

Metadata of the article that will be visualized in OnlineFirst

1	Article Title	The Cuidado music browser: an end-to-end electronic music distribution system
2	Journal Name	Multimedia Tools and Applications
3		Family Name Pachet
4		Particle
5		Given Name François
6	Corresponding	Suffix
7	Author	Organization SONY Computer Science Laboratory
8		Division
9		Address 6, rue Amyot, Paris 75005, France
10		e-mail pachet@csl.sony.fr
11		Family Name Aucouturier
12		Particle
13		Given Name Jean-Julien
14		Suffix
15	Author	Organization SONY Computer Science Laboratory
16		Division
17		Address 6, rue Amyot, Paris 75005, France
18		e-mail jj@csl.sony.fr
19		Family Name Burthe
20		Particle La
21		Given Name Amaury
22		Suffix
23	Author	Organization SONY Computer Science Laboratory
24		Division
25		Address 6, rue Amyot, Paris 75005, France
26		e-mail amaury@csl.sony.fr
27		Family Name Zils
28		Particle
29		Given Name Aymeric
30		Suffix
31	Author	Organization SONY Computer Science Laboratory
32		Division
33		Address 6, rue Amyot, Paris 75005, France
34		e-mail aymeric@csl.sony.fr
35		Family Name Beurive
36		Particle

37		Given Name	Anthony
38		Suffix	
39		Organization	SONY Computer Science Laboratory
40	Author	Division	
41		Address	6, rue Amyot, Paris 75005, France
42		e-mail	beurive@csl.sony.fr
<hr/>			
43		Received	
44	Schedule	Revised	
45		Accepted	
<hr/>			
46	Abstract	<p>The IST project Cuidado, which ran from January 2001 to December 2003, produced the first entirely automatic chain for extracting and exploiting musical metadata for browsing music. The Sony CSL laboratory is primarily interested in the context of popular music browsing in large-scale catalogues. First, we are interested in human-centred issues related to browsing “Popular Music.” Popular here means that the music accessed to is widely distributed, and known to many listeners. Second, we consider “popular browsing” of music, i.e., making music accessible to non-specialists (music lovers), and allowing sharing of musical tastes and information within communities, departing from the usual, single user view of digital libraries. This research project covers all areas of the music-to-listener chain, from music description—descriptor extraction from the music signal, or data mining techniques—similarity based access and novel music retrieval methods such as automatic sequence generation, and user interface issues. This paper describes the scientific and technical issues at stake, and the results obtained.</p>	
<hr/>			
47	Keywords	separated by ' - '	
<hr/>			
48	Foot note	information	
<hr/>			

The Cuidado music browser: an end-to-end electronic music distribution system

4
5

François Pachet · Jean-Julien Aucouturier ·
Amaury La Burthe · Aymeric Zils · Anthony Beurive

6
7

© Springer Science + Business Media, LLC 2006

9

Abstract The IST project Cuidado, which ran from January 2001 to December 2003, produced the first entirely automatic chain for extracting and exploiting musical metadata for browsing music. The Sony CSL laboratory is primarily interested in the context of popular music browsing in large-scale catalogues. First, we are interested in human-centred issues related to browsing “Popular Music.” Popular here means that the music accessed to is widely distributed, and known to many listeners. Second, we consider “popular browsing” of music, i.e., making music accessible to non-specialists (music lovers), and allowing sharing of musical tastes and information within communities, departing from the usual, single user view of digital libraries. This research project covers all areas of the music-to-listener chain, from music description—descriptor extraction from the music signal, or data mining techniques—similarity based access and novel music retrieval methods such as automatic sequence generation, and user interface issues. This paper describes the scientific and technical issues at stake, and the results obtained.

12
13
14
15
16
17
18
19
20
21
22
23
24

Q1 Keywords

25

1 Introduction

27

1.1 Existing popular music access systems

28

There are now many on line searchable music databases. We can classify them in the following categories.

29
30

F. Pachet (✉) · J.-J. Aucouturier · A. La Burthe · A. Zils · A. Beurive
SONY Computer Science Laboratory, 6, rue Amyot, 75005 Paris, France
e-mail: pachet@csl.sony.fr

J.-J. Aucouturier
e-mail: jj@csl.sony.fr

A. La Burthe
e-mail: amaury@csl.sony.fr

A. Zils
e-mail: aymeric@csl.sony.fr

A. Beurive
e-mail: beurive@csl.sony.fr

First, purely editorial systems propose systematic editorial information on popular music, including albums track listings (CDDB,¹ Musicbrainz²), information on artists and songs (AMG³ and Muze⁴). This information is created by music experts, or in a collaborative fashion (CDDB, Musicbrainz). These systems provide useful services for *Electronic Music Distribution* (EMD) systems, but cannot be considered as fully-fledged EMD systems per se, as they provide only superficial and incomplete information on music titles, supposed to exist somewhere else.

The MoodLogic⁵ browser proposes a complete solution for Popular Music access. The core idea of MoodLogic is to associate metadata to songs automatically thanks to two basic techniques: 1) an audio fingerprinting technology able to recognize music titles on personal hard disks, and 2) a database collecting user ratings on songs, which is incremented automatically, and in a collaborative fashion. An ingenious proactive strategy is enforced to encourage users to rate songs, in order to get tokens that allow them to get more metadata from the server. MoodLogic relies entirely on metadata obtained from user ratings and does not perform any acoustic analysis of songs. However, collaborative music rating does not exhaust the description potential of music, and our Browser proposes many other types of metadata.

Other proposals have been made either for fully-fledged music browsers, or for ingredients to be used in browsers (fingerprinting techniques, collaborative filtering systems, metadata repositories, e.g., Wold et al. [20]) that we cannot cover here for reasons of space. We will describe in this paper only the parts of our project that we think are original and may contribute to address the needs of our targeted users.

1.2 The Cuidado music browser

The Cuidado music browser aims at developing all the ingredients of the music-to-listener chain, for a fully-fledge content-based access system. More precisely, the project covers the areas of 1) editorial metadata, 2) acoustic metadata, 3) metadata exploitation and browsing tools, 4) management and share of metadata among users.

The next sections describe the most important results obtained for each of these aspects.

2 Editorial metadata

To manage collections of music titles an application must have access to many information to identify, categorize, index, classify and generally organize music titles.

We consider here two types of data as editorial metadata:

- Consensual information or facts about music titles and artists,
- Content description of titles, albums or artists.

The first category is common to already existing EMD systems and does not raise any particular problem, as this information is universal by nature. It includes for instance: artist

¹ <http://www.gracenote.com/>

² <http://www.musicbrainz.org/>

³ <http://www.allmusic.com/>

⁴ <http://www.muze.com/>

⁵ <http://www.moodlogic.com/>

and songs name, albums and tracks listing, group members, date of recording for a given title, short biography for artists with date of birth, years of activity, etc.

