies. On the other hand the error rates are lower for the second map, and the difference is due to the higher *unclassified* rates for this second map. These results are probably due to the fact

that in the second map the class clusters were more defined, leaving more space to unlabelled units. If we devise a way to assign this unlabelled units a class (based, for example, in taking into account the distances in the weight space for the trained map for assigning a label also to unlabelled units), probably the results would improve.

Note also that the rates for *one label* classification are higher for the larger map. This is due to the lower overlapping rate for this map reported above.

## 4   Conclusions and future works

We have shown the ability of SOM to map symbolic representations of melodies into a set of musical styles using their description in terms of melodic, harmonic and rhythmic features. The best recognition rate has been found with 21 descriptors (the basic set of 18 and 3 about harmony). Rhythmic features like syncopation descriptors have not increased the performance of the system. The best recognition rate has not been achieved when the overlap was minimum, so the overlap ratio does not seem to be a key point when assessing the quality of a map.

Some of the misclassifications can be caused by the lack of a smart method for melody segmentation. The music samples have been arbitrarily restricted to 8 bars, getting just fragments with no relation to musical motives. This fact can introduce artifacts in the descriptors leading to less quality mappings. The main goal was to test the feasibility of the approach, dealing even with incomplete data. Nevertheless a best total recognition rate of 82,9% has been achieved, that is very encouraging keeping in mind these limitations and others like the lack of valuable information for this task like timbre.

A number of possibilities are yet to be explored, like the development and study of new descriptors. A statistical multifactorial study of the whole set of descriptors and other feature selection schemes can aid in the selection of a model that can achieve better results with a minimum subset of them. It is very likely that this subsets is highly dependent to the styles to be discriminated.

To achieve this goal a large music database has to be compiled and tested using our system. Different styles and more melodies are needed to draw significative conclusions. We are working now in a database based on XML representations and descriptions of music, providing more musically meaning information than MIDI files, in a higher level of abstraction, such as hierarchical information about note grouping, note modifiers that modify pitch, duration or dynamics.

Other future lines are based in the integration of time in the description process to capture the evolution of the whole melody. The map activations for a series of fragments of the same melody could be the input to other recognition algo-

rithms in order to increase the classification power of the system, even with a higher number of music styles at the same time. We are working now in the compilation of a large multistyle database for further experimentation.

The activations of the map can be the input to a recognition neural layer that under a supervised learning can be utilised in cooperation to melody generation system to build an automatic composition system specialised in certain musical styles. Works in that direction are currently being developed.

## Acknowledgements

## References

[1] O. A. S. Carpinteiro, 'A self-organizing map model for analysis of musical time series', in *Proceedings 5th Brazilian Symposium on Neural Networks*, eds., A. de Padua Braga and T. B. Ludermir, 140–5, IEEE Comput. Soc, Los Alamitos, CA, USA, (1998).

[2] T. Kohonen, 'Self-organizing map', *Proceedings IEEE*, **78**(9), 1464–1480, (1990).

[3] T. Kohonen, J. Hynninen, J. Kangas, and J. Laaksonen. Som_pak, the self-organizing map program package, v:3.1. Lab. of Computer and Information Science, Helsinki University of Technology, Finland, April, 1995.

[4] A. Rauber and M. Frühwirth, *Automatically analyzing and organizing music archives*, 4–8, 5th European Conference on Research and Advanced Technology for Digital Libraries (ECDL 2001), Springer, Darmstadt, Sep 2001.

[5] Hagen Soltau, Tanja Schultz, Martin Westphal, and Alex Waibel, 'Recognition of music types', in *Proceedings of the IEEE International Conference on Acoustics, Speech, and Signal Processing (ICASSP-1998)*, Seattle, Washington, (May 1998).

[6] Belinda Thom, 'Unsupervised learning and interactive jazz/blues improvisation', in *Proceedings of the AAAI2000*, pp. 652–657, (2000).

[7] Brian Whitman, Gary Flake, and Steve Lawrence, 'Artist detection in music with minnowmatch', in *Proceedings of the 2001 IEEE Workshop on Neural Networks for Signal Processing*, 559–568, Falmouth, Massachusetts, (September 10–12 2001).