

## AUDIOCYCLE: BROWSING MUSICAL LOOPS LIBRARIES

Stéphane Dupont<sup>1</sup>, Nicolas d'Alessandro<sup>1</sup>, Thomas Dubuisson<sup>1</sup>, Christian Frisson<sup>2</sup>, Raphaël Sebbe<sup>1</sup>, Jérôme Urbain<sup>1</sup>

<sup>1</sup> Laboratoire de Théorie des Circuits et Traitement du Signal (TCTS), Faculté Polytechnique de Mons (FPMs), Belgique

<sup>2</sup> Laboratoire de Télécommunications et Télédétection (TELE), Université Catholique de Louvain (UCL), Belgique

### ABSTRACT

This paper presents AudioCycle, a prototype application for browsing through music loop libraries. AudioCycle provides the user with a graphical view where the audio extracts are visualized and organized according to their similarity in terms of musical properties, such as timbre, harmony, and rhythm. The user is able to navigate in this visual representation, and listen to individual audio extracts, searching for those of interest. AudioCycle draws from a range of technologies, including audio analysis from music information retrieval research, 3D visualization, spatial auditory rendering, audio time-scaling and pitch modification. The proposed approach extends on previously described music and audio browsers. Concepts developed here will be of interest to DJs, remixers, musicians, soundtrack composers, but also sound designers and foley artists. Possible extension to multimedia libraries are also suggested.

### 1. INTRODUCTION

Music production is more and more relying on pre-recorded material. One of the dominant musical instruments synthesis paradigm is based on libraries of high-quality recordings of real instruments, covering their different articulations and expressive modes, leading to very large synthesis libraries, close to the gigabyte for a single instruments, with physical modeling probably being at the other extreme in terms of computation/memory tradeoff [36]. Also, the production of many music styles heavily relies on pre-recorded musical phrases, organized into libraries. This is most apparent for instance in styles like hip-hop or remixing, where these phrases, often referred to as “loops” (for instance drum loops), are indeed looped, sequenced, and mixed together to form full audio tracks [40] and compositions. Third, granular synthesis techniques have also been successfully used in music creation, as a way for discovering new musical textures, or in live settings [38, 53].

However, the tools available today for browsing through large musical libraries hinders the creative process. Indeed, loops can be searched through rigid file system hierarchies, or through the use of symbolic descriptors (music style, instruments used) stored as meta-data in the library. Also, this processes, being mostly offline, can not be used during performances. In the music performance area, some DJ products provide assistance to the performer, for instance through real-time tune synchronization capabilities, but searching through the library is still being done the usual way [49]. More generally, with the growing availability of multimedia content, there is a still larger demand for more flexible and efficient tools to access content and search for data. Information indexing and retrieval can rely on automatic technologies to describe contents on one hand, and on the other hand allow formulating queries, and structure the media database to help the user

to navigate through it. The work presented here envision a convenient and fast way of exploring large audio and music extract libraries, relying on similarity analysis and musically relevant audio analysis and transformation, involving rhythm, harmony and timbre [24, 34, 23, 46]. Other features that are specific to the application to music content have also been implemented. Loops can be played in a beat-synchronous fashion, relying on a phase vocoder algorithm to adjust the tempo without affecting the pitch. AudioCycle differs from previously published approaches. Compared to SoundTorch [20] for instance, it provides a range of audio analysis approaches (not only timbral) and enables to weight their importance. On the usability and rendering side, AudioCycle enables to play synchronously any combination of loops, even if they are very distant on the similarity map.

The paper is organized to describe the different blocks of the AudioCycle architecture, schematized in Figure 1. First, an *Audio Analysis* (2) is performed on a set of loops ( $I$ ) loaded by the user. It consists of extracting features representative of the musical properties of rhythm, harmony, and timbre. This is presented in Section 2. The specificities of the algorithms that have actually been implemented are also described. This step generates a database gathering features and meta-data (including tempo and key found in some audio file formats) for each loop (3). Section 3 is interested by *Visualization* (4) techniques that can help users to navigate and search inside large digital media libraries. These techniques rely on the previously created Features Database (3), and approaches for mapping this large dimensional space to a 2D or 3D structured representation of the library content are explored. Users can interact with this representation, influencing the subsequent *3D Video Rendering* (5) - also addressed in Section 3 - and *3D Audio Rendering* (6), which spatializes, synchronizes and plays the selected loops - presented in Section 4. Section 5 provides an overview on the development started to improve the user interaction. A discussion, including improvements to be implemented in terms of visualisation, and suggestions for future activities, is finally proposed in Section 6.

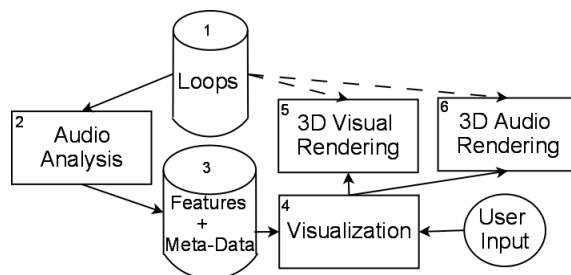


Figure 1: AudioCycle Architecture

## 2. MUSIC ANALYSIS

To characterize music chunks and organize them according to their similarities, a set of relevant mathematical descriptors is needed. They will be introduced hereunder. Then, methods for computing an important rhythmic notion, the tempo (expressed in Beats Per Minutes - *BPM*), and estimating the downbeat positions in a song will be described. Localising the downbeats is useful for appropriately cutting a music file into bars, which are not always similar, and synchronizing musical segments. To conclude this section, results of experiments assessing the relevance of the chosen timbral features are presented.

### 2.1. Feature Extraction

Following the approach of music analysis proposed in [8], features providing information about the timbre, the harmony and the rhythm of music excerpts are extracted from frames whose length is 30ms and whose hopsize is 10ms.