The second category is more problematic. Content description includes such widely needed information as artist style, artist instruments, song mood, song review, song or artist genre and more generally attributes aiming at describing the intrinsic nature of the musical item at stake (artist or song). These descriptions are useful to the extent that they can be used for musical queries in large catalogues. The music browser enables to issue queries for both categories.

Furthermore, the music browser has a tool (see figure 1) devoted to editorial information management. The global architecture of the system is detailed in Section 6. This tool allows editing and adding artists and/or songs properties.

2.1 Editorial metadata philosophy

Editorial metadata are associated distinctly with music titles and artists.

Artists (taken in the most general sense) are key *music identifiers* for many users: Yesterday is by “The Beatles,” and “The 5th symphony” is by Beethoven. Artists are used also for solving ambiguity: “With a Little Help from my Friends” by the Beatles, is definitely not the same tune as the version by Bruce Springsteen. The “Stabat Matter” by Pergolese is not the one by Boccherini, etc. We call these artists “primary artists” as they are most commonly used to identify music titles. These examples show that primary artists are common ways of identifying music titles but also that the role of primary artists changes with styles: in Classical music, primary artists are usually composers. In non-Classical music they are usually performers. In our Browser, we introduced the notion of primary

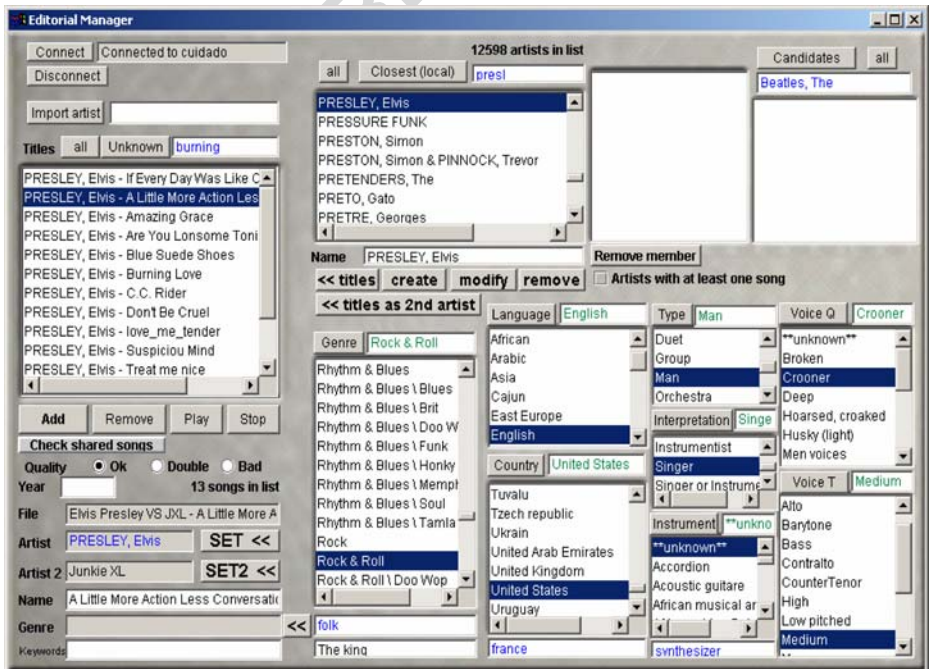


Fig. 1 The editorial data management panel

Print will be in black and white

artists in a deliberate ambiguous way, to cope for Classical and non-Classical music in a uniform way. 89

There are cases where primary artists are not enough for characterizing the identity of a piece. The “1st partita” of Bach has been recorded by Glenn Gould, and also by many other pianists, and this distinction is of course very important: not only for interpreters, but also for conductors (for orchestral pieces). In non-Classical music the need for secondary artists is also obvious, for instance to indicate that the Springsteen version of “A little help” is indeed a Beatles song. 90

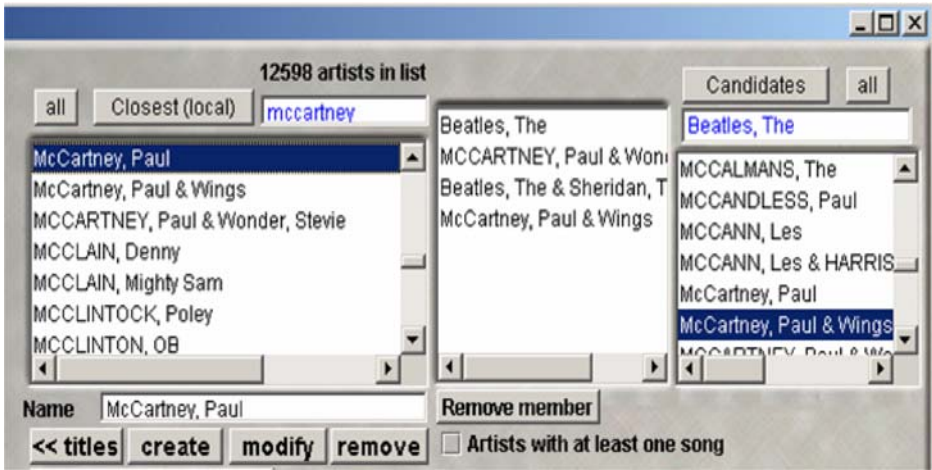
Existing repositories of editorial information do not provide systematic schemes for accessing artists and their relations to songs. This led us to constitute a database of artists, or more generally of “Musical Human Entities” (MHE), including both performers, composers, but also groups (the Beatles), orchestra (the Berlin Philharmonic), duets (Paul McCartney & Michael Jackson). To each artist (or MHE) is associated a limited but useful set of properties in fixed ontologies: type (composer, singer, instrumentist, etc.), country of origin, language (for singers), type of voice (for singers also), main instrument (for instrumentists). Other information concern the relation MHE entertain with each other. For instance, Paul McCartney is a *member of* The Beatles, and artist Phil Collins a *member of* the group Genesis. The Editorial MHE database may be seen more as a knowledge base than a database. 91

Concerning music titles, our tool enables basic editions as title name or keywords, as well as less obvious features such as title genre, primary and secondary artist introduced before. 92

Both artists and songs can be associated with a specific genre. Genres are badly needed for accessing music, and are as badly ill-defined. Our studies on existing taxonomies of genres have shown that there is no consensus, and that a consensus is probably impossible [3]. However, we propose here several ways to partially solve this problem. After several years of trials [13] and errors, we ended up with a simple two-level genre taxonomy consisting of 250 genres. The main property of this taxonomy is flexibility: users can classify artists or songs either in a generic way (Classical, Jazz), more precisely (Jazz/BeBop, Classical/Baroque). However, simpler taxonomies may also produce frustration, as some categories may contain artists or songs that users would consider very different. To make our taxonomy more flexible, we have introduced an optional “keyword” field, which may contain free words. These words may be entered by users to further refine their own classification perspective on artists or songs. This simple yet flexible approach has the advantage of uniformity: artists and songs are classified in the same taxonomy, allowing for various degrees of precision. For instance, The Beatles is classified in “Pop/Brit,” but Beatles songs may be classified in other genres (e.g., “Revolution 9” is “Rock/Experimental”) (figure 2). 93

