

## RetroSpat: a Perception-Based System for Semi-Automatic Diffusion of Acousmatic Music

Joan Mouba, Sylvain Marchand, Boris Mansencal, Jean-Michel Rivet

### ▶ To cite this version:

Joan Mouba, Sylvain Marchand, Boris Mansencal, Jean-Michel Rivet. RetroSpat: a Perception-Based System for Semi-Automatic Diffusion of Acousmatic Music. Sound and Music Computing (SMC), Jul 2008, Berlin, Germany. pp.33-40. hal-00351948

## HAL Id: hal-00351948 https://hal.science/hal-00351948

Submitted on 12 Jan 2009

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers. L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# RetroSpat: a Perception-Based System for Semi-Automatic Diffusion of Acousmatic Music

Joan Mouba, Sylvain Marchand, Boris Mansencal, and Jean-Michel Rivet SCRIME / LaBRI – CNRS, University of Bordeaux 1, France

Abstract-We present the RetroSpat system for the semiautomatic diffusion of acousmatic music. This system is intended to be a spatializer with perceptive feedback. More precisely, RetroSpat can guess the positions of physical sound sources (e.g. loudspeakers) from binaural inputs, and can then output multichannel signals to the loudspeakers while controlling the spatial location of virtual sound sources. Together with a realistic binaural spatialization technique taking into account both the azimuth and the distance, we propose a precise localization method which estimates the azimuth from the interaural cues and the distance from the brightness. This localization can be used by the system to adapt to the room acoustics and to the loudspeaker configuration. We propose a simplified sinusoidal model for the interaural cues, the model parameters being derived from the CIPIC HRTF database. We extend the binaural spatialization to a multi-source and multiloudspeaker spatialization system based on a static adaptation matrix. The methods are currently implemented in a real-time free software. Musical experiments are conducted at the SCRIME, Bordeaux.

#### I. INTRODUCTION

Composers of acousmatic music conduct different stages through the composition process, from sound recording (generally stereophonic) to diffusion (multiphonic). During live interpretation, they interfere decisively on spatialization and coloration of pre-recorded sonorities. For this purpose, the musicians generally use a(n un)mixing console. With two hands, this becomes hardly tractable with many sources or speakers.

The RetroSpat system supports artistically interpretation and technically room calibration. It includes a multisource and multi-loudspeaker spatializer, that adapts to different loudspeaker configurations by "listening to the room". This involves source localization and spatialization in azimuth and distance. Here, we focus on the case of a single source with speakers in the horizontal plane.

First, we enhance the binaural model proposed by Viste [1]. We propose to simplify the spatial cues model, resulting in a new sinusoidal model with better mathematical properties and comparable errors using the CIPIC database [2]. Second, we also consider the distance of the source, with a localization based on the brightness.

Last but not least, we extend the binaural spatialization to a multi-loudspeaker spatialization system. In the classic VBAP [3] approach, the control of the interaurallevel difference (ILD) is done in a frequency-independent and pair-wise way that was previously used for source panning. But this method is suitable only for frequencies up to 600Hz. The RetroSpat system also operates on loudspeakers in a pair-wise manner. But the computation of the coefficients for each channel is based on an adaptation matrix of head-related transfer functions (HRTFs), leading to complex and frequency-dependent coefficients. This paper is organized as follows. In Section II, we present some generalities in acousmatic music and we highlight some practical weaknesses to be improved. After an extensive presentation of the model in Section III, we describe the associated spatialization and localization methods in Sections IV and V, respectively. Section VI is dedicated to the presentation of the RetroSpat software.

#### II. ACOUSMATIC MUSIC

#### A. History

Over centuries, the music has continuously undergone various innovations. In 1948, Schaeffer and Henry at the "Radio Télévision Française" were interested in the expressive power of sounds. They used microphones to capture sounds, discs as supports, and transformation tools. The *musique concrète* was born.

In 1949, Eimer gave birth to *electronic music* in the studios of the German radio "Nordwestdeutscher Rund-funk" in Cologne. This music was produced by frequency generators. Koenig and Stockhausen were among the first to use it.

The merge of *musique concrète* and *electronic music* gave rise to *electro-acoustic music* or *acousmatic music*. Today, many musical pieces are created worldwide. Acousmatic has become a discipline that is taught in universities and conservatories.

