

Contents lists available at [ScienceDirect](http://www.sciencedirect.com)

# Pattern Recognition Letters

journal homepage: [www.elsevier.com/locate/patrec](http://www.elsevier.com/locate/patrec)

## Explicit modeling of temporal dynamics within musical signals for acoustical unit similarity

Mathieu Lagrange<sup>a,\*</sup>, Martin Raspaud<sup>b</sup>, Roland Badeau<sup>a</sup>, Gaël Richard<sup>a</sup>

<sup>a</sup> Institut Telecom, Telecom ParisTech, CNRS LTCI, 46 rue Barrault, 75634 Paris, Cedex 13, France

<sup>b</sup> Linköping University, Bredgatan 33, SE-60174 Norrköping, Sweden

### ARTICLE INFO

#### Article history:

Available online xxx

#### Keywords:

Audio similarity  
Timbre modeling  
Audio analysis  
Temporal dynamics

### ABSTRACT

Timbre is a major cue for the human auditory system to recognize musical sounds. To describe timbre, the temporal dynamics is an important component as well as the widely used spectral envelope.

In this paper, we present new temporal dynamics similarity measures, which will prove valuable for the recognition of timbral patterns. These similarity measures are evaluated, first alone, then in conjunction with spectral envelope similarity measures, for both single tones and solo recordings. Results are provided, showing that the new temporal dynamics features significantly improve timbral pattern recognition.

© 2009 Elsevier B.V. All rights reserved.

### 1. Introduction

The timbre, along with the perceived loudness, pitch, and duration, is an important perceptual attribute of sound. As reported in (McAdams et al., 1993), contemporary research in psychological and cognitive acoustics decomposes this attribute into several perceptual dimensions of temporal, spectral, and spectro-temporal nature.

It is worth noticing that in the audio signal processing area, the spectral nature of timbre has received much more interest than the others. One of the best examples is the set of features called Mel-Frequency Cepstral Coefficients (MFCCs) introduced by Davis and Mermelstein (1980). Those features are widely used in speech and speaker recognition systems as well as in music classifications system such as the one proposed by Tzanetakis and Cook (2002). From a modeling point of view, the spectral envelope is related to the filter part of a source/filter model of the analyzed sound as proposed for speech by Fant (1960). In many cases, the spectral properties of this filter are specific to the vibrating body, *i.e.* the vocal tract or the shape of a musical instrument. This makes the modeling of the spectral envelope particularly interesting for the description of musical sounds. From a technical point of view, the spectral envelope can easily be extracted in a frame-based manner with minimal delay.

However, it is widely known that the temporal dimension of timbre is very important at least from a perceptual point of view as reported in (Grey and Moorer, 1977). As detailed in Section 2, the temporal dynamics of timbre are implicitly modeled by the

frame-to-frame variability of the spectral features. The processing of temporal dynamics thus has to be performed subsequently. One can consider another feature that characterizes the variation over time of a previously computed feature, such as the  $\Delta$ MFCCs. However, recent work of Joder et al. (2009) shows that the use of such differentiators is not sufficiently informative to improve the classification performance. By contrast, we propose in this paper to model explicitly the temporal dynamics of the spectral parameters in order to build spectro-temporal features. The rationale behind the proposed approach is that, compared to feature-level dynamic modeling, the proposed features take into account finer spectro-temporal modulations, which are useful to characterize the audio scene.

In this paper, the spectro-temporal features are used – through the definition of similarity metrics – to discriminate audio signals produced by different musical instruments. The discrimination performance of the similarity metric in turn gives us an evaluation of the ability of the spectro-temporal features to describe the analyzed sound in a meaningful way.

The remainder of the paper is organized as follows: related work is reviewed and discussed in Section 2. The proposed approach is motivated in Section 3 to propose several features that model the temporal dynamics. Those features are considered for the definition of similarity metrics between audio signals in Section 4. The proposed metrics are next evaluated in Section 5. In light of those experiments, the benefits of spectro-temporal features proposed in this paper are discussed in Section 6.

### 2. Previous work

Sinusoidal modeling is one of the first attempts to model explicitly the temporal dynamics of sound. The sinusoidal model

\* Corresponding author. Tel.: +33 1 45 81 73 24; fax: +33 1 45 81 71 44.  
E-mail address: [lagrange@enst.fr](mailto:lagrange@enst.fr) (M. Lagrange).

represents pseudo-periodic sounds as sums of sinusoidal components – so-called partials – controlled by parameters that evolve slowly with time as considered by McAulay and Quatieri (1986) and Serra and Smith (1990). More formally put, the audio signal  $s$  can be calculated from the controlling parameters using Eqs. (1) and (2), where  $N$  is the number of partials and the functions  $f_p$ ,  $a_p$ , and  $\phi_p$  are the instantaneous frequency, amplitude, and phase of the  $p$ th partial, respectively. The  $N$  pairs  $(f_p, a_p)$  are the parameters of the additive model and represent points in the frequency–amplitude plane at time  $t$ :

$$s(t) = \sum_{p=1}^N a_p(t) \cos(\phi_p(t)) \quad (1)$$

$$\phi_p(t) = \phi_p(0) + 2\pi \int_0^t f_p(u) du \quad (2)$$

This can also be written from the set point of view:

$$p^k(n) = \{f^k(n), a^k(n), \phi^k(n)\} \quad (3)$$

where  $f^k(n)$ ,  $a^k(n)$ , and  $\phi^k(n)$  are, respectively, the frequency, amplitude, and phase of the partial  $p^k$  at frame index  $n$ . These parameters are valid for all  $n \in [b^k, \dots, b^k + l^k - 1]$ , where the  $b^k$  and  $l^k$  are, respectively, the starting index and the length of the partial. These sinusoidal components are called *partials* because they are only a part of a more perceptively coherent entity that will be noted in this article an *acoustical unit*.

### 2.1. The common variation cue

When a sound-generating object changes its properties so that its fundamental frequency gets higher or lower, all the partials of the sound also change synchronously. In several experiments McAulay (1989) studied the influence of this phenomenon in the auditory system. Whether the common variation cue is an important cue for the fusion and segregation capacity of the human auditory system is still an open issue as reported by McAulay et al. (1993).