3 Acoustic metadata 127

The main type of metadata that the MB proposes for songs besides editorial information is acoustic metadata, i.e., information extracted from the audio signal. The Mpeg7 standard aims at providing a format for representing these information, and a specialized audio group produces specific constructs to represent musical metadata [1, 10]. However, music metadata in Mpeg7 refers in general to low-level, objective information that can be extracted automatically in a systematic way. Typical descriptors (called LLD for Low-Level Descriptors in the Mpeg7 jargon) proposed by Mpeg7 concern superficial signal 128



Print will be in black and white

Fig. 2 The “member_of” predicate

characteristics such as means and variance of amplitude, spectral frequencies, spectral centroid, ZCR (zero crossing rate), etc. 135 136

Concerning high-level descriptors that can be mapped to high-level perceptual categories, Mpeg7 is strictly concerned with the format for representing this information, and not the extraction process per se. 137 138 139

3.1 Extracting high-level music percepts 140

We have conducted in the project several studies focusing on particular dimensions of music that are relevant in our context. 141 142

3.1.1 Rhythm 143

We have proposed a rhythm extractor [23], that is able to extract the time series of percussive sounds in music signals of popular music. Rhythm information is a useful extension of tempo or beat, as proposed by Scheirer in [17]. However, many things remain to be done in the field of rhythm. One key issue seems to rely not so much in how to extract rhythm, but how to exploit the information: most people are unable to describe rhythm with words, and even less to produce rhythm (our attempts at designing a query by rhythm did not prove successful). 144 145 146 147 148 149 150

3.1.2 Energy 151

In [22], we have addressed another dimension of music pertaining to popular music access, the perceptual energy, i.e., whether a song is thrilling and exciting (e.g., hard rock, dance music), or relaxing and calm (e.g., a piano piece by Schumann). 152 153 154

We have studied the correlation of experimental measures (user tests) with a variety of signal features, such as tempo, raw signal energy, spectral analysis, the associated variances, correlations... as well as their linear combinations (using discrimination analysis) and their possible compositions with signal operators (filters, etc...). The most discrimina- 155 156 157 158

tive parameter we found is $\log_{10}(\text{var}(\text{diff}(x^2)))$, which gave a classification error of 22% on the validation set.

3.1.3 *Timbre*

In [5], we have proposed to describe music titles based on their global *timbral quality*. Our motivation is that, although it is difficult to define precisely music taste, it is quite obvious that music taste is often correlated with timbre. Some sounds are pleasing to listeners, other are not. Some timbres are specific to music periods (e.g., the sound of Chick Corea playing on an electric piano), others to musical configurations (e.g., the sound of a symphonic orchestra). In any case, listeners are sensitive to timbre, at least in a global manner.

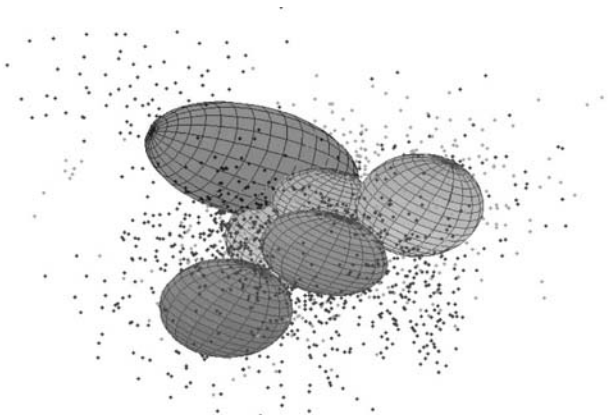
We model the global “sound” of a music title as a distribution in the space of mel cepstrum coefficients (MFCC). MFCCs provide a compact representation of the signal’s spectral envelopes, which are a good correlate of the timbre. By comparing timbre distributions between titles, it is then possible to match music titles of possible very different genres based solely on their timbre color. Figure 3 shows a 3D projection of the feature space (which is originally of dimension 8), showing two distributions of MFCCs, each modelled with a mixture of three Gaussian distributions (GMM). The light-grey GMM is the timbre model of the song “The Beatles—Yesterday,” and the dark-grey GMM is the timbre model of the song “Joao Gilberto—Besame Mucho.” This two songs have a very similar “sound” (acoustic guitar and a string quartet, plus a gentle and melancholic male voice), and indeed we see that their MFCC distributions are very close. As explained in Section 4, timbre models are used in the MusicBrowser to compute similarities between songs.

3.1.4 *Instrumental/voice presence*

A fourth descriptor which is currently available in the Music Browser describes whether a given tune contains singing voice or only instrumental sounds. This property is useful e.g., to either access particular “genres” of music (“opera” falls in the first category, while “piano sonatas” falls in the instrumental category), or to differentiate different versions of the same song (e.g., “Dub” instrumental versions of “reggae” songs).

There has been a large number of studies about the issue of speech/music discrimination (see e.g., [18]), which has received successful solutions, but the detection of singing voice

Fig. 3 Comparison of the timbre models of two songs: “The Beatles—Yesterday” and “Joao Gilberto—Besame Mucho”



has proved a more difficult problem. Berenzweig in [7] proposes to use complex features (output probabilities of a speech recognizer system) combined with hidden Markov models (HMMs). The extractor currently used in the Music Browser was designed automatically by the EDS system, described in the next section. It has a classification error of 19% on the validation set.

3.2 EDS: a general framework for extracting extractors 194

These various studies in descriptor extraction from acoustic signals have shown that the design of an efficient acoustic extractor is a very heuristic process, which requires sophisticated knowledge of signal processing, intuitions, and experience. Indeed, most approaches in feature extraction as published in the literature consist in using statistical analysis tools to explore spaces of combinations of LLD. The approaches proposed by Peeters [16], Scheirer [18] and Tzanetakis [19] typically fall in this category. However, these approaches are not capable of yielding precise extractors, and depend on the nature of the palette of LLD, which usually do not capture the relevant, often intricate and hidden characteristics of audio signals. Consequently, designing extractors is very expensive and hazardous.

On the other hand, user studies have shown that there is a virtually infinite number of extractors of musical attributes that could be useful in EMD systems. Different users have different needs: one—say, a jazz musician might be interested in listening to songs which exhibit a particular chord sequence, another may be interested by the sound (“some saturated guitar with a little bit of chorus”), while another simply wants to find “funky” music for his birthday party. Even when talking about the same attribute, the definitions (i. e., in terms of pattern recognition, the training sets) vary a lot. The perception of “harmonic complexity” of a tune for instance highly depends on the musical expertise of the listener.

These experiments have given rise to a systematic approach to feature extraction, embodied in the EDS system [14]. Departing from the usual LLD approach, the idea of EDS is to automate—in part or totally—the process of designing extractors. EDS searches in a richer and more complex space of signal processing functions, much in the same way than experts do: by inventing functions, computing them on test databases, and modifying them until good results are obtained.