#### 2.1.1. Timbre

Timbre is defined as the quality of tone distinctive of a particular singing voice or musical instrument [27]. It is linked to the relative intensities of a sound harmonics, independantly from its pitch and intensity. It has been shown in many studies ([23, 24, 33]) that timbre can be efficiently characterized by the instantaneous spectral envelope of the sound, using a non-linear frequency scale. A popular non-linear frequency scale is the Mel scale, built to model the human perceptions of pitches. Based on this scale, the frequency content of a frame can be summarized by the Mel-Frequency Cepstral Coefficients (MFCCs) [52]. To characterize the timbre in our application, we use a filterbank of 20 filters covering the frequency range between 0Hz and the Nyquist frequency (22050Hz for the music excerpts we used), and keep only the first 12 MFCCs (+ the energy).

#### 2.1.2. Harmony and Melody

Harmony can be defined as the combination of simultaneous musical notes in a chord [27]. In addition, the melody is defined as a rhythmic succession of single tones organized as an aesthetic whole [27]. In consequence, to characterize the harmony and melody of a musical signal, features revealing the presence of particular notes, their evolution and combination must be extracted. Since notes are defined by frequency values, harmony and melody analysis is generally performed in the frequency domain. A widely used method is to project the power spectrum of each frame on a chromatic scale, which decomposes each octave into 12 equally broad bands on a base 2 logarithmic scale, called "semitones". The "chromas" characterizing the harmony and melody of a frame are then the 12 sums of the energetic contributions of each semitone on all the (considered) octaves. They represent the combination of notes played at a given time, independently from their octave.

#### 2.1.3. Rhythm

Rhythm is defined as the aspect of music comprising all the elements (as accent, meter, and tempo) that relate to forward movement [27]. From a signal processing point of view, it is closely related to the notion of periodicity [8]. Various analyses of a signal periodicity can be proposed. For characterizing the rhythm of musical excerpts, one popular way is to determine onsets instants

- which can also serve to decompose a musical chunk into bars, beats, etc. To estimate these onset instants, it is common to compute the Perceptual Spectral Flux (PSF) [23], which can be seen as an integration of the local derivatives of the spectrogram frequency bins. The algorithm proposed in [24] for computing the PSF, also called "onset function" is:

- To compensate for human ear attenuation, weight the power spectrum of each frame with the inverse of the Equal Loudness Curve (EQL) defined at 40 phons [21] (see Eq. 1 where  $f$  denotes the frequency).

$$EQL(f) = \left( \frac{f^2}{f^2 + 1.6 * 10^5} \right)^2 * \left( \frac{f^2 + 1.44 * 10^6}{f^2 + 9.61 * 10^6} \right) \quad (1)$$

- Compute the Perceptual Spectral Flux:

$$PSF(n) = \sum_{k=1}^{N_b/2} ((a_n^k)^{1/3} - (a_{n-1}^k)^{1/3}) \quad (2)$$

where  $n$  stands for the frame index,  $N_b$  for the number of frequency bins and  $a_n^k$  for the compensated power spectrum at the  $k^{th}$  frequency bin of current frame  $n$ . The exponent  $1/3$  is introduced to simulate the intensity-loudness power law. As in [12], the power spectrum was obtained by applying a Mel filterbank decomposition on top of the usual Fourier Transform.

In [24] and [8], the rhythmic information is computed by calculating local autocorrelation functions on overlapping frames of the PSF.

Another approach has been designed here, based on the notions of bar, beat and tatum, which are the onset events related to the perception of rhythm (see Section 2.2). We have chosen to save as rhythmic features the onset values in the neighbourhood of the beats and tatums positions for one bar. The window length is set to 11 frames: PSF values are summed from 5 frames before and until 5 frames after the metrical event position. Negative onset values are discarded from the sum, since rhythmic events are perceived when a particular sound is striking, not disappearing. Since beats are mainly characterized by low-frequency contributions of the onset function, it has been decided to split the onset function in three frequency bands: low frequencies (0 to 100 Hz), medium frequencies (100 to 8000 Hz) and high frequencies (8000 Hz to the Nyquist frequency).

A 4/4 rhythmic signature is assumed, meaning a bar contains 4 beats - the first one being the downbeat - and 4 tatums, located halfway between two beats. Thus, to capture the beat and tatum information of one 4/4 bar, features around 8 particular instants must be computed. The vector characterizing the rhythm of each music excerpt is in consequence composed of 8 values for each frequency band. The frequency band decomposition of the onset function and mask applied for computing the rhythm features are illustrated in Figure 2. It can be seen that the mask - generated independently from the onset function - will catch the main peaks of the onset functions, corresponding to beats and tatums. Only the first 8 steps of the mask are used to compute the rhythm features.

The described method requires the knowledge of the downbeat position and the interval between two successive beats. The latter can be easily computed from the Beats Per Minute (BPM) value, frequently available as meta-data with professional music files. Loops used for DJing applications generally start at the downbeat, solving the problem of localizing this event. When the BPM and

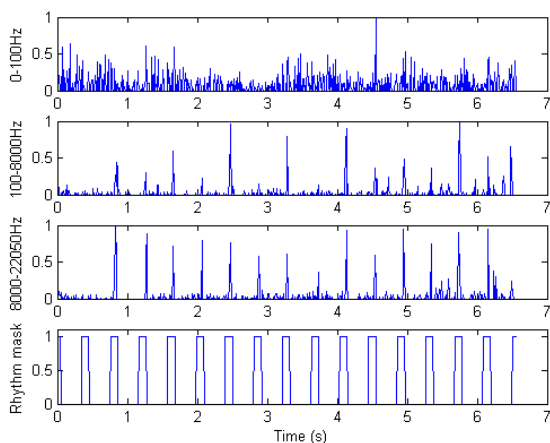


Figure 2: Top 3 graphs: band frequency decompositions of the onset function (normalized amplitudes). Bottom graph: Mask applied for computing the rhythm features.

downbeat position are a priori unknown, it is still possible to compute them and use the proposed method for extracting the rhythm features thanks to the algorithms that will be described in section 2.2.