However, from a physical point of view, this phenomenon can be measured and can therefore be used to perform the clustering

of the partials of the same acoustical unit, as proposed by Lagrange (2005). Let us consider the case of a harmonic set of partials modulated by a vibrato. The frequencies  $f^k(n)$  are periodically modulated at the same rate, and the depth of the vibrato is a function of the rank of the partial in the harmonic set. An extra care should be taken while considering the induced modulation of the amplitude. Indeed, depending on the sign of the spectral envelope slope at the partial frequency location, the modulation phase can be shifted. This phenomenon is illustrated by Fig. 1 where the lowest amplitude partial has its amplitude modulated at the same rate but with a  $\pi/2$  delay.

### 2.2. Integration of Frequency-Axis Features

As with most model-based approaches, the sinusoidal model tends to be brittle when applied to real-world sounds as studied by Lagrange (2004). Therefore, most practical approaches are based on spectral Fourier representations computed in frame-based manner and summarized by numerous means, one of the most famous being the MFCCs. The MFCCs are coefficients that describe the short-term power spectrum of a sound, based on a linear cosine transform of a log power spectrum on a nonlinear Mel scale of frequency. By selecting the first coefficients, one can estimate a “smooth” version of the spectrum, usually considered as an approximate of the spectral envelope. This kind of feature is a static observation of the spectral content of the signal and will therefore be termed a *Frequency-Axis Feature* (FAF).

The temporal aspect of the analyzed sound is not completely neglected by the frame-based approach. Indeed, the temporal dynamics are in this case implicitly encoded by the frame-to-frame variability of the frame-based features. This variability can potentially be captured at later stages by the following two approaches previously studied by Aucouturier and Pachet (2007).

The first one, known as *classifier* or *late integration*, does not try to explicitly extract feature dynamics, but rather operates at the classifier level, usually a supervised classifier with sequentiality constraints, like Hidden Markov Models considered in (Eronen, 2003; Kitahara et al., 2006) or Sequence kernel-based Support Vector Machines considered in (Scaringella and Zoia, 2005; Joder et al.,

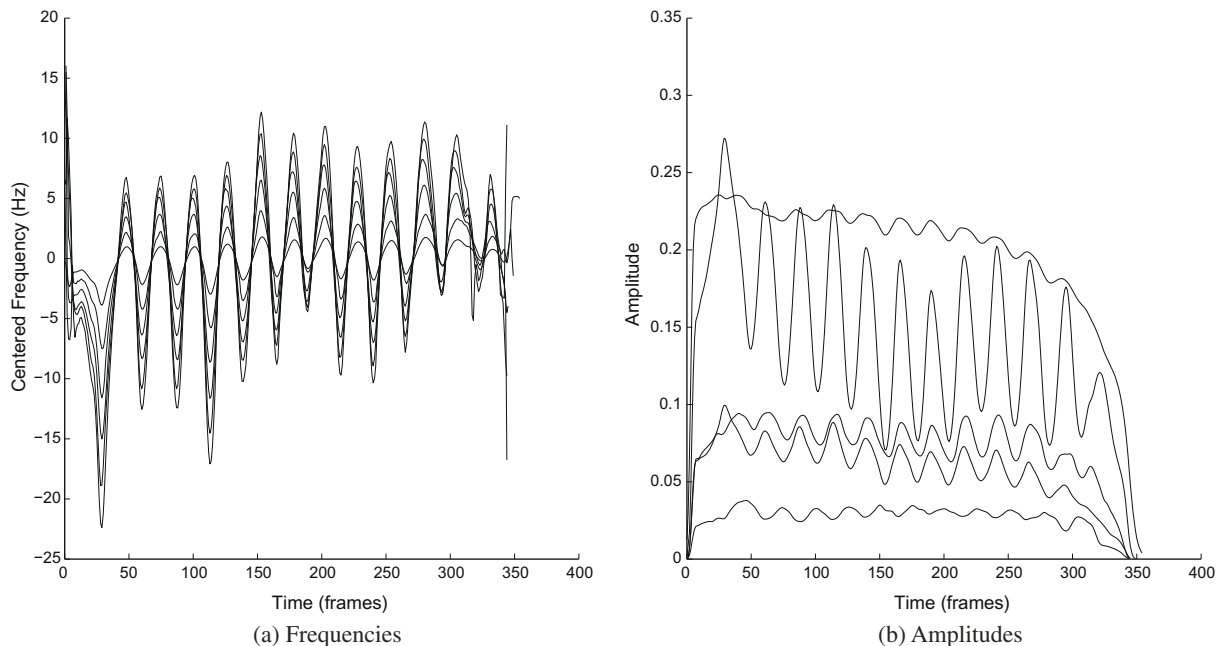


Fig. 1. Mean-centered frequencies and amplitudes of some partials of a saxophone tone with vibrato.

2009; Shimodaira et al., 2002; Cuturi et al., 2007). As these techniques do not deal with explicit modeling of the temporal dynamics of the sound, they will not be discussed further in this paper.

A second approach, known as *feature-level integration*, refers to the computation of a new feature vector that characterizes the evolution of a given set of features at a larger time scale. The most commonly used feature-level integration is the  $\Delta$ MFCs. Meng et al. (2007) have studied more complex models, like high-order auto-regressive models for genre classification. However, Joder et al. (2009) reported that the use of this kind of feature integration actually degrades the classification performance in an instrument recognition task. The poor performance may be due to the nature of FAFs which tends to smooth away potentially meaningful information about the temporal dynamics of the sound.

We introduce in the next section an alternative approach which aims at modeling the temporal dynamics of the spectrum and is consequently termed *spectral-level integration*.

### 3. Proposed features

The spectral envelope is an important piece of information if one wants to characterize the timbre of an audio signal. However, we believe that this should be complemented by *Time-Axis Features* (TAFs) that explicitly model the temporal dynamics of the spectral components of the sound.

Let us consider a simplified version of the source/filter model in order to better motivate our approach. The FAFs mainly model the filter part of the sound production chain. In order to complement these features, it is important to avoid any redundancy and therefore focus on a different aspect of the sound production chain. Following the source/filter dichotomy, we propose to root the TAFs on the source part.

#### 3.1. Frequency Evolution Features

We design Frequency Evolution Features (FEFs) such that they encode the modulation of the frequency of the main components of the spectrum over time. It is therefore natural to choose the sinusoidal model as a signal representation.