To reach this goal EDS uses a genetic programming engine, augmented with fine grained typing system, which allows to characterize precisely the inputs and outputs of functions. EDS also uses rewriting rules to simplify complex signal processing functions (see the example of the Perceval equality being used by EDS to simplify the expression in figure 4). Finally EDS uses expert knowledge to guide its search, in the form of heuristics.

Typical heuristics include “do not try functions which contain too many repetition of the same operator,” or “apply twice a FFT on a signal is interesting, but not three times,” or also “spectral coefficients are particularly useful when applied on signals in the temporal domain, possible filtered,” etc.

The signal operators available in the EDS which serve as basic bricks for building extractors include the full set of MPEG7 LLDs, but also typical signal operators like filters, FFT, time windowing, and higher level operators like pitch detection, partial tracking or mel filterbank. These operators are selected from the literature and our experiments of designing extractors manually. The features designed and discovered by the system can be further combined, manually or automatically, by statistical models like GMMs or HMMs, or classifiers like neural networks. The output of the whole process is an executable file, which can be directly integrated in applications like the Music Browser.

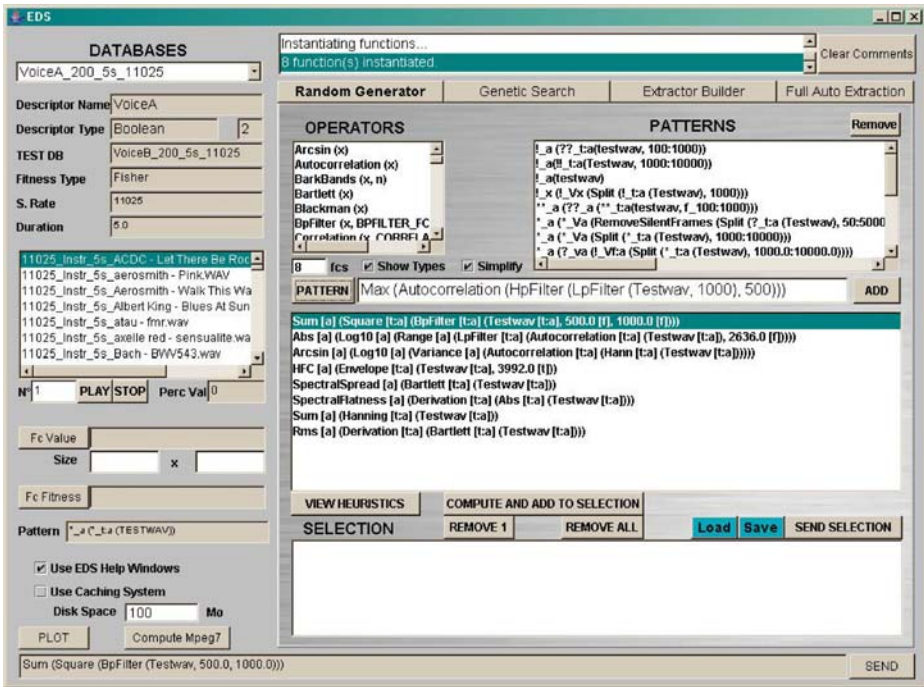


Fig. 4 Screenshot of the EDS system

The current extractors targeted by EDS are perceptual energy (or a refinement of the descriptor we designed by hand), discrimination between songs and instrumental (already described in the previous section), discrimination between studio and live versions of songs, harmonicity vs noisiness, percussivity, harmonic complexity, etc. The ambitious goal of EDS makes it a project in itself, as it aims at capturing complex knowledge, in an expanding field. However, we think that the contribution to the MIR community is potentially important as it is a first step towards a unified vision of high level audio feature extraction.

4 Similarity

The notion of similarity is of utter importance in the field of music information retrieval, and the expectation to have systems that find songs that are “similar” to one or several seed songs is now second nature. However, here again, similarity is ill-defined, and it can be of many different sorts. For instance, one may consider all the titles by a given artist as similar. And they are, of course, artist-wise. Similarity can also occur at the feature level. For instance, one may consider that Jazz saxophone titles are all similar. Music similarity can yet occur at a larger level, and concern songs in their entirety. For instance, one may consider Beatles titles as similar to titles from, say, the Beach Boys, because they were recorded in the same period, or are considered as the same “style.” Or two titles may be considered similar by a user or a community of users for no objective reason, simply because they think so.

Print will be in black and white

236
 237
 238
 239
 240
 241
 242
 243

 244
 245
 246
 247
 248
 249
 250
 251
 252
 253
 254
 255

4.1 Acoustic similarity

256

Feature-based similarity is trivially obtained by defining similarity measures from the metadata obtained and described above, either editorial or acoustic. Most descriptors yield implicit similarity measures that can be useful in some circumstances, e.g., similarity of tempo, of energy, or similarity based on artist relationships, etc.

257
258
259
260

One very interesting type of similarity that we already mentioned is based on the global “timbre” of the songs. The distance analysis is based on Gaussian models of Cepstrum coefficients as described in [5]: a first model is sampled and then the likelihood of the samples is computed given the other model. Figure 5 shows a screenshot of the “Find by Similarity” panel in the Music Browser. Here, the user has select a jazz piano song (“Ahmad Jamal-L’instant de Vérité”), and asked the system to return “songs that sound the same.” The result lists contains songs of many genres, which all contain romantic-styled piano: Jazz (Hank Jones, Alain Jean-Marie), Classical piano pieces (Brahms, Chopin), and even a “Variety” song (William Sheller, a French singer and pianist who had a classical training).