#### **B.** Actual Practices

Composers of acousmatic music use both electronic and natural sounds recorded close to a microphone, such as wind noise, voices, wrinkling paper, etc. The sounds are then processed by a computer and organized by editing and mixing. The result is a *musical composition*.

However, the creation gets its full value when it is played in concert using an *acousmonium*: an orchestra of loudspeakers. The acousmonium consists of a highly variable number of loudspeakers with different characteristics. The interpreter of the piece controls the acousmonium from a special (un)mixing console.

The originality of such a device is to map the two stereo channels at the entrance to 8, 16, or even hundreds of channels of projection. Each channel is controlled individually by knobs and equalization systems. The channel is assigned to one or more loudspeakers positioned according to the acoustical environment and the artistic strategy.

#### C. Expected Improvements

Behind his/her console, the interpreter of acousmatic music acts in real time on various sound parameters such as spatial location, sound intensity, spectral color. He/She broadcasts a unique version of the music fixed on a medium. The acousmatic diffusion requires some skills. RetroSpat intends to facilitate the work of the interpreter by improving the following embarrassing practices:

- two wheels needed to spatialize one source;
- stereo sources as inputs;
- no individual source path, only one global mix path;
- the distance spatialization requires some expertise.

#### III. BINAURAL MODEL

We consider a punctual and omni-directional sound source in the horizontal plane, located by its  $(\rho, \theta)$  coordinates, where  $\rho$  is the distance of the source to the head center and  $\theta$  is the azimuth angle. Indeed, as a first approximation in most musical situations, both the listeners and instrumentalists are standing on the (same) ground, with no relative elevation.

The source s will reach the left (L) and right (R) ears through different acoustic paths, characterizable with a pair of filters, which spectral versions are called Head-Related Transfer Functions (HRTFs). HRTFs are frequency- and subject-dependent. The CIPIC database [2] samples different listeners and directions of arrival.

A sound source positioned to the left will reach the left ear sooner than the right one, in the same manner the right level should be lower due to wave propagation and head shadowing. Thus, the difference in amplitude or Interaural Level Difference (ILD, expressed in decibels – dB) [4] and difference in arrival time or Interaural Time Difference (ITD, expressed in seconds) [5] are the main spatial cues for the human auditory system [6].



Fig. 1. Frequency-dependent scaling factors:  $\alpha$  (top) and  $\beta$  (bottom).

#### A. Interaural Level Differences

After Viste [1], the ILDs can be expressed as functions of  $sin(\theta)$ , thus leading to a sinusoidal model:

$$ILD(\theta, f) = \alpha(f)\sin(\theta) \tag{1}$$

where  $\alpha(f)$  is the average scaling factor that best suits our model, in the least-square sense, for each listener of the CIPIC database (see Figure 1). The overall error of this model over the CIPIC database for all subjects, azimuths, and frequencies is of 4.29dB. The average model error and inter-subject variance are depicted in Figure 2.



Fig. 2. Average ILD model error (top) and inter-subject variance (bottom) over the CIPIC database.

Moreover, given the short-time spectra of the left  $(X_L)$  and right  $(X_R)$  channels, we can measure the ILD for each time-frequency bin with:

$$ILD(t,f) = 20\log_{10} \left| \frac{X_L(t,f)}{X_R(t,f)} \right|.$$
(2)

#### B. Interaural Time Differences

Because of the head shadowing, Viste uses for the ITDs a model based on  $\sin(\theta) + \theta$ , after Woodworth [7]. However, from the theory of the diffraction of an harmonic plane wave by a sphere (the head), the ITDs should be proportional to  $\sin(\theta)$ . Contrary to the model by Kuhn [8], our model takes into account the inter-subject variation and the full-frequency band. The ITD model is then expressed as:

$$ITD(\theta, f) = \beta(f)r\sin(\theta)/c \tag{3}$$

where  $\beta$  is the average scaling factor that best suits our model, in the least-square sense, for each listener of the CIPIC database (see Figure 1), r denotes the head radius, and c is the sound celerity. The overall error of this model over the CIPIC database is 0.052ms (thus comparable to the 0.045ms error of the model by Viste). The average model error and inter-subject variance are depicted in Figure 3.