#### 2.1.4. Summarizing the information

The targeted application must be able to deal with thousands of musical extracts, fastly compare them and display them in a suitable way (see Section 3). This objective has two implications on the feature sets we provide. First, the characteristic vectors must be of reasonable length, otherwise computation times would dramatically increase and the application would not be efficient in real-time. Second, characteristic vectors of each musical excerpt should have a constant length, independently from its duration, in order to facilitate comparisons of the signal properties.

Regarding the rhythm characterization, the method aforementioned already provides a constant number of 24 rhythmic features per music excerpt. However, as the number of timbre and chroma features is proportional to the music excerpt duration, this information is summarized into fixed size vectors by storing only the mean and variance of each timbre and harmony feature. It is shown in section 2.3 that even if information is lost when summarizing, the mean and the variance of the features still contain enough information to efficiently recognize the kind of instruments involved in the extract, which is the aim of timbre characterization.

## 2.2. Beat and Downbeat Estimation

In the Music Information Retrieval domain, the notion of "beat" is often expressed by musicians and listeners through the act of foot tapping along to the music. Downbeat is defined as the first beat mark at the beginning of a bar.

Beat tracking systems form the basis of applications like automatic accompaniment, transcription, computer-assisted audio editing and evaluation of music similarity.

### 2.2.1. State of the Art

The onset function (see Section 2.1) is commonly used to extract the tempo (expressed in BPM). The method proposed in [12] estimates it by computing the autocorrelation of the onset function and by favoring autocorrelation values located in a range defined by the user. It is performed by weighting the autocorrelation function by a Gaussian window centered around the most likely BPM value and whose spread is set by the user. Once the largest peak is located in the weighted autocorrelation, a second value of BPM is found as the largest value of autocorrelation at 0.33, 0.5, 2 or 3 times the first value of BPM. This system was tested on the database of MIREX 2006 contest [29] and achieved 77 % of agreement with the manual BPM.

Among the other beat tracking methods based on the onset function, the system exposed in [32] computes the onset function from a reassigned spectrogram, in which the energy of the bins (at frequency  $\omega_k$  at time  $t_i$ ) are reassigned to the frequency  $\omega_r$  and time  $t_r$  corresponding to their center of gravity. BPM is then estimated by combining periodicity measures of the onset function (Discrete Fourier Transform and Autocorrelation). In [15] the structure of a musical excerpts (in terms of bar, mid-bar and beat) is estimated by computing an onset function from different frequency bands in the spectrogram. Chord-change possibilities and drum patterns are then estimated at beat locations and the structure of the music excerpt is inferred from these informations.

Fewer studies focus on the difficult problem of downbeat estimation although it has a perceptual sense for us. A method for downbeat extraction based on beat locations and BPM is presented in [9]. The signal is cut into segments synchronized at beat locations and whose length is the beat period. The difference between the spectrum of successive segments is then computed. For each bar, the downbeat is located where value of this difference is highest. In [22] Support Vector Machines are trained in order to predict the downbeat locations by using spectral features synchronized on the local maxima of the onset function. In the testing phase, the system takes as input the spectral features of several music segments and provides an estimation of the downbeat locations in the next segments.

### 2.2.2. Implementation

As the onset function appears to be very popular for beat period estimation, it has been chosen in the present study. Onset functions from [12, 10, 15, 32] have been compared and appeared quite similar. Because Ellis' method provides good performance and because the whole source code is provided on his website, it has been chosen for beat period estimation in our system.

Concerning the downbeat estimation, a method inspired by previous literature [22] has been developed. The idea is to find the best alignment between a template of bar, beat and tatum locations and the onset function. As in [22], it has been chosen to consider only templates for 4/4 time signature in this study.

The template is built using the estimated BPM value. An example is shown in Fig.3 for 144 BPM. The best alignment is achieved by computing the cross-correlation between the onset function and a lagged version of the template. This correlation is computed for lags from 0 to one bar period. The lag corresponding to the highest correlation value is used to set the downbeat.

The algorithms implemented for computing the beat period and localizing the downbeat position also enable us to perform "slicing" (i.e. decomposition into bars) of audio files. Not only loops but also full songs can thus be analyzed by the application.

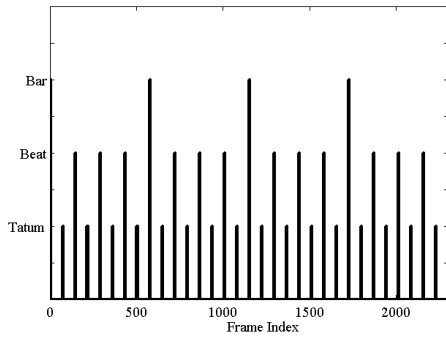


Figure 3: Template for 144 BPM

The algorithm designed for slicing consists of first cutting the files into 15s-windows (with a 10s-overlap), then launch the computation of the beat period and finally localize the downbeats position in the central 5s of the window. Resulting bars are then processed by the application as if they were loops.

### 2.3. Instruments Classification

In order to define which set of features is the most appropriate for the similarity estimation and visualization in AudioCycle, their performance for classifying instruments has been assessed. The mean and variance of the timbre features have been considered in this evaluation.

#### 2.3.1. Database

The features are extracted from musical loops of the ZeroG ProPack ACID database [51]. This database is a product of ZeroG Ltd. and contains 11428 loops and samples of various instruments in diverse musical styles. We manually tagged the files using 7 instruments classes: Brass, Drums, Vocals, Percussion, Electric Bass, Acoustic Guitar and Electric Guitar. Only the files for which the instrument class is clear have been taken into account in this study.