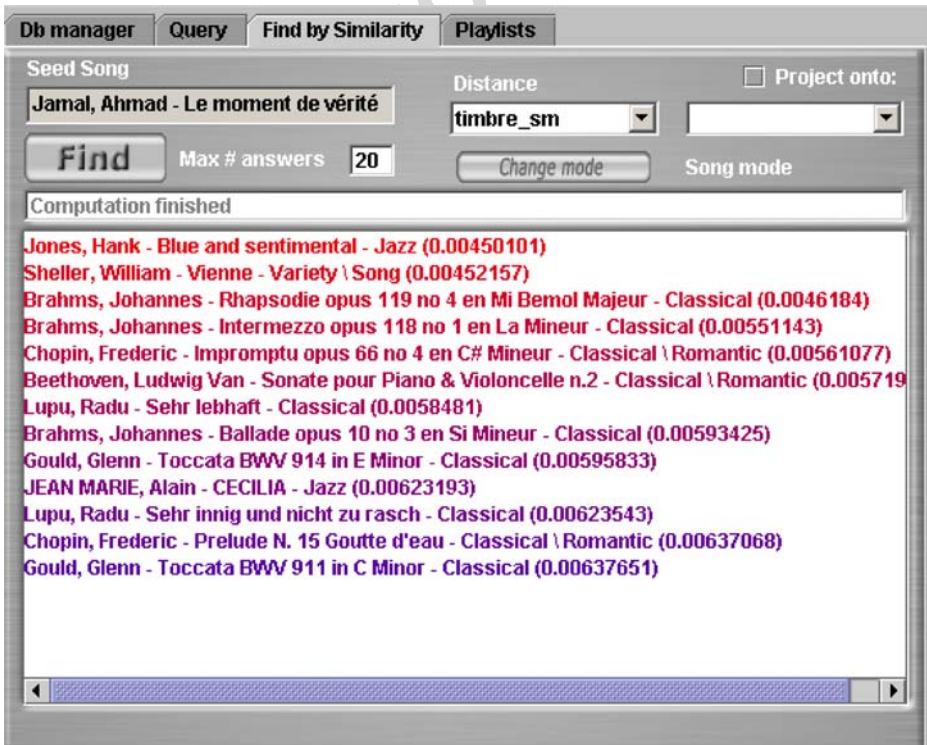
261
262
263
264
265
266
267
268
269

4.2 Cultural similarity

270

Cultural similarity is based on a well-known technique used in statistical linguistics: co-occurrence analysis. Co-occurrence analysis is based on a simple idea: if two items appear in the same context, it is obvious that there is some kind of similarity between them. In linguistics, co-occurrence analysis based on large corpora of written and spoken text has

271
272
273
274



Print will be in black and white

Fig. 5 The “Find by Similarity” panel in the Music Browser

been used to extract clusters of semantically related words. Similarity measurements based on co-occurrence counts have been demonstrated to be cognitively plausible [8]. We have identified several interesting corpora:

- The web, 275
- Radio programs, 276
- Compilations. 277

In the framework of Cuidado we are currently exploiting the web with a crawler specifically designed for this task. 281
282

4.2.1 The Cuidado Crawler 283

It is a multi-thread software designed to crawl the web. Its goal is to gather as many web pages as possible, parsing every word and every link on each page. Each crawled web page is given a score according to the presence of keywords. Each URL gathered on the page is given the score of the page. Several crawling modes are available from blind crawling (no keywords, only a few starting URLs) to narrow crawling (specific keywords that can be changed dynamically) 284
285
286
287
288
289

The Cuidado Crawler can create/handle several crawling database. Each user can create as many databases as his hard drive can contain. Therefore, users can create database on specific topics or according to specific tastes. For example, if you interested in “intelligent techno.” There are over 118,000 hits in Google⁶ for this query and probably more when you will read this. You can start crawling using the first answers provided by Google as well as specific keywords you entered like “new, research, noise, click and avant-garde.” Therefore you construct an “intelligent-techno” oriented database which favours your vision of intelligent techno thanks to the keywords. 290
291
292
293
294
295
296
297

The second part of this software is devoted to the distance computation. The various formula used here were introduced in [15]. We are looking for occurrences of words in the same page, taking into account the number of pages where each word is found. 298
299
300

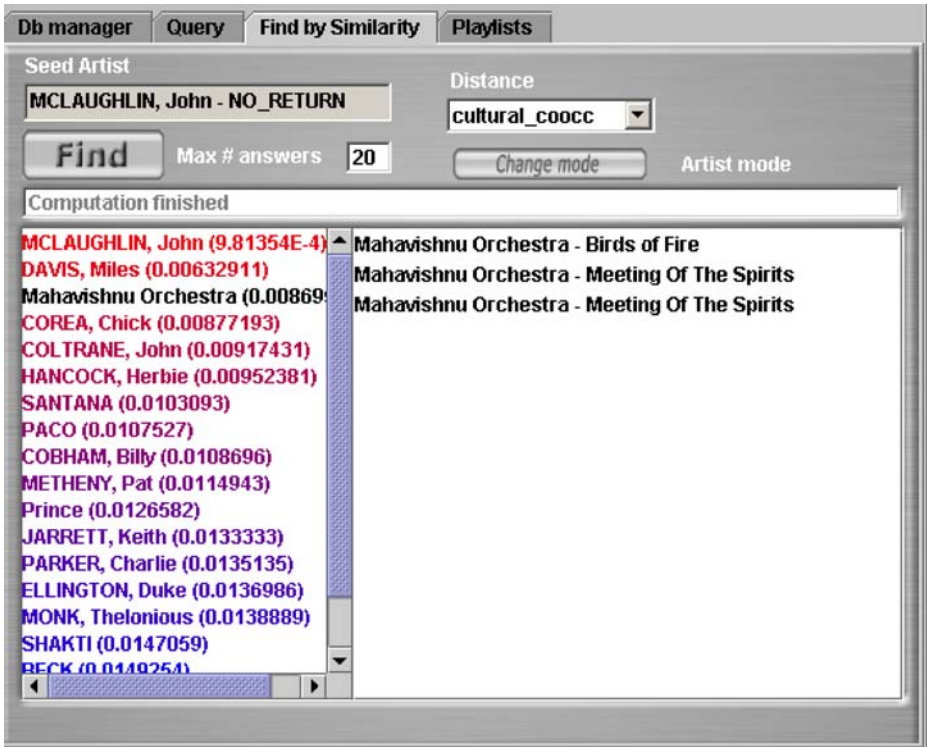
4.2.2 Integration in the Music Browser 301

To ensure the compatibility with the Music browser users can import any data coming from Cuidado tables. The distance is then computed for each entry and is exported back to the Music Browser as a new distance table. Figure 6 shows the results of a cultural similarity query on the jazz guitarist John McLaughlin. The closest artists include Miles Davis (McLaughlin played on two of his records, “In a Silent Way” and “Bitches Brew” in 1969), the Mahavishnu Orchestra (a fusion band formed by McLaughlin in 1971, including drummer Billy Cobham, also present on the list), jazz pianist Chick Corea (who played with McLaughlin and Miles Davis in the 1969 records), jazz guitarists Pat Metheny and Paco de Lucia (who McLaughlin played in trio with), etc. 302
303
304
305
306
307
308
309
310

5 Exploitation 311

We have covered so far the core technologies for producing content descriptions of music titles. A key issue is the exploitation of these information on the user side. The graphical 312
313

⁶ <http://www.google.com>



Print will be in black and white

Fig. 6 The similarity panel showing artists culturally similar to jazz guitarist John McLaughlin

interface issue is problematic because of the great variety of behaviours of users, and because the actual devices that will be used for large scale access to music catalogues are still unknown (computers? set-top boxes? PDAs? telephones? hard-disk Hi-Fi?). Many user interfaces have been proposed for music access systems, from straightforward feature-based search systems to innovative graphical representations of play lists. For instance, Gigabeat⁷ display music titles in spirals to reflect similarity relations titles entertain with each other. The gravitational model of SmarTuner,⁸ represent titles as mercury balls moving graciously on the screen, to or from “attractors” representing the descriptors selected by the user. The IBM GlassEngine⁹ proposes to browse a collection of pieces by minimalist composer Philip Glass, using a set of sliding cursors which rearrange the collection according to several criteria simultaneously (joy, sorrow, density, velocity, etc.). However gracious, these interfaces impose a fixed interaction model, and assume a constant attitude of users regarding exploration: either non-explorative—music databases in which you get exactly what you query—or very exploratory. But the users may not choose between the two, even less adjust this dimension to their wish. The current interface of the Music Browser aims at allowing users to choose between many modes of music access: explorative, precise, focused or hazardous.