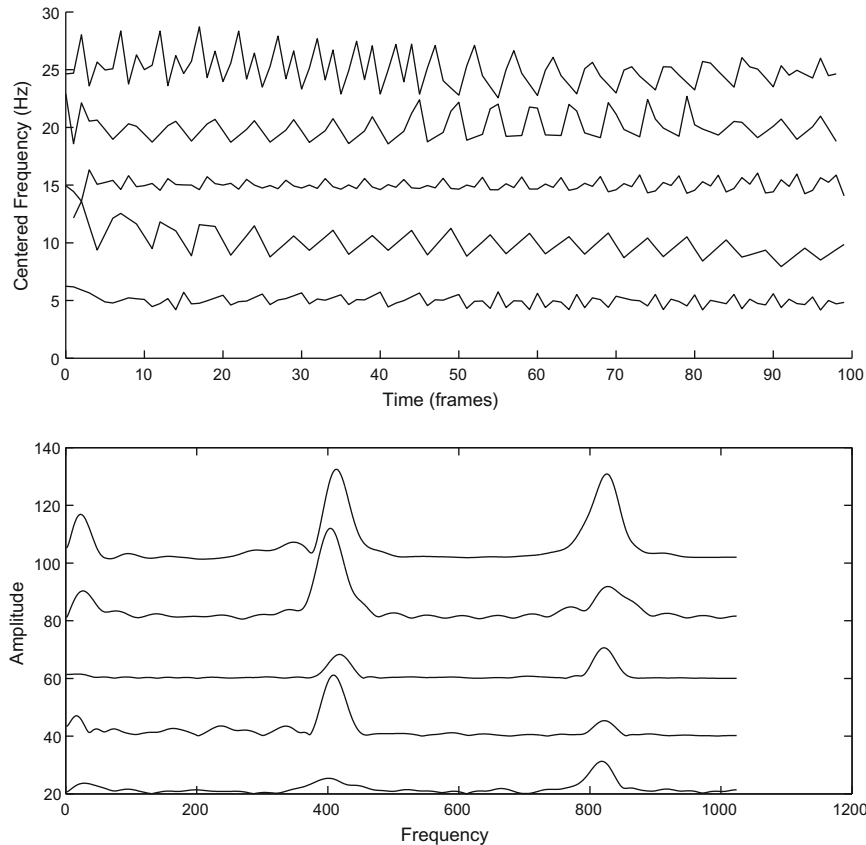
In a first approach, the temporal evolution of the frequency of the partials  $f^k$  can be considered directly. Alternatively, since the Fourier transform is based on the periodicity of the input signal, using a spectrum of the evolutions of partial parameters might show common periodicities of the partials. This will be useful for the modulations of the partials created by vibrato and tremolo, since we can assimilate these modulations to sinusoidal ones over a short period of time as studied by Mellody and Wakefield (2000) and Marchand and Raspaud (2004). It can also be interesting for micro-modulations such as the ones produced by vibrating strings such as the strings of a piano (see Fig. 2).

Let us define the following operator, based on the complex modulus of the short-time Fourier transform:

$$X(k) = \left| \sum_{n=0}^{N-1} (x(n) - \bar{x}) h(n) e^{-\frac{2j}{N}\pi kn} \right|^2 \quad (4)$$

$$|\mathcal{X}|_{n_l}^{n_h}(x(n)) = \{X(m) | n_l < m < n_h\} \quad (5)$$

where  $N$  is the size of the Fourier transform,  $h$  is an analysis window, and  $\bar{x}$  denotes the mean of  $x$ .  $n_l$  and  $n_h$  are, respectively, the minimal and maximal frequency indexes considered. Thanks to the complex modulus applied to the spectrum, this operator is phase-invariant. We can then define the *frequency evolution features* (FEFs), as:



**Fig. 2.** Centered frequencies (top) of a piano note from the IOWA database and their corresponding spectra (bottom). The spectra are interpolated using zero-padding for clarity sake.

$$F^k = |\mathcal{X}|_1^{N-1}(f^k) \quad (6)$$

$$Fc^k = |\mathcal{X}|_1^{n_c}(f^k) \quad (7)$$

where  $f^k$  is the frequency of the partial  $k$  and  $n_c$  is chosen in order to reduce the dimensionality of the feature while keeping low frequency content information.

### 3.2. Amplitude Evolution Features

In order to encode the temporal dynamics of the amplitude, we can in a similar fashion consider the evolution of the amplitude  $a^k$  of the partials  $p^k$  or the corresponding spectral features, called *Amplitude Evolution Features* (AEFs):

$$A^k = |\mathcal{X}|_1^{N-1}(a^k) \quad (8)$$

$$Ac^k = |\mathcal{X}|_1^{n_c}(a^k) \quad (9)$$

In order to consider only the modulated part of the amplitude signal, leaving aside the global envelope, it is relevant to decompose the signal in two components, one being polynomial, and the other being pseudo-periodic as proposed by Raspaud et al. (2005):

$$a^k(t) = \Pi(t) + \sum_i \alpha_i(t) \cos(\psi_i(t)) \quad (10)$$

where  $\Pi(t)$  is a polynomial, and  $\alpha_i(t)$  and  $\psi_i(t)$  are the parameters of sinusoidal components, see Fig. 3a.

Indeed, while subtracting the mean of the signal – as performed by the operator  $|\mathcal{X}|$  – is enough to center the oscillations of the evolution of the frequency of partials, it is not the case for the evolution of the amplitude of partials. As studied by Raspaud (2007), the

idea behind this polynomial subtraction is that the envelope of a sound (seen as attack, decay, sustain and release) can be approximated by a 9th degree polynomial. Then, we define the following two other AEFs based solely on the oscillating part of the partials amplitude:

$$ap^k = a^k - \tilde{\Pi}(a^k) \quad (11)$$

$$Ap^k = |\mathcal{X}|_1^{N-1}(ap^k) \quad (12)$$

$$Apc^k = |\mathcal{X}|_1^{n_c}(ap^k) \quad (13)$$

where  $\tilde{\Pi}(x)$  is the envelope polynomial computed from signal  $x$  using a simple least-squares method. As an approximation of the polynomial removal, one can consider removing the DC component in the following way:

$$Ad^k = |\mathcal{X}|_{n_d}^{N-1}(ap^k) \quad (14)$$

$$Adc^k = |\mathcal{X}|_{n_d}^{n_c}(ap^k) \quad (15)$$

where  $n_d$  is chosen so that periodicities like the tremolo are preserved, see Fig. 3b.