#### 2.3.2. Dimension Reduction

Before classifying the feature vectors, they have to be manipulated in order to assess if their dimension can be reduced. For this purpose, the Generalized Fisher criterion and Principal Component Analysis (*PCA*) have been tested.

The Generalized Fisher criterion [14] is a class separation criterion based on the features and the class labelling. For a given feature  $k$ , its discriminant power between  $C$  classes is defined as

$$D_k = \frac{\sum_{c=1}^C p(\omega_c) (\mu_{ck} - \mu_k)^2}{\sum_{c=1}^C p(\omega_c) \sigma_{ck}^2} \quad (3)$$

where  $p(\omega_c)$  stands for the percentage of representation of class  $c$  in the database,  $\mu_{ck}$  for the mean of the feature  $k$  in the class  $c$ ,  $\mu_k$  the mean of feature  $k$  for all the classes and  $\sigma_{ck}$  for the standard deviation of feature  $k$  in the class  $c$ . A dimension reduction is possible by selecting the features associated to the highest values of discriminant power.

Principal Component Analysis (*PCA*) [14] prepares the features in a classification system by linearly transforming them in order to find the best representation space (in terms of least square

error). If  $X$  represents the normalized features matrix, the new features matrix  $Z$  is obtained by

$$Z = U^T X \quad (4)$$

where  $U$  is the linear transformation matrix. The matrix  $U$  leading to the best final representation consists of the eigen vectors of the covariance matrix  $XX^T$ . The dispersion of features around each new axis of representation is given by its associated eigen value. A reduction of features dimensionality is possible by selecting the axis of representation associated to the highest eigen values.

#### 2.3.3. Results

Three ways of reducing the dimension of features vector have been investigated: applying the Fisher criterion alone, applying *PCA* alone and applying Fisher criterion followed by *PCA*. These different approaches have been applied on the mean and variance of the timbre features. For each configuration, a Multi-Layer Perceptron was trained with 60% of the objects, validated with 10% and tested with the remaining 30%. The operation was repeated 100 times and the objects were always randomly selected. Table 1 shows the average classification performance for each configuration of features.

Table 1: Performance of classification (in %)

Description	Mean	Var
No selection	73.1	50.8
6 most discriminant after Fisher	66.2	52.1
2 most discriminant after Fisher	54.3	47.2
No selection after <i>PCA</i>	93	70.9
6 most discriminant after <i>PCA</i>	85	67.7
2 most discriminant after <i>PCA</i>	66	52.6
No selection after Fisher and <i>PCA</i>	93	68.1
6 most discriminant after Fisher then <i>PCA</i>	86.9	59.8
6 most discriminant after Fisher then 2 most discriminant after <i>PCA</i>	75	56.3
2 most discriminant after Fisher then 2 most discriminant after <i>PCA</i>	64.1	51.9

Classification performance for timbre information is very good (up to 93%), supporting our decision to summarize the timbre information 2.1.4 and confirming that these vectors still contain enough information to distinguish instruments with a very high accuracy. Since performance is significantly better using means than variances and further tests shown that no improvement was achieved by combining means and variances, only timbre means are finally used for the AudioCycle application. Similar tests are under way to assess the opportuneness of summarizing the chroma information the same way.

## 3. VISUALIZATION

### 3.1. Visualization Techniques

The representation of large multimedia database has been a subject of research in various cases of media content. As stated in [45], the visualization has to meet three requirements: overview (faithful

overview of the distribution of excerpts in the collection), structure preservation (the relations between music excerpts should be preserved in the projection in the visualisation space) and visibility (visually understanding of the content of each music excerpts). Thus if the similarity between music excerpts is used, the ideal is to build a representation in which the similarity distances are preserved, what our choice of the node-link visualization paradigm, great for emphasizing two Gestalt laws (similarity and connectedness) [48], ensures.

When displaying the whole collection of music excerpts is desired, the idea is to associate to an object  $x_i$  a vector  $y_i$  aiming at transposing the similarity between  $x_i$  and its neighbours in a 2D or 3D representation. Three kinds of approaches are described hereafter.

### 3.1.1. Dimensionality Reduction

The goal of this approach is to preserve the distance between points in a high dimensional space when they are projected in a subspace of two dimensions. This can be achieved by e.g. PCA or by Linear Discriminant Analysis (LDA) [11]. This latter projects the features vectors from a  $d$ -dimensional space to a  $D$ -dimensional space (with  $D < d$ ) with the constraint to maximize the ratio between the between-class covariances and the within-class covariances of the new features. This is done by applying the linear transformation

$$Z = W^T X \quad (5)$$

where  $W$  is the transformation matrix whose dimension is  $d \times D$ .

Among the LDA techniques, one can cite the Fisher method which consists on finding  $C - 1$  linear discriminant functions in a problem of  $C$  classes. It can be proved that the columns of  $W$  that maximises the ratio between the between-class covariance and the within-class covariance correspond to the highest eigen values of

$$S_B W_i = \lambda_i S_W W_i \quad (6)$$

where  $W_i$  stands for a column of  $W$  and  $\lambda_i$  for the  $i^{th}$  eigen value. As only  $C - 1$  eigen values are different of 0, the new system of representation consists of only  $C - 1$  dimensions.

### 3.1.2. Graphs and Multidimensional Scaling

When using graphs for visualization, each excerpt is represented by a vertex in the graph and the arcs between neighbouring excerpts are used for navigation. The objective of the graphs is to preserve the neighbourhood of excerpts and to enhance data structures. For retrieval and analysis purposes, the radial tree structure and global graph are particularly adapted. The radial tree structure provides a focus on a part of a graph, analogous to a hierarchical representation. The global graph aims at representing a large amount of excerpts with various levels of details, but with a larger point of view than radial tree structures.