314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330

⁷ <http://www.gigabeat.com>
⁸ <http://www.mzz.com>
⁹ <http://www.philipglass.com/glassengine>

5.1 Focused interfaces

331

The query panel (figure 7) is mostly dedicated to focused search in the database. In this panel users can issue queries on all available artists and songs metadata.

332
333

These metadata can be editorial: artists' names, songs' names, voice quality, etc. as well as computed: subjective energy, tempo, etc. The result of a query is a music titles list. Then this result set can be further filtered to return only songs with fast tempo, or only songs with a male singer. This result list can be transferred to the player for listening/exporting purpose

334
335
336
337

5.2 Explorative interfaces

338

5.2.1 Sliding between similarities

339

An interesting issue resulting from the studies on feature-based and cultural similarities is the comparison between these different sorts of similarity. For instance in figure 5, a starting title such as "Le moment de vérité" played by Ahmad Jamal, is considered by the MB as similar timbre-wise to "Humoresque Op. 20" by Schumann or "Blue and sentimental" by Hank Jones, but culturally, it is closer to "Ahmad's blues" by Miles Davis, because of the strong relationship between these two players, captured by the web crawler. Of course, there is no grounded truth here, and all these similarities are relevant. The next issue to solve is to aggregate these similarities, or at least propose users simple and meaningful ways of exploiting these different techniques.

340
341
342
343
344
345
346
347
348



Print will be in black and white

Fig. 7 Screenshot of the query panel in the Music Browser

In [5], we have proposed an interface, the “aha slider,” which allows the user to rank the results of a query according to two possibly orthogonal types of similarity. The slider is simply a way to filter the result set of one similarity according to the values of the second similarity measure. For instance, one can ask for “timbrally” similar songs which are also very close according to cultural similarity (e.g., “Ahmad’s blues” by Miles Davis), or, on the contrary, filter the result set so that it only contains songs which are culturally very distant from the query (e.g., Schumann or William Sheller).

This interface attempts to give the user full control over the degree of surprise and freedom in the way the system satisfies his request. A non-exploratory behavior (e.g., culturally similar) implies that the system should return exactly the answer to the query, or an answer that is as expected as possible (same title, same artist). An exploratory behavior (e.g., culturally distant) consists in letting the system try different regions of the catalogue rather than strictly match the query.

5.2.2 Playlist generation

An original feature introduced by the Browser is a powerful playlist generation system, based on constraint satisfaction techniques [2]. This technique allows user to get entire music playlists from a catalogue, by specifying only abstract properties on the playlist, such as:

- the playlist should contain 12 different titles,
- the playlist should not last more than 76 min,
- the genre of a title should be *close* to the genre of the next title,
- the playlist should contain at least 60% of *instrumental* titles,
- the sequence should contain titles with increasing tempo, etc.

The problem of generating such playlists given a very large title catalogue with musical metadata, and a set of arbitrary constraints is a NP-hard combinatorial problem. Moreover, in the case of a contradictory set of constraints, there may not be an exact solution. An ideal system should therefore be able to generate good approximate compromises. The Cuidado Music Browser is able to generate such playlists automatically (figure 8), using a fast algorithm based on adaptive search [2].

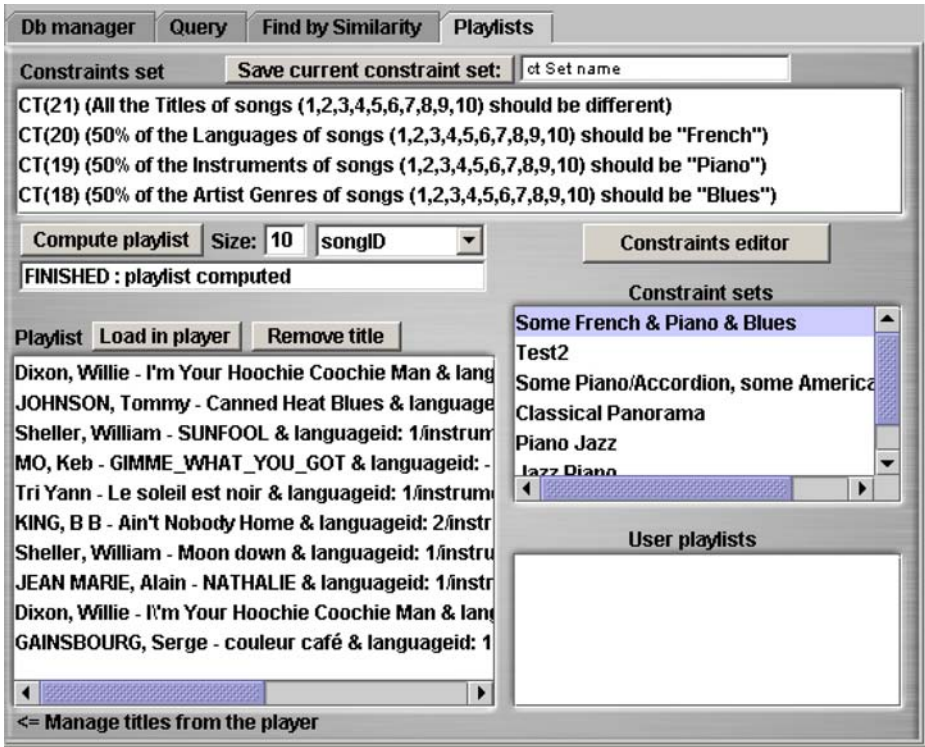
We give here an example of a five-title playlist with the following constraints:

1. Timbre continuity: the playlist should be “timbrally” homogeneous, and shouldn’t contain abrupt changes of textures.
2. Genre Cardinality: the playlist should contain 30% of Rock pieces, 30% of Folk, and 30% of Pop
3. Genre Distribution: the titles of the same genre should be as separated as possible.