### 3.3. Magnitude Evolution Features

As stated previously, the sinusoidal model provides a meaningful representation for analyzing the temporal dynamics. However, the estimation of this model from complex signals can hardly be done in a fully automatic fashion. As an approximation of the AEFs, we also consider *Magnitude Evolution Features* (MEFs) that rely on the spectrogram only. Taking into account the evolution of the magnitude in the spectrogram is a non parametric way to account for the temporal dynamics, *i.e.* without relying on the sinusoidal model.

Let us consider  $X(k, n)$  the spectral bin  $k$  of the frame  $n$  of the spectrogram of the signal  $x$ . For a given spectral bin of frequency index  $k$ , the MEFs correspond to the magnitude evolution of a frequency line in the spectrogram:

$$m^k = X(k, n) \quad (16)$$

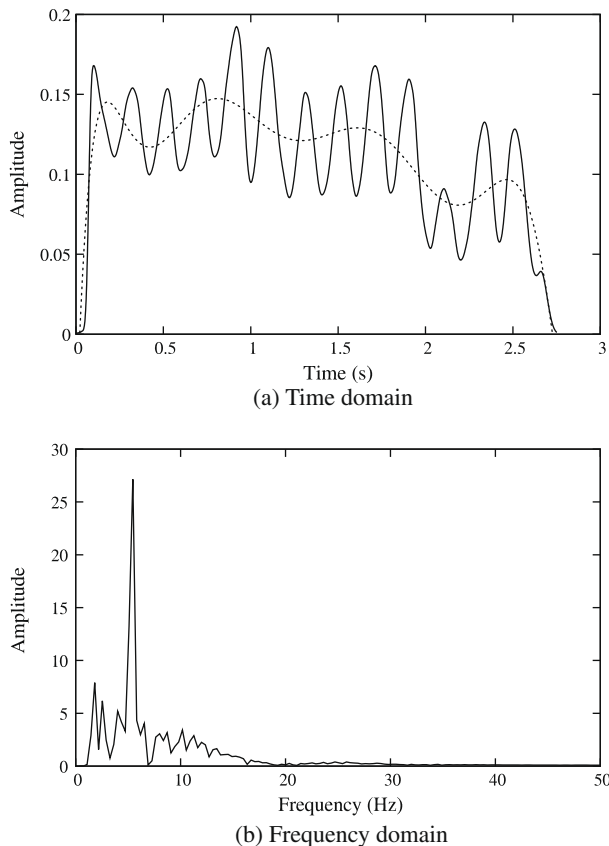
where  $n$  varies within a given horizon of observation. The other MEFs  $M^k$ ,  $Mc^k$ ,  $Mp^k$ ,  $Mpc^k$ ,  $Md^k$  and  $Mdc^k$  are computed as described in the previous section.

## 4. Acoustical units similarity

We evaluate the proposed features by their ability to express the similarity between acoustical units. We define an acoustical unit to be a musical tone or a sequence of musical tones performed by a unique musical instrument, within a limited time interval. The task is then to decide whether two acoustical units have been played by the same musical instrument or not. This decision is made according to the information given by the features computed from the acoustical unit. This evaluation task is chosen for two reasons.

From a practical application point of view, there is an increase of interest towards recommendation systems that are not based on an ontology such as genre as used by Tzanetakis and Cook (2002) or instrument type considered by Martin and Kim (1998). Alternatively, one can consider a recommendation system that states “show me tunes that are comparable to the ones I like”. In this case, one needs to define the similarity between musical audio signals. For the definition of such a similarity, the timbre is an interesting dimension.

From a scientific point of view, this allows us to propose a much simpler evaluation framework than the one required by classification-based systems (such as the previously cited genre and instru-



**Fig. 3.** (a) Amplitude of a partial, its estimated polynomial envelope (dotted line) and (b) the corresponding frequency domain representation of the polynomial-removed amplitude evolution.

ment type). Indeed, the latter rely on training complex classifiers which may have constraints over the statistical properties of the features and can consequently introduce a bias over the evaluated features.

#### 4.1. Features integration

One issue with the proposed scheme is that the dimensionality of the problem is the square of the number of elements to be sorted. Consequently, we are interested in efficiently describing longer acoustical units than those that are usually considered in frame-based classification systems. The experiments reported in this paper will demonstrate the usefulness of the TAFs in such a case.

When considering isolated notes, the acoustical unit duration is adapted to the actual duration of the note. When considering solos excerpt, the audio signal is arbitrarily segmented at a regular sampling rate. Two fixed integration sizes are considered in the experiments. The first one, termed “texture”, considers 20 frames, each of 22 ms. Consequently, the texture duration is 440 ms. The second one, termed “event”, is 1 s long. For the largest database considered in this article, the number of acoustical units is 0.2 million, leading to a number of similarities that have to be computed to about 50 billions which took about a month to compute on a state-of-the-art computer.

Within an acoustical unit, the features are extracted from a spectrogram computed in the following way:

1. The input signal is first filtered with a DC block filter and a pre-emphasis filter.
2. The spectra are next computed with a frame size of 40 ms and a hop size of 10 ms.

Concerning the FAFs features, the MFCCs are computed thanks to the Matlab implementation of Ellis (2005) within each analysis frame with 40 Mel sub-bands and only the first 12 coefficients after the zeroth coefficient are considered. The  $\Delta$ MFCCs are computed using a 5-point derivation filter. Finally, those coefficients are summed over the integration interval, as proposed in (Joder et al., 2009).

The TAFs features are extracted within each acoustical unit in the following way:

1. The spectrogram is first integrated over time and split in 20 Mel sub-bands, leading to the selection of a maximum of 20 sinusoidal components.
2. Within each sub-band, the bin with maximal amplitude is selected, indicating the potential frequency location of a partial.

The magnitude evolution within each of those bins is considered for the computation of the MEFs. In order to extract the FEFs and the AEFs, partials are tracked in a frequency interval around each of the maximal amplitude bins using a standard partial tracker proposed by Ellis (2003). The width of this interval is set to  $\approx 220$  Hz. The partial with maximal cumulative amplitude within each frequency interval is then selected to compute the TAFs.