Practically, our model is easily invertible, which is suitable for sound localization, contrary to the  $\sin(\theta) + \theta$  model by Viste which introduced mathematical errors at the extreme azimuths (see [9]).

Given the short-time spectra of the left  $(X_L)$  and right  $(X_R)$  channels, we can measure the ITD for each time-frequency bin with:

$$\operatorname{ITD}_{p}(t,f) = \frac{1}{2\pi f} \left( \angle \frac{X_{L}(t,f)}{X_{R}(t,f)} + 2\pi p \right).$$
(4)

The coefficient p outlooks that the phase is determined up to a modulo  $2\pi$  factor. In fact, the phase becomes ambiguous above 1500Hz, where the wavelength is shorter than the diameter of the head.



Fig. 3. Average ITD model error (top) and inter-subject variance (bottom) over the CIPIC database.

#### C. Distance Cues

The distance estimation or simulation is a complex task due to dependencies on source characteristics and the acoustical environment. Four principal cues are predominant in different situations: intensity, direct-to-reverberant (D/R) energy ratio [10], spectrum, and binaural differences (noticeable for distances less than 1m, see [11]). Their combination is still an open research subject. Here, we focus effectively on the intensity and spectral cues.

In ideal conditions, the intensity of a source is halved (decreases by -6dB) when the distance is doubled, according to the well-known Inverse Square Law [12]. Applying only this frequency-independent rule to a signal has no effect on the sound timbre. But when a source moves far from the listener, the high frequencies are more attenuated than the low frequencies. Thus the sound spectrum changes with the distance. More precisely, the spectral centroid moves towards the low frequencies as the distance increases. In [13], the authors show that the frequency-dependent attenuation due to atmospheric attenuation is roughly proportional to  $f^2$ , similarly to the ISO 9613-1 norm [14]. Here, we manipulate the magnitude spectrum to simulate the distance between the source and the listener (see Section IV). Conversely, we measure the spectral centroid (related to brightness) to estimate the source's distance to listener (see Section V).

#### IV. SPATIALIZATION

#### A. Relative Distance Effect

In a concert room, the distance is often simulated by placing the speaker near / away from the auditorium, which is sometimes physically restricted in small rooms. In fact, the architecture of the room plays an important role and can lead to severe modifications in the interpretation of the piece.

Here, simulating the distance is a matter of changing the magnitude of each short-term spectrum X. More precisely, the ISO 9613-1 norm [14] gives the frequencydependent attenuation factor in dB for given air temperature, humidity, and pressure conditions. At distance  $\rho$ , the magnitudes of X(f) should be attenuated by  $D(f, \rho)$  decibels:

$$D(f,\rho) = \rho \cdot a(f).$$
(5)

where a(f) is the frequency-dependent attenuation, which will have an impact on the brightness of the sound (higher frequencies being more attenuated than lower ones).

More precisely, the total absorption in decibels per meter a(f) is given by a rather complicated formula:

$$\frac{a(f)}{P} \approx 8.68 \cdot F^2 \left\{ 1.84 \cdot 10^{-11} \left( \frac{T}{T_0} \right)^{\frac{1}{2}} P_0 + \left( \frac{T}{T_0} \right)^{-\frac{5}{2}} \right. \\ \left. \left[ 0.01275 \cdot e^{-2239.1/T} / [F_{r,O} + (F^2/F_{r,O})] \right. \\ \left. + 0.1068 \cdot e^{-3352/T} / [F_{r,N} + (F^2/F_{r,N})] \right] \right\} (6)$$

where F = f/P,  $F_{r,O} = f_{r,O}/P$ ,  $F_{r,N} = f_{r,N}/P$  are frequencies scaled by the atmospheric pressure P, and  $P_0$  is the reference atmospheric pressure (1 atm), f is the frequency in Hz, T is the atmospheric temperature in Kelvin (K),  $T_0$  is the reference atmospheric temperature (293.15K),  $f_{r,O}$  is the relaxation frequency of molecular oxygen, and  $f_{r,N}$  is the relaxation frequency of molecular nitrogen. See [13] for details.

#### B. Binaural Spatialization

In binaural listening conditions using headphones, the sound from each earphone speaker is heard only by one ear. Thus the encoded spatial cues are not affected by any cross-talk signals between earphone speakers.