One solution found by the system is the following playlist:

- Rolling Stones—You Can’t always get what you want—Genre=Pop/Blues
- Nick Drake—One of these things first—Genre=Folk/Pop
- Radiohead—Motion Picture Soundtrack—Genre=Rock/Brit
- The Beatles—Mother Nature’s Son—Genre=Pop/Brit
- Tracy Chapman—Talkin’ about a Revolution—Genre=Folk/Pop

Our current research regarding playlist generation aims at designing simple user interfaces to specify arbitrary constraints in a more intuitive way than in the current implementation,



Print will be in black and white

Fig. 8 Screenshot of the playlist generation system

which based on a crude mix of lists and multiple choices. A possible direction towards this is 392
 the use of simple drawings or gestures as a way to describe dynamical behaviours 393
 (“increasing”), or distribution properties (“a lot of,” “from here ... to here”). 394

6 Architecture 395

This section describes the general architecture of the Music Browser (figure 9). The central 396
 element of the architecture is the metadata server. This server is a MySQL database hosted 397
 on a SQL server. The server acts both as a server for PHP scripts and servlets. The 398
 MusicBrowser is implemented in Java and communicates with the MySQL database using 399
 JDBC drivers. The metadata server runs a PHP server accessible over the Internet. Specific 400
 PHP scripts allow client applications to fetch and submit metadata to this server. 401

The music browser contains four panels aimed at music title access: the player, the query 402
 panel, the similarity panel and the playlist panel. 403

Additionally, the browser includes two management tools: the editorial data management 404
 tool and the extractor and computation management tool. The purpose of the computation 405
 management tool is to compute descriptors for the songs in the database as well as similarity 406
 measures. It can use any stand-alone extractor (exe or bat files) developed by third party. 407

The editorial metadata management tool is used to manage artists and songs properties. 408
 It provides choice lists for each property and enables basic editions such as title name or 409

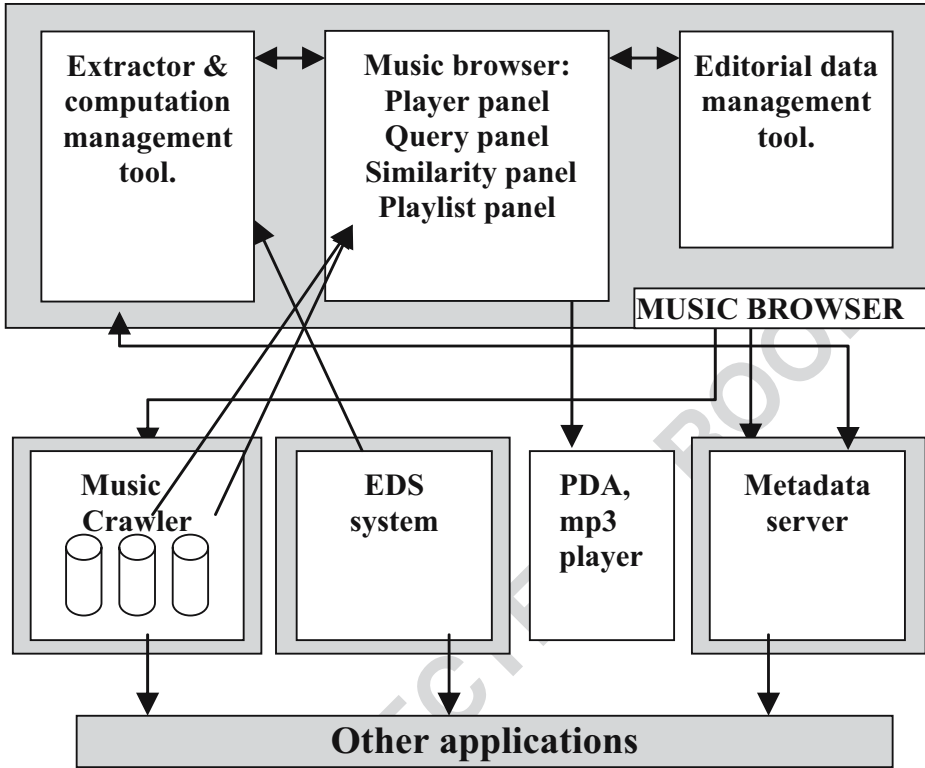


Fig. 9 Interaction between the different components of the Music Browser

keywords, as well as title genre, primary and secondary artist, as described in Section 2.1. 410
 This tool interacts on-line with our metadata server. 411

Lastly, with the apparition of ad hoc networks, users can share their data easily with 412
 other users and in a transparent way. This situation raises an issue in the management and 413
 synchronization of metadata. We describe in [11] a solution to allow both private and 414
 shared metadata to coexist in a single environment. 415

7 Future work: towards an API

416

Our experience in designing a large-scale EMD system such as the Music Browser shows 417
 that the main difficulty is to combine several systems/languages in a seamless manner: a 418
 database (SQL), an object-oriented engine to manage “multimedia items,” like songs, 419
 artists, albums, etc. (JAVA), user interfaces/interaction modules (JAVA), signal processing 420
 algorithms and extractors (Matlab, C), music rendering (JMF). All of these aspects 421
 interoperate closely, e.g., the interface calls an executable which computes a value, which is 422
 stored in the db, and re-used in another interaction module. 423

This architecture, although it does not present any particular technical difficulty, is 424
 expensive to design, and requires much incremental “doodling” both to specify and to 425
 build. On the other hand, such an architecture is needed for many other applications than 426
 the Music Browser, virtually every application concerned with content-based interaction, 427

access, browsing of large multimedia collections. Among other Sony CSL projects, the Musing [21], a composition tool to create sequences of samples according to high-level properties on their metadata (e.g., a steady tempo, with some voice samples, a given energy profile, etc.), and Personal Radio [12], an automatic, customized radio station, are based on the same type of architecture.

Moreover, the overhead of building such an architecture is often a limiting factor for many subtasks like evaluating content-extraction algorithms, a problem which is hotly debated in the music information retrieval community [9]. As described in [4], in order to evaluate and fine-tune algorithms like the timbre similarity used in the Music Browser, one needs to be able to:

- access and manage the collection of music signals the measures should be tested on
- store each result for each song (or rather each duplet of songs as we are dealing with a binary operation $\text{dist}(a,b)=d$) and each set of parameters
- compare results to a ground truth, which should also be stored
- build or import this ground truth on the collection of songs according to some criteria
- easily specify the computation of different measures, and to specify different parameters for each algorithm variant, etc.

Following these experiments, we have started developing a more general API, the so-called MCM (multimedia content management). MCM is a set of java classes, which offer the following data structures and functionalities:

- multimedia *items* (e.g., songs or artists), existing synchronously both in db and in memory.
- *fields* or metadata for each of these items (e.g., song's tempo or artist's name).
- *field values* for each item are read/written in db, and can be cached in memory for applications which require more CPU power, like playlist generation.
- items may link one to another (e.g., song items can be associated with artist items, video clip items, etc.). These associations are treated like fields (the "artist" item is a metadata of the "song" item), which values are the corresponding items.
- some fields are computable, i.e., their value is the output of an extractor, either computed online or offline, in batch mode.
- items can link to other items with *relations*, e.g., timbre or cultural similarity.
- items, fields, relations can be added (e.g., add a new directory of mp3s in the Browser, add a third-party extractor, etc.), updated, retrieved or deleted from the db.

Using MCM, all the architectural difficulties of creating databases, synchronizing data, calling extractors are hidden out. Applications like the Music Browser can be developed very quickly, by concentrating only on meaningful, higher-level concepts. Like for the EDS, we think that this is a potentially important contribution to the Music Information Retrieval community as it is a first step towards a unified vision of content based interaction and access systems.