The TAFs are associated to 20 partials belonging to an acoustical unit. However, in this task, we need to define the similarity between acoustical units. Consequently, the TAFs describing each partial of a given acoustical unit are integrated by summation:

$$f(n) = \frac{1}{N} \sum_{k=1}^N \tilde{f}^k(n) \quad (17)$$

where  $\tilde{f}^k(n)$  is the resampling version of  $f^k(n)$  of length 20. The features  $a$  and  $m$  are defined similarly. An equivalent operation is carried out for spectral-based features:

$$F(n) = \frac{1}{N} \sum_{k=1}^N \log_{10}(F^k(n)) \quad (18)$$

The other spectral-based features are computed similarly with the following parametrization:  $N = 64$ ,  $n_c = 16$  and  $n_d = 2$ .

#### 4.2. Similarity metrics

The features described above are next considered to build similarity metrics after being normalized to have zero mean and unity variance. In order to demonstrate the ability of the proposed TAFs features to complement a FAF feature like the MFCCs, the combination with the MFCCs is also considered. For that purpose, the radial basis function is considered:

$$s(U_i(k), U_j(l)) = e^{-\|V_i(k) - V_j(l)\|} \quad (19)$$

where  $\|x\|$  is the Euclidean norm of  $x$  and  $V_i(k)$  is a feature vector or a stacking of feature vectors describing the acoustical unit  $U_i(k)$  played by the instrument  $i$ . In order to account for the discrepancy between feature dimensionality, the features are divided by the square root of their dimensions prior to stacking.

### 5. Experiments

In this section, the similarity metrics described in the previous section are considered to evaluate the potential of the proposed features. The evaluation metrics are first introduced and two evaluation scenarios are studied. The first one considers databases of isolated notes and the second one continuous musical solos excerpts.

#### 5.1. Evaluation metrics

In order to evaluate a set of features describing an acoustical unit, the value of the similarity metric is first computed for every pair of acoustical units. To compute performance indicators, we propose to consider a classifying task which corresponds to “Are those acoustical units played by the same instrument?”. As a consequence, at a given classifying threshold  $T_c$ , we can define the *False-Alarm Rate* (FAR) as:

$$\text{FAR} = \frac{\#\{s(U_i(k), U_j(l)) > T_c\}}{\#U \cdot (\#U - 1)} \quad \text{for } i \neq j \quad (20)$$

where  $\#X$  denotes the cardinal of  $X$  and  $U$  is the overall set of acoustical units in the evaluation database. Similarly, the *Miss Detection Rate* (MDR) is defined as:

$$\text{MDR} = \frac{\#\{s(U_i(k), U_i(l)) < T_c\}}{\sum_i \#U_i \cdot (\#U_i - 1)} \quad (21)$$

The *Detection Error Trade-Off* (DET) curve proposed by Martin et al. (1997) is used to visualize the performance of the classifier corresponding to the evaluated feature or combination of features at a varying classification threshold, see Fig. 4. In order to summarize the behavior of the classifier depending on the chosen threshold, we consider three criteria. The first one called *Equal Error Rate* (EER) corresponds to the crossing of the DET curve with a line that starts from (0,0) coordinates with unity slope. The second one called *Minimum Cost Point* (MCP) indicates the performance of the classifier for an optimal error trade-off between a balanced weighting of the MDR and FAR as described in (Martin et al., 1997). The MCP for each feature is plotted in Fig. 4 with a small circle. The last criterion, called AREA, computes the area under the curve and indicates the general behavior of the classifier. For all of those criteria, a lower value indicates a better performance.

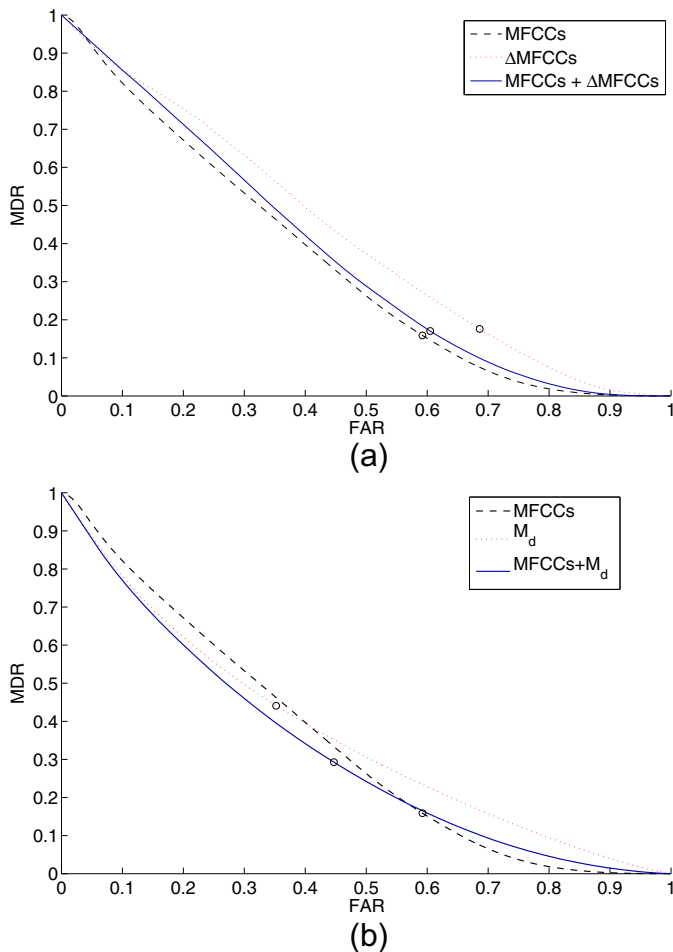


Fig. 4. DET curves of selected features computed over the IOWA database. On both axis, a lower value indicates better performance.

### 5.2. Isolated musical tones

We consider in a first experiment, two databases of isolated tones, the IOWA database and the RWC database proposed by Goto et al. (2003). The IOWA database features 3637 tones of mean duration 4.3 s and standard deviation 7.5 s from 15 musical instruments for a total of 4.3 h. The RWC database features 6138 tones of mean duration 2.6 s and standard deviation 0.8 s for a total of 4.7 h.