To spatialize a sound source to an expected azimuth  $\theta$ , for each short-term spectrum X, we compute the pair of left  $(X_L)$  and right  $(X_R)$  spectra from the spatial cues corresponding to  $\theta$ , using Equations (1) and (3), and:

$$X_L(t,f) = X(t,f) \cdot 10^{+\Delta_a(f)/2} e^{+j\Delta_\phi(f)/2},$$
(7)

$$X_R(t,f) = X(t,f) \cdot 10^{-\Delta_a(f)/2} e^{-j\Delta_\phi(f)/2}$$
(8)

(because of the symmetry among the left and right ears), where  $\Delta_a$  and  $\Delta_{\phi}$  are given by:

$$\Delta_a(f) = \operatorname{ILD}(\theta, f)/20, \tag{9}$$

$$\Delta_{\phi}(f) = \text{ITD}(\theta, f) \cdot 2\pi f. \tag{10}$$

The control of both amplitude and phase should provide better audio quality [15] than amplitude-only spatialization<sup>1</sup> (see below).

Indeed, we reach a remarkable spatialization realism through informal listening tests with AKG K240 Studio headphones. The main problem which remains is the classic front / back confusion [16].

#### C. Multi-Loudspeaker Spatialization

In a stereophonic display, the sound from each loudspeaker is heard by both ears. Thus, the stereo sound is filtered by a matrix of four transfer functions  $(C_{ij}(f, \theta))$ between loudspeakers and ears (see Figure 4). Here, we generate the paths artificially using the binaural model. The best panning coefficients under CIPIC conditions for the pair of speakers to match the binaural signals at the ears (see Equations (7) and (8)) are then given by:

$$K_L(t,f) = C \cdot (C_{RR}H_L - C_{LR}H_R), \quad (11)$$
  

$$K_R(t,f) = C \cdot (-C_{RL}H_L + C_{LL}H_R) \quad (12)$$

<sup>1</sup>see URL: http://dept-info.labri.fr/~sm/SMC08/

with the determinant computed as:

$$C = 1/(C_{LL}C_{RR} - C_{RL}C_{LR}).$$
 (13)

In extreme cases where |C| = 0 (or close to zero) at any frequency, the matrix C is ill-conditioned, and the solution becomes unstable. To avoid unstable cases, attention should be paid during the loudspeakers configuration stage, before live diffusion.

During diffusion, the left and right signals  $(Y_L, Y_R)$  to feed left and right speakers are obtained by multiplying the short-term spectra X with  $K_L$  and  $K_R$ , respectively:

$$Y_L(t,f) = K_L(t,f) \cdot X(t,f), \qquad (14)$$

$$Y_R(t,f) = K_R(t,f) \cdot X(t,f).$$
(15)

In a setup with many speakers we use the classic pairwise paradigm [17], consisting in choosing for a given source only the two speakers closest to it (in azimuth): one at the left of the source, the other at its right.



Fig. 4. Stereophonic loudspeaker display.

#### D. Analysis of Panning Coefficients

We used the speaker pair  $(-30^{\circ}, +30^{\circ})$  to compute the panning coefficients at any position (between the speakers) with the two techniques: our approach and the classic vector-based amplitude panning (VBAP) approach [3]. VBAP was elaborated under the assumption that the incoming sound is different only in amplitude, which holds for frequencies up to 600Hz. In fact, by controlling correctly the amplitudes of the two channels, it is possible to produce resultant phase and amplitude differences for continuous sounds that are very close to those experienced with natural sources [16]. We restrict our comparisons to the [0, 800]Hz frequency band.

1) Comparisons of Panning Coefficients: The panning coefficients of the two approaches are very similar until 600Hz (see Figure 5), and can differ significantly above. In fact, our coefficients are complex values, and their imaginary parts can contribute in a significant way (see Figure 6).



Fig. 5. Amplitude of the panning coefficients from VBAP (plain) and our approach (dotted), for the left (top) and right (bottom) channels of the panning pair for  $-15^{\circ}$ , in the [0, 800]Hz band.



Fig. 6. Phase of the panning coefficients from our approach, for the left (dotted) and right (plain) channels of the panning pair for  $-15^{\circ}$ , in the [0, 800]Hz band.