Multi Dimensional Scaling (MDS) aims at trading with the high dimension of data by using graphs. The objective is to iteratively optimize a function that measures the difference between the distance of two excerpts in the 2D space and their original dissimilarity. In order to build a correct visualization. The original method requiring a high computational cost, other approaches have been developed. Among them, one can cite Isomap, aiming at clustering first the data before applying MDS or Visumap, which maps the excerpts in a 2D map so that Euclidian distance visually approaches as much as possible the dissimilarities of the excerpts.

### 3.1.3. Clustering

Clustering is an other way to organize a large amount of data. An original approach has been developed here. It relies on a K-Means algorithm. Clusters of loops are built by iteratively constructing partitions through association of each loop to the closest centroid. The Euclidian distance has been chosen and based on the extracted features defined earlier. The emphasis of the clustering can be changed by the user by scaling differently the feature space dimensions related to timbre, harmony, and rhythm. The number of cluster is chosen in advance by the user and, in order to preserve near real-time interaction, the clustering is stopped before convergence after a limited number of iterations. The initial cluster centers are randomly chosen among the loops but this random sequence is made reproducible to provide back and forth navigation.

Clusters can be represented visually in a flower-like scheme, where the center is the currently selected loop and other loops are gathered into groups (clusters) around the selected one. The 2D radial distance between a loop and the selected centered element is proportional to the distance between them (in a feature space with scaling factors complementary to the one used for clustering), while the distance between the loop and the cluster center (in the scaled space used for clustering) is used to compute the 2D angular coordinate to locate it in the cluster. This scheme allows to explore organization schemes where the radial and angular coordinates are associated to complementary musical properties, e.g. timbre and rhythm.

## 3.2. 3D Imaging, OSG and OpenGL

Visualization of the audio loops is made with Open Scene Graph (OSG) [31], a scene graph visualization library running on top of OpenGL and permitting the definition of high level objects. These objects can either render geometry or represent 3D transformations, and, when associated as a graph, they define a rendered image.

The scene graph is built by adding cubes representing each loop. Each of these cubes has a parent node that specifies its position, which can change over time, when the user changes the clustering properties for instance. Loops can be activated by clicking on the associated objects (cubes) with the mouse. This process is known as picking and is implemented in the usual OSG way by computing intersection with a ray fired at the mouse location. When a loop is activated, the scene graph is enhanced with a waveform display as well as a cursor indicating the current playback position. A screenshot of the AudioCycle OSG view is displayed in Fig. 4. OSG will allow us to extend this work to 3D representations very naturally.

## 4. AUDIO RENDERING

In AudioCycle, several audio rendering features have been designed to facilitate the browsing through large libraries of audio loops. First, the loops are positioned in space in relation with the position of their visual representation on the screen. This tight coupling between the audio and the image is expected to make it easier to locate the area of interest in the visual representation [13, 20], when the user is looking for specific sounds while several loops are playing together. This also allows to play simultaneously several sounds. This has been implemented using OpenAL (described hereafter) with a specific extension providing HRTF-based (Head

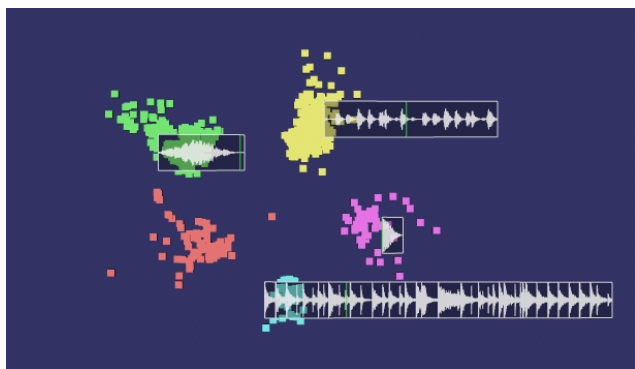


Figure 4: AudioCycle OSG View

Related Transfer Functions) rendering. Then, as we are dealing with music content, it makes sense to provide the application with facilities to allow beat synchronous playback (through audio time-scaling) of loops that have originally not been played at the same tempo, and key normalization (through pitch shifting) for loops that have originally not been played in the same keys. This is expected to reduce the cognitive load when several loops are playing together, hence facilitating the search for specific sounds. Eventually, this may also provide the basis for extending AudioCycle to a performance tool, where live composition and performance is made possible through beat and key matching.

OpenAL [30] is a cross-platform audio API that has been designed for rendering of multichannel 3D audio. Designed with a style similar to OpenGL, it has mostly been used for 3D modeling and game engines. The API allows to define *source* objects and a *listener*. Sources can be specified using several properties, including position, velocity, sound direction and intensity. The listener is also characterized by properties including position, velocity, direction, as well as a global gain. The buffers containing the audio samples in PCM format are associated to the sources by other API calls. The OpenAL engine then performs the rendering and simulates physical effects such as the attenuation caused by distance between the sources and the listener, as well as the Doppler effect.

Additional functionalities can be provided through extensions to the API. Among the available plugins, an AU (Audio Unit) plugin making use of HRTF when rendering for a stereo output has been used in the application. This conveys an enhanced realism to the sound scene.

Coming back to sources, it's of significant importance to note that they be of two types: *static* or *streaming*. As static sources are not appropriate for applications where expressive control is required, streaming sources have been preferred. It allows to pass buffers that are queued for real-time playback. Since the API does not provide a callback mechanism that requests audio buffer, this aspect has hence to be managed by the application. In this work, this is done using a thread with the necessary priority level to periodically check whether buffers will be needed by the OpenAL engine.

## 5. USER INTERACTION

### 5.1. Gestural input

Gestural input theory has been long studied, concerning generic computerized applications [3] to more specific computer music

applications [47]. We have inspired ourselves by the techniques developed in these fields.