Considering isolated notes is helpful to study the performance of the evaluated features in a controlled environment. Indeed, in

this case, the observation interval is at an optimal location and duration and no frequency resolution issues may arise except for very low pitched tones, making the tracking of partials relatively non-ambiguous.

We first evaluate the features solely. The results are shown in Table 1. For each of the TAFs, the spectral features perform best. For the amplitude and magnitude ones, removing the polynomial before the spectral transformation is helpful. Even though there should be no tracking issues, the magnitude features clearly dominate the other TAFs, leading to comparable results with the MFCCs. It should be noticed that there is a large correlation between the different evaluation criteria (EER, AREA, and MCP) demonstrating a good behavior of the evaluated classifiers.

To fully understand the behavior of the TAFs when considered jointly with the MFCCs, we plot in Fig. 4 the DET curves for several features computed over the IOWA database. In Fig. 4a, the DET curves of the MFCCs,  $\Delta$ MFCCs are considered solely or jointly. The curve of the MFCCs shows a typical behavior with a linear slope on the high MDR range. The DET curve of the  $\Delta$ MFCCs shows similar evolution properties with worst performance. As they are highly correlated, considering them jointly leads to an averaging of their performance.

On the contrary, the DET curve of the  $M_d$  feature shows a different behavior, almost symmetrical towards the unity slope. Considering them jointly leads to a very balanced classifier with good properties. Empirical experiments show that a better joint feature can be obtained by weighting the features prior to similarity calculation.

Considering the proposed TAFs jointly with the MFCCs leads to an improvement with respect to the performance of the TAFs. However, only the use of the MEFs leads to an improvement with respect to the sole use of the MFCCs, see Table 2.

### 5.3. Solo performance music

We consider in a second experiment a database of solo recordings used in several musical instruments classification experiments done by Essid et al. (2004) and Joder et al. (2009). The SOLOS database features 505 solos recordings where each one is of mean duration 110 s and standard deviation 162 s performed by 20 different instruments for a total of 15.56 h.

As the tracking is more difficult when considering less constrained sounds, the performance of the FEFs and the AEFs do not improve compared to the previous experiment. Consequently, only the spectral MEF are shown in Table 2. Also, the use of the DC or polynomial removal does not lead to a significant difference. The loss of relevance of the polynomial removal may be due to the fact

Table 1

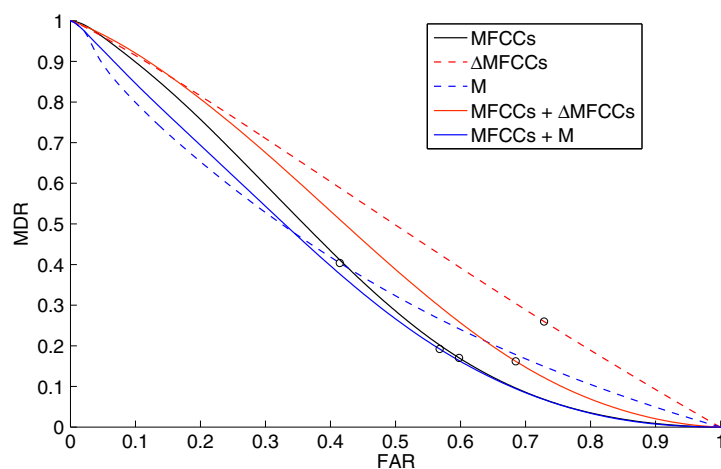
Results for single tones and single features. Best results for each feature group are displayed in bold characters.

Criterion	Database	$f$	$F$	$F_c$	$a$	$A$	$A_c$	$ap$	$Ap$	$Apc$	
EER	IOWA	0.699	<b>0.6</b>	0.621	0.659	0.655	0.656	0.69	0.639	<b>0.618</b>	
	RWC	0.72	<b>0.687</b>	0.69	0.696	0.68	0.677	0.726	0.677	<b>0.675</b>	
AREA	IOWA	0.491	<b>0.401</b>	0.418	0.458	0.451	0.45	0.484	0.434	<b>0.411</b>	
	RWC	0.512	<b>0.48</b>	0.482	0.491	0.474	0.47	0.518	0.472	<b>0.465</b>	
MCP	IOWA	0.928	<b>0.61</b>	0.629	0.66	0.656	0.656	0.701	0.645	<b>0.619</b>	
	RWC	0.928	<b>0.687</b>	0.691	0.723	0.684	0.677	0.977	0.679	<b>0.675</b>	
Criterion	Database	MFCCs	$\Delta$ MFCCs	$m$	$M$	$M_c$	$mp$	$Mp$	$Mpc$	$M_d$	$Mdc$
EER	IOWA	<b>0.56</b>	0.62	0.71	0.56	0.62	0.71	<b>0.54</b>	0.61	0.56	0.63
	RWC	<b>0.58</b>	0.65	0.73	0.69	0.69	0.72	<b>0.68</b>	0.69	0.69	0.7
AREA	IOWA	<b>0.34</b>	0.41	0.49	0.35	0.41	0.49	<b>0.34</b>	0.40	0.36	0.42
	RWC	<b>0.37</b>	0.44	0.52	0.48	0.49	0.51	<b>0.48</b>	<b>0.48</b>	0.48	0.49
MCP	IOWA	<b>0.61</b>	0.70	0.92	0.56	0.64	0.87	<b>0.54</b>	0.61	0.56	0.65
	RWC	<b>0.67</b>	0.69	0.96	0.80	0.79	0.96	0.75	<b>0.75</b>	0.80	0.77

**Table 2**

Results for single tones for joint features. Best results for each feature group are displayed in bold characters.