2) Comparisons of the Ratio of Panning Coefficients: Generally, inter-channel differences are perceptually more relevant (*e.g.* ILD, ITD) than absolute values.

Given the left and right panning coefficients,  $K_L$  and  $K_R$ , we compute the *panning level difference* (PLD):

$$PLD = 20 \log_{10} \left| \frac{K_L}{K_R} \right|.$$
 (16)

We computed the absolute difference between the PLDs of both VBAP and our approach. The maximal PLD difference (in the considered frequency band) has a linear trend, and its maximum does not exceed 3dB. Thus, the two approaches seem to be consistent in the [0, 800]Hz band (see Figure 7). For higher frequencies, the new approach should yield better results, as confirmed perceptively in our preliminary and informal listening tests.



Fig. 7. Maximum difference per azimuth between PLDs of VBAP and the proposed method in the [0, 800]Hz band.

#### V. LOCALIZATION

#### A. Azimuth Estimation

In Auditory Scene Analysis (ASA), ILDs and ITDs are the most important cues for source localization. Lord Rayleigh mentioned in his Duplex Theory [18] that the ILDs are more prominent at high frequencies (where phase ambiguities are likely to occur) whereas the ITDs are crucial at low frequencies (which are less attenuated during their propagation).

Obtaining an estimation of the azimuth based on the ILD information (see Equation (2)) is just a matter of inverting Equation (1):

$$\theta_L(t, f) = \arcsin\left(\frac{\operatorname{ILD}(t, f)}{\alpha(f)}\right).$$
(17)

Similarly, using the ITD information (see Equation (4)), to obtain an estimation of the azimuth candidate for each p, we invert Equation (3):

$$\theta_{T,p}(t,f) = \arcsin\left(\frac{c \cdot \text{ITD}_p(t,f)}{r \cdot \beta(f)}\right).$$
(18)

The  $\theta_L(t, f)$  estimates are more dispersed, but not ambiguous at any frequency, so they are exploited to find the right modulo coefficient p that unwraps the phase. Then the  $\theta_{T,p}(t, f)$  that is nearest to  $\theta_L(t, f)$  is validated as the final  $\theta$  estimation for the considered frequency bin, since it exhibits a smaller deviation:

$$\theta(t,f) = \theta_{T,m}(t,f),\tag{19}$$

with  $m = \operatorname{argmin}_{p} |\theta_{L}(t, f) - \theta_{T,p}(t, f)|$ .

Practically, the choice of  $\hat{p}$  can be limited among two values  $(\lceil p_r \rceil, |p_r|)$ , where

$$p_r = \left( f \cdot \text{ITD}(\theta_L, f) - \frac{1}{2\pi} \angle \frac{X_L(t, f)}{X_R(t, f)} \right).$$
(20)

An estimate of the azimuth of the source can be obtained as the peak in an energy-weighted histogram (see [9]). More precisely, for each frequency bin of each discrete spectrum, an azimuth is estimated and the power corresponding to this bin is accumulated in the histogram at this azimuth. For the corresponding bin frequency f, the power  $|X(f)|^2$  is estimated by inverting Equations (7) and (8) for the left and right spectra, respectively, then the square of the estimate of the loudest – supposedly most reliable – channel is retained for the power estimate.

Thus, we obtain a power histogram as shown in Figure 8. This histogram is the result of the localization of a Gaussian white noise of 0.5s spatialized at azimuth  $-45^{\circ}$ . On this figure, we can clearly see two important local maxima (peaks), one around azimuth  $-45^{\circ}$ , the other at azimuth  $-90^{\circ}$ . The first (and largest) one corresponds to the sound source; the second one is a spurious peak resulting from extreme ILDs (a problem we have to solve in our future research).



Fig. 8. Histogram obtained with a source at azimuth  $-45^{\circ}$ . One can clearly see two important local maxima (peaks): one around azimuth  $-45^{\circ}$ , the other at azimuth  $-90^{\circ}$ . The first (and largest) one corresponds to the sound source; the second one is a spurious peak resulting from extreme ILDs.



Fig. 9. Histogram obtained with a real source positioned at azimuth  $30^{\circ}$  in a real room, with binaural signals recorded at the ears of the musician.