#### 5.1.1. Available modalities and controllers

Here is an overview of miscellaneous modalities, or ways to interact with a computer, that have been employed in applications dedicated to navigation in multimedia databases:

- 3D positioning with remote controllers [20, 43] such as the cheap Nintendo Wii;
- Glove [19], by wavering fingers;
- Query by beat-boxing [25], by indicating rhythm and gross timbre mimicked by the mouth;
- Haptics and force-feedback [4, 1], controllers that response to the user movements by applying forces conditioning it;
- Tactile feedback [28], vibration sensed through the skin;
- Tangibles [2, 5, 16, 17], in other words objects that can be grasped, moved, manipulated.

#### 5.1.2. Bimanual interaction

The gestural control for the navigation in an audio database has mostly been introduced and improved by DJ's [42, 49] and the development of these practises is still ongoing [26], with a tight relation to the use of the hands.

Now that a vast content of audio files can be stored on the computer, we can define two modes of interaction for the use case of AudioCycle:

1. navigation: the action of selecting and annotating files to be stored in a library and/or to be retrieved from it;
2. performance: the action of listening to one or more files with an optional addition of real-time temporal and/or spectral effects for both artistic/aesthetic and mining purposes.

Bimanual interaction, the use of both hands to perform tasks, is deeply rooted in the history of musical instruments and is being applied in the computer music field, for instance in the following recent works [7, 37]. To design the user interaction for the AudioCycle project, we have decided to focus on and study this paradigm: the hands can be assigned, according to the CARE properties [6], each to one of the two interaction modes (one hand selects an audio file, the other performs an audio file); or freely used for any mode. In the case of an audio mining task, in order to obtain the fastest search possible, assigning each hand to one mode could help to cognitively accelerate the process. In the case of a musical performance, swapping hands on controllers is left to the performers' will.

## 5.2. Our prototype

#### 5.2.1. Chosen controllers

We have decided to try using two categories of gestural controllers to build our hardware interface, as illustrated in figure 5:

1. 3D mice (3Dconnexion Space Navigator) to navigate in the library;
2. "Jog wheels" (Contour Design Shuttle) to alter the time and pitch of the listened files.



Figure 5: Gestural controllers chosen to be tested in AudioCycle: 3D mice on the left and jog wheels on the right

For an audio mining task, one of each can be featured on the setup and assigned to one hand, to reduce the span of the user movements on his desktop. In the case of a musical performance, the performer could decide to use as many 3D mouse / jog wheel couples as necessary in order to be able to perform a certain number fo loops simultaneously.

Additionally, standard sliders available on “MIDI fader boxes” can be assigned to set the percentage of similarity of each feature.

### 5.2.2. OSC namespace

In order to prototype the user interface of Audio Cycle, that is testing if the use of the aforementioned hardware controller couple is relevant, and trying other controllers, we have started implementing an OpenSoundControl (OSC) [50] bridge in the AudioCycle application, so that to allow pre-formated messages to be sent through a network (that can be hosted on the same computer), from the interaction prototyping application to the AudioCycle application that will interpret theses messages.

As it is now done in many applications requiring gestural input [39], we have defined our own OSC namespace to control the AudioCycle application as depicted in figure 6:

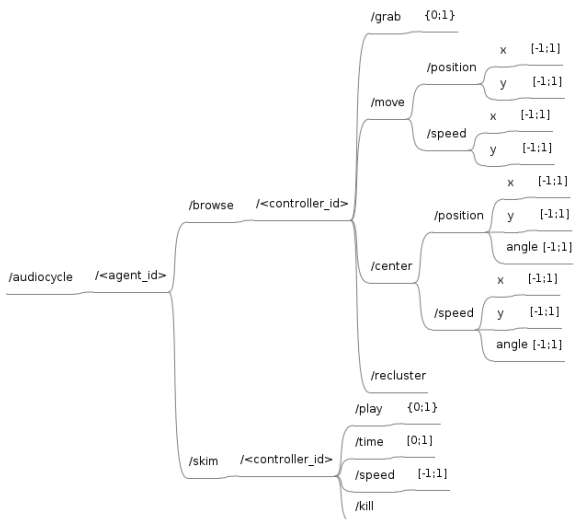


Figure 6: Mindmap view of the AudioCycle OSC namespace

For example:

- `/audiocycle/1/browse/2/grab 1` selects the closest audio file to the position set by navigation controller number 2 on the graph,
- `/audiocycle/1/skim/3/play` plays the audio file previously selected and assigns to it performance controller number 3,
- `/audiocycle/1/browse/2/recluster` reorganizes the view around the audio file selected previously,
- `/audiocycle/1/skim/3/time 0.5` repositions the playhead to half of the duration length of the selected audio file,
- and so on...

At the current stage, AudioCycle supports interaction with a single user (`<agent_id>` is set to one) but the namespace is open to the implementation to collaborative and concurrent use.

### 5.2.3. Prototyping the user interface

To quickly implement and verify gestural control prototypes, modular environments such as PureData [35], using `[hid]` and `[hidio]` objects under Linux and OSX [41], or OpenInterface [44] under Linux and Windows, soon OSX, are being considered for the rapid prototyping. Both environments can act as OSC server and generate OSC messages to be sent to the AudioCycle application. This configuration allows to study various mappings [18], or ways to interconnect gestural parameters and application events, applied to our use case.

## 6. DISCUSSIONS

Some aspects of the visualisation require further study [48], for instance, respecting the symetry of the 3D node-link diagram paradigm and condensing the mapping of the waveforms around each loop representation.

In order to achieve the fastest possible audio mining, we initiated a user-centered usability approach, featuring context enquiries with experts (sound designers [36], DJ's [49], etc...) and usability tests with software/hardware prototypes, leading to the definition and support of different user and hardware profiles.