Criterion	Database	MFCCs+								
		$f$	$F$	$F_c$	$a$	$A$	$Ac$	$ap$	$Ap$	$Apc$
EER	IOWA	0.609	<b>0.567</b>	0.568	0.604	0.586	0.586	0.613	0.579	<b>0.56</b>
	RWC	0.649	<b>0.613</b>	0.619	0.634	0.607	<b>0.606</b>	0.637	0.608	0.609
AREA	IOWA	0.398	<b>0.357</b>	0.36	0.392	0.377	0.376	0.4	0.369	<b>0.349</b>
	RWC	0.441	<b>0.409</b>	0.413	0.426	0.404	<b>0.401</b>	0.43	0.406	<b>0.401</b>
MCP	IOWA	0.645	0.592	<b>0.578</b>	0.614	0.611	0.614	0.695	0.602	<b>0.577</b>
	RWC	0.726	<b>0.627</b>	0.644	0.651	0.632	<b>0.626</b>	0.75	0.63	0.63
		MFCCs+								
		$\Delta$ MFCCs	$m$	$M$	$Mc$	$mp$	$Mp$	$Mpc$	$Md$	$Mdc$
EER	IOWA	<b>0.578</b>	0.608	0.526	0.551	0.617	<b>0.522</b>	0.558	0.526	0.557
	RWC	<b>0.621</b>	0.651	<b>0.6</b>	0.613	0.644	0.604	0.622	0.601	0.613
AREA	IOWA	<b>0.364</b>	0.395	0.32	0.345	0.404	<b>0.316</b>	0.346	0.321	0.35
	RWC	<b>0.411</b>	0.443	<b>0.398</b>	0.41	0.435	0.399	0.414	<b>0.398</b>	0.41
MCP	IOWA	<b>0.629</b>	0.663	0.534	0.556	0.677	<b>0.53</b>	0.573	0.534	0.565
	RWC	<b>0.696</b>	0.73	0.625	0.639	0.73	0.636	0.668	<b>0.624</b>	0.637

**Fig. 5.** DET curves for selected features computed over the SOLOS database with a texture observation interval.**Table 3**

Results for solo recordings. Best results are displayed in bold characters.

Criterion	Integration type	MFCCs	$\Delta$ MFCCs	$M$	$Mc$	MFCCs+		
						$\Delta$ MFCCs	$M$	$Mc$
EER	Texture	0.523	0.661	0.552	0.584	0.579	<b>0.519</b>	0.526
	Event	<b>0.538</b>	0.665	0.615	0.623	0.59	0.56	0.57
AREA	Texture	0.314	0.456	0.354	0.378	0.357	<b>0.297</b>	0.304
	Event	<b>0.317</b>	0.45	0.408	0.425	0.363	0.334	0.341
MCP	Texture	0.509	0.658	0.552	0.579	0.557	<b>0.501</b>	0.509
	Event	<b>0.515</b>	0.651	0.61	0.622	0.563	0.536	0.544

that the boundaries of the acoustical units considered in this experiment are arbitrary set. As shown in Fig. 5, the properties of the different features discussed in the previous section are rather similar when considering real-world sounds. The use of the dimensionality reduction over the  $M$  feature leads to a slight decrease in the performance. The longer duration of the acoustical unit leads to a general performance decrease as shown in Table 3. This decrease is small concerning the MFCCs. On contrary, it is significant for the MEFs. This demonstrates that the length of the observation interval is crucial for capturing meaningful periodicities in the variations of the spectral parameters. If the interval is too short, no modulations

can be captured. On contrary, if the interval is too large, the observed modulations corresponds to the transition between several notes, which is typical from the score played by the instrumentalist and not from the actual timbre of the instrument.

## 6. Conclusion

We have proposed in this paper several approaches for extracting the evolution of the spectral parameters over time and for modeling them in a meaningful way. We show that the spectral

description of the evolution of the parameters of the partials is relevant as well as the removal of the polynomial part of the amplitude evolution when considering isolated notes databases. When facing less constrained scenarios, such as the solos database, the removal of the polynomial part does not improve the results whereas the spectral description is still in favor, leading to a relevant feature set when combined with the MFCCs.

As a conclusion, the proposed features are found to be more adapted to the tasks considered in this paper than the standard feature-level temporal dynamic features usually considered, the  $\Delta$ MFCCs. The results obtained by the FEFs and AEFs may be limited by the partial tracker considered in the experiments. The use of more advanced algorithm like those introduced by Lagrange et al. (2007) and by Robel (2006) may reduce the gap between model-based features (FEF and AEF) and transform-based features (MEF).

### Acknowledgments

The authors want to thank Dr. Ellis for providing the community with the code for MFCCs computation and partial tracking. This work has been partly funded by the Quaero project within the task 6.4: “Music Search by Similarity” and the French GIP ANR under contract ANR-06-JCJC-0027-01, “Décompositions en Éléments Sonores et Applications Musicales” – DESAM.

### Appendix A. Notation glossary

- $\Delta$ MFCCs:  $\Delta$  Mel-Frequency Cepstral Coefficients.
- AEF: Amplitude Evolution Feature.
- AREA: Area under DET curve.
- DET: Detection Error Trade-Off.
- EER: Equal Error Rate.
- FAF: Frequency-Axis Feature.
- FAR: False-Alarm Rate.
- FEF: Frequency Evolution Feature.
- MCP: Minimum Cost Point.
- MDR: Miss Detection Rate.
- MEF: Magnitude Evolution Feature.
- MFCCs: Mel-Frequency Cepstral Coefficients.
- TAF: Time-Axis Feature.

### Appendix B. Features glossary

The parameters amplitude, frequency and magnitude are implicitly of the evolution of the given parameter of a given acoustical unit through time.

- **f**: Time domain feature encoding the frequency of the partials, see Eq. (17).
- **a**: Time domain feature encoding the amplitude of the partials.
- **ap**: Time domain feature encoding the polynom-removed amplitude of the partials.
- **m**: Time domain feature encoding the magnitude of some bins of the Fourier spectrum.
- **F**: Frequency domain feature encoding the frequency of the partials of a given acoustical unit through time, see Eq. (18).
- **A**: Frequency domain feature encoding the amplitude of the partials.
- **Ac**: Frequency domain feature encoding the amplitude of the partials where only the lowest frequency bins are considered.
- **Ap**: Frequency domain feature encoding the polynom-removed amplitude of the partials.
- **Apc**: Combination of the two previous features.

- **M**: Frequency domain feature encoding the magnitude of some bins of the Fourier spectrum.
- **Mc**: Frequency domain feature encoding the magnitude of some bins of the Fourier spectrum. Only the lowest frequencies are considered.
- **Mp**: Frequency domain feature encoding the polynom-removed magnitude of some bins of the Fourier spectrum.
- **Mpc**: Combination of the two previous features.
- **Md**: Frequency domain feature encoding the magnitude of some bins of the Fourier spectrum. The first spectral bins are discarded.
- **Mdc**: Combination of the features *Md* and *Mc*.