For our localization tests, we spatialized a Gaussian white noise using convolutions with the HRTFs of the KEMAR manikin (see [2]), since they were not part of the database used for the learning of our model coefficients and thus should give results closer to those expected with a real – human – listener. Indeed, in our first experiments with real listeners (see Figure 9), the same trends as in Figure 8 were observed: a rather broader histogram but

still with a local maximum close to the azimuth of the sound source, plus spurious maxima at extreme azimuths  $\pm 90^{\circ}$ .

To verify the precision of the estimation of the azimuth, we spatialized several noise sources at different azimuths in the horizontal plane, between  $-80^{\circ}$  and  $+80^{\circ}$ , and we localized them using the proposed method. The results are shown in Figure 10. We observe that the absolute azimuth error is less than  $5^{\circ}$  in the  $[-65, +65]^{\circ}$  range.



Fig. 10. Absolute error of the localization of the azimuth from Gaussian white noise spatialized at different azimuths using convolutions with the HRTFs of the KEMAR manikin.

In real reverberant environments, due to more superpositions at the microphones, an amplitude-based method is not really adapted; in contrast, generalized crosscorrelation based ITD estimation should be more robust [19].

#### B. Distance Estimation

As a reference signal for distance estimation, we use a Gaussian white noise spatialized at azimuth zero, since pure tones are not suitable for distance judgments [20]. The distance estimation relies on the quantification of the spectral changes during the sound propagation in the air.

To estimate the amplitude spectrum, we first estimate the power spectral density of the noise using the Welch's method [21], [22]. More precisely, we compute the mean power of the short-term spectra over L frames, then take its square root, thus:

$$|X| = \sqrt{\frac{1}{L} \sum_{l=-(L-1)/2}^{l=+(L-1)/2} |X_l|^2}.$$
 (21)

In our experiments, we consider L = 21 frames of N = 2048 samples, with an overlap factor of 50% (and with a CD-quality sampling rate of 44.1kHz, thus the corresponding sound segment has a length < 0.5s).

Then we use this amplitude spectrum to compute the spectral centroid:

$$C = \frac{\sum_{f} f \cdot |X(f)|}{\sum_{f} |X(f)|}.$$
(22)

The spectral centroid moves towards low frequencies when the source moves away from the observer. The related perceptive brightness is an important distance cue. We know the reference distance since the CIPIC speakers were positioned on a 1-m radius hoop around the listener. By inverting the logarithm of the function of Figure 11, obtained thanks to the ISO 9613-1 norm and Equations (5), (6), and (22), we can propose a function to estimate the distance from a given spectral centroid:

$$\rho(\log(\mathcal{C})) = -38.89044\mathcal{C}^3 + 1070.33889\mathcal{C}^2 - 9898.69339\mathcal{C}^+ 30766.67908$$
(23)

given for the air at  $20^{\circ}$  Celsius temperature, 50% relative humidity, and 1 atm pressure.



Fig. 11. Spectral centroid (related to perceptive brightness) as a function of distance at 20° Celsius temperature, 50% relative humidity, and 1 atm atmospheric pressure (for white noise played at CD quality).

Up to 25m, the maximum distance error is theoretically less than 4mm, if the noise power spectral density is known. However, if the amplitude spectrum has to be estimated using Equation (21), then the error is greater, though very reasonable until 50m. Figure 12 shows the results of our simulations for Gaussian white noise spatialized at different distances in the [0, 100]m range.



Fig. 12. Absolute error of the localization of the distance from Gaussian white noise spatialized at different distances.

#### VI. RETROSPAT SOFTWARE SYSTEM

The RetroSpat system is being implemented as a realtime musical software under the GNU General Public License (GPL). The actual implementation is based on C++, Qt4<sup>2</sup>, JACK<sup>3</sup>, FFTW<sup>4</sup> and works on Linux and MacOS X.

Currently, RetroSpat implements the described methods (*i.e.* localization and spatialization) in two different modules: *RetroSpat Localizer* for speaker setup detection and *RetroSpat Spatializer* for the spatialization process. We hope to merge the two functionalities in one unique software soon.

#### A. RetroSpat Localizer

RetroSpat Localizer (see Figure 13) is in charge of the automatic detection of the speakers configuration. It also allows the user to interactively edit a configuration, which has been just detected or loaded from an XML file.