This work is part of an ongoing series of projects intending to extend our approach to multimedia content (images, videos, text, etc...) with scalability to very large archives.

## 7. CONCLUSIONS

A novel approach and prototype implementation of a music loop library browser has been presented. It allows to display the music extracts organized on a map according to timbre, harmony, rhythm, or a combination of these. Spatial sound is also used to enhance the experience when playing several extracts at the same time. Features relevant to music production have also been integrated, including the possibility to play the loops in a beat-synchronous fashion, as well as alter their pitch. Some avenues for extending and improving the concept and technologies involved are also proposed.

## 8. ACKNOWLEDGMENTS

This work has been supported by the numediart research project, funded by Région Wallonne, Belgium (grant N°716631).

## 9. REFERENCES

### 9.1. Scientific references

- [1] Timothy Beamish, Karon Maclean, and Sidney Fels. “Manipulating Music: Multimodal Interaction for DJs”. In: *Proceedings of CHI’04*. 2004. P.: 124.
- [2] Henry Bernard. “Tangible Media Browsing: Exploring And Browsing Large Multimedia Collections With A Tabletop”. MA thesis. Music Technology Group, UPF-IUA, Barcelona and INPG-ICA, Grenoble, 2007. P.: 124.
- [3] S. K. Card, J. D. Mackinlay, and G. G. Robertson. “The design space of input devices.” In: *Proceedings of the SIGCHI conference on Human factors in computing systems (CHI’90)*. 1990. Pp. 117–124. P.: 124.
- [4] Lonny L. Chu. “Haptic Interactions For Audio Navigation”. PhD thesis. Stanford University, 2003. P.: 124.
- [5] E. Costanza, S. B. Shelley, and J. Robinson. “Introducing Audio D-Touch: A Tangible User Interface For Music Composition And Performance”. In: *Proc. of the 6th Int. Conference on Digital Audio Effects (DAFX-03)*. 2003. P.: 124.
- [6] Joëlle Coutaz et al. “Four Easy Pieces for Assessing the Usability of Multimodal Interaction: The CARE properties”. In: *Proceedings of the INTERACT’95 conference*. 1995. Pp. 115–120. URL: [http://iihm.imag.fr/publs/1995/Interact95\\_CARE.pdf](http://iihm.imag.fr/publs/1995/Interact95_CARE.pdf). P.: 124.
- [7] Jean-Michel Couturier and Magnolya Roy. “Grapholine, Instrument Audiovisuel De “Dessin Musical””. In: *Articles des Journées d’Informatique Musicale*. 2008. P.: 124.
- [8] Laurent Couvreur et al. “Audio Thumbnailing”. In: *QPSR of the numediart research program*. Ed. by Thierry Dutoit and Benoît Macq. Vol. 1. 2. <http://www.numediart.org>. 2008. Pp. 67–85. P.: 120.
- [9] M. Davies and M. Plumbley. “A Spectral Difference Approach to Downbeat Extraction in Musical Audio”. In: *Proceedings of EUSIPCO 2006*. 2006. P.: 121.
- [10] M. Davies and M. Plumbley. “Beat Tracking with a Two State Model”. In: *Proceedings of ICASSP 2005*. 2005. P.: 121.
- [11] R.O. Duda and P.E. Hart. *Pattern Classification and Scene Analysis*. John Wiley and Sons, 1973. P.: 123.
- [12] D. Ellis. “Beat Tracking with Dynamic Programming”. In: *Proceedings of MIREX 2006*. 2006. Pp.: 120, 121.
- [13] M. Fernström and E. Brazil. “Sonic browsing: An auditory tool for multimedia asset management”. In: *International Conference on Auditory Display (ICAD2001)*. Helsinki, Finland 2001. P.: 123.
- [14] K. Fukunaga. *Introduction to Statistical Pattern Recognition*. Ed. by Werner Rheinboldt. Academic Press - Second Edition, 1990. P.: 122.
- [15] M. Goto. “An Audio-based Real-Time Beat Tracking System for Music With or Without Drums-sounds”. In: *Journal of New Music Research* 30.2 (2001). Pp. 159–171. P.: 121.
- [16] Kjetil F. Hansen, Marcos Alonso, and Smilen Dimitrov. “Combining Dj Scratching, Tangible Interfaces And A Physics-Based Model of Friction Sounds”. In: *Proceedings of ICMC07*. 2007. P.: 124.
- [17] Kjetil Falkenberg Hansen and Marcos Alonso. “More DJ techniques on the reactable”. In: *Proceedings of the International Conference on New Interfaces for Musical Expression (NIME)*. 2008. P.: 124.
- [18] Kjetil Falkenberg Hansen and Roberto Bresin. “Mapping strategies in DJ scratching”. In: *Proceedings of the 2006 International Conference on New Interfaces for Musical Expression (NIME06)*. 2006. P.: 125.
- [19] Kouki Hayafuchi and Kenji Suzuki. “MusicGlove: A Wearable Musical Controller for Massive Media Library”. In: *Proceedings of the International Conference on New Interfaces for Musical Expression (NIME)*. 2008. P.: 124.
- [20] Sebastian Heise, Michael Hlatky, and Jörn Loviscach. “SoundTorch: Quick Browsing in Large Audio Collections”. In: *125th Audio Engineering Society Convention*. 7544. 2008. Pp.: 119, 123, 124.
- [21] H. Hermansky, N. Morgan, and P. D. Korn. “Auditory Model for Parametrization of Speech”. Pat. 5,450,522. 1995. P.: 120.
- [22] T. Jehan. “Downbeat prediction by listening and learning”. In: *Proceedings of WASPAA 2005*. 2005. P.: 121.
- [23] Tristan Jehan. “Creating Music by Listening”. PhD thesis. Massachusetts Institute of Technology, 2005. Pp.: 119, 120.
- [24] K. Jensen. “Multiple Scale Music Segmentation Using Rhythm, Timbre and Harmony”. In: *EURASIP Journal on Advances in Signal Processing 2007 7* (2007). Pp.: 119, 120.
- [25] Ajay Kapur, Manj Benning, and George Tzanetakis. “Query-By-Beat-Boxing: Music Retrieval For The Dj”. In: *Proceedings of ISMIR’04*. 2004. P.: 124.
- [26] Takuro Mizuta Lippit. “Turntable Music in the Digital Era: Designing Alternative Tools for New Turntable Expression”. In: *Proceedings of the 2006 International Conference on New Interfaces for Musical Expression (NIME06)*. 2006. P.: 124.
- [28] Roderick Murray-Smith et al. “Stane: Synthesized Surfaces for Tactile Input”. In: *CHI 2008*. 2008. P.: 124.
- [32] G. Peeters. “Time Variable Tempo Detection and Beat Marking”. In: *Proceedings of ICMC 2005*. 2005. P.: 121.
- [33] G. Peeters. “Toward Automatic Music Audio Summary Generation from Signal Analysis”. In: *Proceedings of the International Symposium on Music Information Retrieval (ISMIR)*. Paris, France 2002. P.: 120.
- [34] G. Poliner et al. “Melody Transcription from Music Audio: Approaches and Evaluation”. In: *IEEE Transactions on Audio, Speech and Language Processing* 15(4) (2007). Pp. 1247–1256. P.: 119.
- [36] Martin Russ. *Sound Synthesis and Sampling*. 3rd ed. Music Technology. Focal Press, 2008. ISBN: 9780240521053. Pp.: 119, 125.
- [37] To Yip Sang and Kam Wong. “Multi-User Hand Gesture Based Musical Element Mapping With Traditional Musical Elements”. In: *Proceedings of ICMC’08*. 2008. P.: 124.