### References

- Aucouturier, J.-J., Pachet, F., 2007. The influence of polyphony on the dynamical modelling of musical timbre. *Pattern Recognition Lett.* 28 (5), 654–661.
- Cuturi, M., Vert, J.-P., Birkenes, O., Matsui, T., 2007. A kernel for time series based on global alignments. In: *Proc. IEEE ICASSP*, vol. 2, pp. 413–416.
- Davis, S., Mermelstein, P., 1980. Comparison of parametric representations for monosyllabic word recognition in continuously spoken sentences. *IEEE Trans. Acoust. Speech Signal Process.* 28, 357–366.
- Ellis, D.P.W., 2003. Sinewave and Sinusoid + Noise Analysis/Synthesis in Matlab. <<http://www.ee.columbia.edu/dpwe/resources/matlab/sinemodel/>>, online web resource.
- Ellis, D.P.W., 2005. PLP and RASTA (and MFCC, and inversion) in Matlab. <<http://www.ee.columbia.edu/dpwe/resources/matlab/rastamat/>>, online web resource.
- Eronen, A., 2003. Musical instrument recognition using ICA-based transform of features and discriminatively trained HMMs. In: 7th Internat. Symposium on Signal Processing and its Applications.
- Essid, S., Richard, G., David, B., 2004. Musical instrument recognition based on class pairwise feature selection. In: *Proc. ISMIR*. <<http://ismir2004.ismir.net/proceedings/p102-page-560-paper194.pdf>>.
- Fant, G., 1960. *Acoustic Theory of Speech Production*. Mouton, The Hague.
- Goto, M., Hashiguchi, H., Nishimura, T., Oka, R., 2003. RWC music database: Music genre database and musical instrument sound database. In: *Proc. 4th Internat. Conf. on Music Information Retrieval (ISMIR 2003)*.
- Grey, J.M., Moorer, J.A., 1977. Perceptual evaluations of synthesized musical instrument tones. *J. Acoust. Soc. Am.* 62 (3), 454–462.
- Joder, C., Essid, S., Richard, G., 2009. Temporal integration for audio classification with application to musical instrument classification. *IEEE Trans. Acoust. Speech Signal Process.* 17 (1), 174–186.
- Kitahara, T., Goto, M., Komatani, K., Ogata, T., Okuno, H.G., 2006. Musical instrument recognizer “Instrogram” and its application to music retrieval based on instrumentation similarity. In: *Proc. IEEE Internat. Symposium on Multimedia*, ISBN 0-7695-2746-9. <<http://dx.doi.org/10.1109/ISM.2006.113>>.
- Lagrange, M., 2004. *Sinusoidal Modeling of Polyphonic Sounds*. Ph.D. Thesis, University of Bordeaux 1, LaBRI (in French).
- Lagrange, M., 2005. A new dissimilarity metric for the clustering of partials using the common variation cue. In: *Proc. ICMC, organization ICMA, Barcelona, Spain*.
- Lagrange, M., Marchand, S., Rault, J., 2007. Enhancing the tracking of partials for the sinusoidal modeling of polyphonic sounds. *IEEE Trans. Acoust. Speech Signal Process.* 28, 357–366.
- Marchand, S., Raspaud, M., 2004. Enhanced time-stretching using order-2 sinusoidal modeling. In: *Proc. DAFx, Federico II University of Naples, Italy*, pp. 76–82.
- Martin, K.D., Kim, Y.E., 1998. Musical instrument identification: A pattern-recognition approach. In: *Proc. 136th Meeting of the Acoustical Society of America*.
- Martin, A., Doddington, G., Kamm, T., Ordowski, M., Przybocki, M., 1997. The DET curve in assessment of detection task performance. In: *Proc. EuroSpeech*.
- McAdams, S., 1989. Segregation of concurrent sounds: Effects of frequency modulation coherence. *JAES* 86 (6), 2148–2159.
- McAdams, S., Bigand, E., 1993. *Thinking in Sound*. Oxford Science Publications (Chapter 2.5).
- McAulay, R.J., Quatieri, T.F., 1986. Speech analysis/synthesis based on a sinusoidal representation. *IEEE Trans. Acoust. Speech Signal Process.* 34 (4), 744–754.
- Mellody, M., Wakefield, G., 2000. The time–frequency characteristic of violin vibrato: Modal distribution analysis and synthesis. *J. Acoust. Soc. Am.* 107, 598–611.
- Meng, A., Ahrendt, P., Larsen, J., Hansen, L., 2007. Temporal feature integration for music genre classification. *IEEE Trans. Acoust. Speech Signal Process.* 15 (5), 1654–1664.
- Raspaud, M., 2007. *Hierarchical Spectral Models for Sound and Applications*. Ph.D. Thesis, University of Bordeaux 1, 2007.
- Raspaud, M., Marchand, S., Girin, L., 2005. A generalized polynomial and sinusoidal model for partial tracking and time stretching. In: *Proc. DAFx, Universidad Politecnica de Madrid*. ISBN: 84-7402-318-1, pp. 24–29.
- Robel, A., 2006. Adaptive additive modeling with continuous parameter trajectories. *IEEE Trans. Acoust. Speech Signal Process.* 14 (4), 1440–1453.



- Scaringella, N., Zoia, G., 2005. On the modelling of time information for automatic genre recognition systems in audio signals. In: Proc. 6th Internat. Conf. on Music Information Retrieval (ISMIR).
- Serra, X., Smith, J.O., 1990. Spectral modeling synthesis: A sound analysis/synthesis system based on a deterministic plus stochastic decomposition. *CMJ* 14 (4), 12–24.
- Shimodaira, H., Noma, K., Nakai, M., Sagayama, S., 2002. Dynamic time-alignment kernel in support vector machine. *Advances in Neural Information Processing Systems*. MIT Press, Cambridge, MA <[citeseer.ist.psu.edu/shimodaira01dynamic.html](http://citeseer.ist.psu.edu/shimodaira01dynamic.html)>.
- Tzanetakis, G., Cook, P., 2002. Musical genre classification of audio signals. *IEEE Trans. Acoust. Speech Signal Process.* 10 (5), 293–302.