The automatic detection of the positions (azimuth and distance) of the speakers connected to the soundcard is of great importance to adapt to new speaker setups. Indeed, it will be one of the first actions of the interpreter in a new environment.

For room calibration, the interpreter carries headphones with miniature microphones encased in earpieces (see Figure 15, where Sennheiser KE4-211-2 microphones have been inserted in standard headphones). The interpreter orients the head towards the desired zero azimuth. Then, each speaker plays in turn a Gaussian white noise sampled at 44.1kHz. The binaural signals recorded from the ears of the musician are transferred to the computer running RetroSpat Localizer. Each speaker is then localized in azimuth and distance. The suggested configuration can be adjusted or modified by the interpreter according to the rooms characteristics.

#### B. RetroSpat Spatializer

For sound spatialization, mono sources are loaded in RetroSpat, parameterized, and then diffused. The settings include the volume of each source, the initial localization, the choice of special trajectories such as circle, arc, etc. A loudspeaker-array configuration is the basic element for the spatialization (see Section VI-A).

The snapshot of Figure 14 depicts a 7-source mix of instruments and voices (note icons), in a 6-speaker front-facing configuration (loudspeaker icons), obtained from RetroSpat Localizer.

During the diffusion, the musician can interact individually with each source of the piece, change its parameters (azimuth and distance), or even remove / insert a source from / into the scene. In this early version, the interaction with RetroSpat is provided by a mouse controller.

Thanks to an efficient implementation using the JACK sound server, RetroSpat can diffuse properly simultaneous sources even within the same speaker pair (see Figure 14, three sources in speaker pair (2,3)). All the speaker pairs have to stay in synchrony. To avoid sound perturbation, the Qt-based user interface runs in a separated thread with less priority than the core signal processing process. We tested RetroSpat on a MacBook Pro, connected to 8 speakers, through a MOTU 828 MKLL soundcard, and were able to play several sources without problems.

However, further testing is needed to assess scalability limits.



Fig. 15. The "phonocasque" used for the binaural recordings: standard headphones where microphone capsules have been inserted.

#### C. Musical Applications

In a live concert, the acousmatic musician interacts with the scene through a special (un)mixing console.

With RetroSpat, the musician has more free parameters on one single controller (mouse):

- only mouse movement to control simultaneously the azimuthal and distance location;
- mono sources as inputs;
- many sources can be spatialized to different locations at the same time;
- a dynamic visualization of the whole scene (source apparition, movement, speed, etc.) is provided.

We believe that RetroSpat should greatly simplify the interpreter interactions and thus should allow him / her to focus more on the artistic performance.

#### VII. CONCLUSIONS AND FUTURE WORK

In this paper, we have introduced a flexible multisource, multi-loudspeaker system: RetroSpat. This realtime system implements our proposed binaural to multiloudspeaker spatialization method. The system can also locate the loudspeakers azimuths and distances.

Several experiments at the SCRIME studio on an octophonic setup justify the utility of the system for live performance by composers of electroacoustic music. Next, we should enhance the source localization in real – reverberant – environments, and possibly evolve to source control through gesture or a more intuitive hardware controller. Also, a major scientific challenge would be to separate the different sources present in a binaural mix (for a semi-automatic diffusion from a compact disc as support).

#### VIII. ACKNOWLEDGMENTS

This research was carried out in the context of the SCRIME (Studio de Création et de Recherche en Informatique et Musique Électroacoustique) and was supported by the Conseil Régional d'Aquitaine, the Ministère de la Culture, the Direction Régionale des Actions Culturelles d'Aquitaine, and the Conseil Général de la Gironde.

<sup>&</sup>lt;sup>2</sup>see URL: http://trolltech.com/products/qt

<sup>&</sup>lt;sup>3</sup>see URL: http://jackaudio.org

<sup>&</sup>lt;sup>4</sup>see URL: http://www.fftw.org



Fig. 13. RetroSpat Localizer graphical user interface with a 6-speaker configuration.



Fig. 14. RetroSpat Spatializer graphical user interface, with 7 sources spatialized on the speaker setup presented on Figure 13.