8 Conclusion

The Cuidado music browser is the first large scale, fully content-based music access system. It includes all the technologies needed to extract descriptors, create similarity relations, and make these information easily available to users. The system is fully operational, and user tests have started to assess the usability of content information for

music access. Two side projects emerged from the design of this system: the EDS, a general framework for the automatic design of extractors, and MCM, an API to speed up the design of applications concerned with extracting and exploiting musical metadata for browsing music. Both projects constitute a first step towards a unified vision of content-based interaction and access systems.

References

1. Allamanche E, Herre J, Helmuth O, Frba B, Kasten T, Cremer M (2001) Content-based identification of audio material using MPEG-7 low level description. In: Proc. of the 2nd International Symposium on Music Information Retrieval, (ISMIR 01), Bloomington, Indiana, USA 478
2. Aucouturier J-J, Pachet F (2002) Scaling up playlist generation. In: Proc. of the IEEE International Conference on Multimedia and Expo (ICME 02), Lauzanne, Switzerland 479
3. Aucouturier J-J, Pachet F (2003) Musical genre: a survey. *J New Music Res* 32:1 480
4. Aucouturier J-J, Pachet F (2004) Improving timbre similarity: how high's the sky? *Journal of Negative Results in Speech and Audio Sciences (JNRSAS)*, submitted 481
5. Aucouturier J-J, Pachet F, Sandler M (2004) The way it sounds: timbre models for structural analysis and retrieval of music signals. *IEEE Trans Multimedia*, submitted 482
6. Belkin N (2000, August) Helping people find what they don't know. *Commun ACM* 43(8):58–61 483
7. Berenzweig A, Ellis D (2001) Locating singing voice segments within music signals. In: Proc. IEEE Workshop on Applications of Signal Processing to Acoustics and Audio (WASPAA 01), Mohonk, NY, USA 484Q3/Q4
8. Cohen W, Fan W (2000) Web-collaborative filtering: recommending music by crawling the web. In Proc. 9th International World Wide Web Conference (WWW9), Amsterdam, The Netherlands 485
9. Downie S (2003) Toward the scientific evaluation of music information retrieval systems. In: Proc. International Symposium on Music Information Retrieval (ISMIR 03), Baltimore, Maryland, USA 486
10. Herrera P, Serra X, Peeters G (1999) Audio descriptors and descriptors schemes in the context of MPEG-7. In: Proceedings of the International Computer Music Conference (ICMC 99), Beijing, China 487Q4
11. La Burthe A, Pachet F, Aucouturier JJ (2003) Editorial metadata in the Cuidado Music Browser: between universalism and autism. In: Proc. 3rd International Conference of Web Delivering of music (WedelMusic 03), Leeds, UK 488Q5
12. Pachet F (2003) Content management for electronic music distribution: the real issues. *Commun ACM* 2003 489
13. Pachet F, Cazaly D (2000). A taxonomy of musical genres. In: Proc. Content-Based Multimedia Information Access (RIAO), Paris, France 490
14. Pachet F, Zils A (2003) Evolving automatically high-level music descriptors from acoustic signals. Springer, Berlin Heidelberg New York LNCS, 2771 491
15. Pachet F, Westerman G, Laigre D (2001) Musical data mining for electronic music distribution. In: Proceedings of First International Conference of Web Delivering of Music (WedelMusic 01), Firenze, Italy 492
16. Peeters G, Rodet X (2002) Automatically selecting signal descriptors for sound classification. In: Proc. of the International Computer Music Conference (ICMC 02), Goteborg (Sweden) 493
17. Scheirer ED (1998, January) Tempo and beat analysis of acoustic musical signals. *J Acoust Soc Am (JASA)* 103(1):588–601 494
18. Scheirer E, Slaney M. Construction and evaluation of a robust multifeature speech/music discriminator. In: Proc. IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 97), Munich, Germany 495
19. Tzanetakis G, Perry C (2002, July) Musical genre classification of audio signals. *IEEE Trans Speech Audio Process* 10(5) 496
20. Wold E, Blum T, Keislar D, Wheaton J (1996) Content-based classification, search, and retrieval of audio. *IEEE Multimed* 3(3):27–36 497
21. Zils A, Pachet F (2001) Musical mosaicing. In: Proc. of COST-G6 Conference on Digital Audio Effects (DAFX01), Limerick, Ireland 498
22. Zils A, Pachet F. Extracting automatically the perceived intensity of music titles. In: Proc. of the COST-G6 Conference on Digital Audio Effects (DAFX03), London, UK 499
23. Zils A, Pachet F, Delerue O, Gouyon F (2002) Automatic extraction of drum tracks from polyphonic music signals. In: Proc. 2nd International Conference of web Delivering of Music (WedelMusic 02), Darmstadt, Germany 500



François Pachet, Civil Engineer (Ecole des Ponts et Chaussées), Ph.D. from University of Paris 6, Assistant Professor in Artificial Intelligence and Computer Science at Paris 6 University until 1997. In 1997, he set up the music research team at Sony Computer Science Laboratory (Paris), and developed the vision that metadata can greatly enhance the musical experience in all its dimensions from listening to performance. His team conducts research in interactive music listening, computer-aided performance and musical metadata, and has developed several innovative technologies and award-winning systems (MusicSpace—constraint-based spatialization, PathBuilder—intelligent music scheduling using metadata, the Continuator—interactive music improvisation). François Pachet is the author of more than 60 scientific publications in the fields of musical metadata and interactive instruments.



Jean-Julien Aucouturier graduated from Ecole Supérieure d'Electricité (France) as an electronics and computer engineer and holds a M.Sc. in Music Signal processing from King's College, University of London, U.K. He is currently a Ph.D. student in University of Paris, France and research assistant in Sony Computer Science Laboratory (Paris), where he investigates interaction systems with large music databases.



Amaury La Burthe is a Research Engineer and holds a M.Sc in Signal Processing, Acoustics and Computer Science applied to music from Ircam, Paris. His research interests concern music access system, music descriptions, and representations. He is now working as a lead sound designer for a video game company where he focuses on sound design for games and tools towards the implementation of complex music system.



Aymeric Zils is a Signal Processing Engineer from Ecole Centrale (Nantes, France), and holds a M.Sc. in Acoustics from Laboratoire d'Acoustique de University du Maine (Le Mans, France). He obtained a Ph.D. in computer music from University of Paris 6, for his research at Sony Computer Science Laboratory (Paris) on the automatic extraction of musical descriptors.



Anthony Beurivé holds a M.Sc. in Computer Science from the Bordeaux 1 University, Talence, France. He is currently a research assistant in Sony Computer Science Laboratory (Paris), where he helps in the development of new systems for multimedia management.

UNCORRECTED PROOF

AUTHOR QUERIES

AUTHOR PLEASE ANSWER ALL QUERIES.

- Q1. Please provide keywords.
- Q2. Please check insertion of Figure 2 citation if appropriate.
- Q3. Please provide journal abbreviated title.
- Q4. Please update bibliographical information if this has already been published.
- Q5. Reference [6] is uncited. Please check.
- Q7. Please provide pages
- Q7. Please provide year of publication.

UNCORRECTED PROOF