- [38] Diemo Schwartz. “Musical Applications of Real-Time Corpus-Based Concatenative Synthesis”. In: *ICMC 2007*. Copenhagen 2007. P.: 119.
- [39] Stephen Sinclair and Marcelo M. Wanderley. “Defining a control standard for easily integrating haptic virtual environments with existing audio/visual systems”. In: *Proceedings of NIME’07*. 2007. P.: 125.
- [41] Hans-Christoph Steiner, David Merrill, and Olaf Matthes. “A Unified Toolkit for Accessing Human Interface Devices in Pure Data and Max/MSP”. In: *Proceedings of the 2007 Conference on New Interfaces for Musical Expression (NIME07)*. 2007. P.: 125.
- [42] John Steventon. *DJing for Dummies*. John Wiley & Sons, 2007. ISBN: 978-0470032756. P.: 124.
- [43] Rebecca Stewart, Mark Levy, and Mark Sandler. “3D interactive environment for music collection navigation”. In: *Proceedings of the Conference on Digital Audio Effects (DAFx)*. 2008. P.: 124.
- [45] Anne Veroust-Blondet. *VITALAS - State of the Art on Advanced Visualisation Methods*. Tech. rep. 2007. P.: 122.
- [46] H. Vinet, P. Herrera, and F. Pachet. “The Cuidado Project: New Applications Based on Audio and Music Content Description”. In: *Proc. of ICMC*. Göteborg, Sweden 2002. Pp. 450–454. P.: 119.
- [47] Marcelo Wanderley and Marc Battier, eds. *Trends In Gestural Control Of Music*. Ircam - Centre Pompidou, 2000. ISBN: 2-8442-6039-X. P.: 124.
- [48] Colin Ware. *Information Visualization: Perception for Design*. 2nd ed. Interactive Technologies. Morgan Kaufmann, 2004. ISBN: 1-55860-819-2. Pp.: 123, 125.
- [49] Stephen Webber. *DJ Skills: The essential guide to Mixing and Scratching*. Focal Press, 2007. ISBN: 978-0240520698. Pp.: 119, 124, 125.
- [50] Matthew Wright. “Implementation and Performance Issues with OpenSound Control”. In: *Proceedings of ICMC 1998*. 1998. URL: <http://www.cnmat.berkeley.edu/ICMC98/papers-pdf/OSC.pdf>. P.: 125.
- [52] F. Zheng, G. Zhang, and Z. Song. “Comparison of Different Implementations of MFCC”. In: *Journal of Computer Science and Technology* 16 (2001). Pp. 582–589. P.: 120.
- [53] A. Zils and F. Pachet. “Musical Mosaicing”. In: *COST G-6 Conference on Digital Audio Effects (DAFX-01)*. Limerick, Ireland 2001. P.: 119.
- [40] Sony Creative Software. *ACID*. <http://www.sonycreativesoftware.com/products/acidfamily.asp>. URL: <http://www.sonycreativesoftware.com/products/acidfamily.asp>. P.: 119.
- [44] *The OpenInterface Platform*. URL: <http://www.openinterface.org>. P.: 125.
- [51] Zero-G. *Pro Sample Library*. URL: <http://www.zero-g.co.uk/>. P.: 122.

## 9.2. Software and technologies

- [27] *Merriam-Webster Dictionary Online*. <http://www.merriam-webster.com/>. Consulted on January 14, 2009. P.: 120.
- [29] *Music Information Retrieval Evaluation eXchange*. URL: [http://www.music-ir.org/mirex/2006/index.php/Main\\_Page](http://www.music-ir.org/mirex/2006/index.php/Main_Page). P.: 121.
- [30] *OpenAL*. URL: <http://www.openal.org/>. P.: 124.
- [31] *OpenSceneGraph (OSG)*. URL: <http://www.openscenegraph.org/>. P.: 123.
- [35] *PureData*. URL: <http://www.puredata.info>. P.: 125.