#### REFERENCES

- [1] H. Viste, "Binaural Localization and Separation Techniques," Ph.D. dissertation, École Polytechnique Fédérale de Lausanne, Switzerland, 2004.
- [2] V. R. Algazi, R. O. Duda, D. M. Thompson, and C. Avendano, "The CIPIC HRTF Database," in *Proceedings of the IEEE Work*shop on Applications of Signal Processing to Audio and Acoustics
- (WASPAA), New Paltz, New York, 2001, pp. 99–102. V. Pulkki, "Virtual Sound Source Positioning using Vector Base Amplitude Panning," *Journal of the Acoustical Society of America*, [3] vol. 45, no. 6, pp. 456-466, 1997.
- J. W. Strutt (Lord Rayleigh), "On the Acoustic Shadow of a [4] Sphere," Philosophical Transactions of the Royal Society of London, vol. 203A, pp. 87-97, 1904.
- "Acoustical Observations I," Philosophical Magazine, vol. 3, [5] pp. 456-457, 1877.
- [6] J. Blauert, Spatial Hearing, revised ed. Cambridge, Massachusetts: MIT Press, 1997, translation by J. S. Allen.
- [7] R. S. Woodworth, Experimental Psychology. New York: Holt, 1954.
- [8] G. F. Kuhn, "Model for the Interaural Time Differences in the Azimuthal Plane," Journal of the Acoustical Society of America, vol. 62, no. 1, pp. 157–167, 1977.
- [9] J. Mouba and S. Marchand, "A Source Localization / Separation / Respatialization System Based on Unsupervised Classification of Interaural Cues," in *Proceedings of the Digital Audio Effects* (DAFx) Conference, Montreal, 2006, pp. 233-238.
- [10] A. W. Bronkhorst and T. Houtgast, "Auditory Distance Perception in Rooms," *Nature*, vol. 397, pp. 517–520, 1999.
  [11] D. Brungart and W. Rabinowitz, "Auditory Localization of Nearby
- Sources," Journal of the Acoustical Society of America, vol. 106, pp. 1465–1479, 1999.

- [12] R. E. Berg and D. G. Stork, The Physics of Sound, 2nd ed. Prentice Hall, 1994.
- [13] H. Bass, L. Sutherland, A. Zuckerwar, D. Blackstock, and D. Hester, "Atmospheric Absorption of Sound: Further Developments," Journal of the Acoustical Society of America, vol. 97, no. 1, pp. 680-683, 1995.
- [14] ISO 9613-1:1993: Acoustics Attenuation of Sound During Propagation Outdoors - Part 1: Calculation of the Absorption of Sound by the Atmosphere, International Organization for Standardization, Geneva, Switzerland, 1993.
- C. Tournery and C. Faller, "Improved Time Delay Analy-sis/Synthesis for Parametric Stereo Audio Coding," *Journal of the Audio Engineering Society*, vol. 29, no. 5, pp. 490–498, 2006.
- [16] F. Rumsey, Spatial Audio, 1st ed. Oxford, United Kingdom: Focal Press, 2001, reprinted 2003, 2005. [17] J. M. Chowning, "The Simulation of Moving Sound Sources,"
- Journal of the Acoustical Society of America, vol. 19, no. 1, pp. 2-6, 1971.
- [18] J.
- J. W. Strutt (Lord Rayleigh), "On Our Perception of Sound Direction," *Philosophical Magazine*, vol. 13, pp. 214–302, 1907. C. H. Knapp and G. C. Carter, "The Generalized Correlation Method for the Estimation of Time Delay," *IEEE Transactions* [19] on Signal Processing, vol. 24, no. 4, pp. 320-327, 1976.
- [20] J. Molino, "Perceiving the Range of a Sound Source When the Direction is Known," *Journal of the Acoustical Society of America*, vol. 53, no. 5, pp. 1301–1304, 1973.
- [21] P. D. Welch, "The Use of Fast Fourier Transform for the Estimation of Power Spectra: A Method Based on Time-Averaging over Short, Modified Periodograms," *IEEE Transactions on Audio and Electroacoustics*, vol. 15, no. 22, pp. 70–73, 1967.
- [22] M. H. Hayes, Statistical Digital Signal Processing and Modeling. New Jersey: John Wiley & Sons, 1996.