# Effcient Multichannel Audio Resynthesis by Subband-Based Spectral Conversion

Athanasios Mouchtaris, Shrikanth S. Narayanan, and Chris Kyriakakis *
mouchtar@sipi.usc.edu, shri@sipi.usc.edu, ckyriak@imsc.usc.edu

## ABSTRACT

Multichannel audio offers significant advantages for music reproduction that include the ability to provide better localization and envelopment, as well as reduced imaging distortion. On the other hand, multichannel audio is one of the most demanding media types in terms of transmission requirements. A novel architecture was previously proposed, allowing delivery of uncompressed multichannel audio over high-bandwidth communications networks. In most cases, however, bandwidth limitations prohibit transmission of multiple audio channels. In such cases, an alternative would be to transmit only one or two reference channels and recreate the rest of the channels at the receiving end. In this paper, we propose a system that is capable of synthesizing the required signals from a smaller set of signals recorded in a particular venue. These synthesized "virtual" microphone signals can be used to produce multichannel recordings that accurately capture the acoustics of the particular venue. Applications of the proposed system include transmission of multichannel audio over the current Internet infrastructure and, as an extension of the methods proposed here, remastering of existing monophonic and stereophonic recordings for multichannel rendering.

## 1 Introduction

Multichannel audio is one of the most demanding media types in terms of transmission requirements. A novel architecture allowing delivery of uncompressed multichannel audio over high-bandwidth communications networks has been presented in [1]. As suggested there, for applications in which bandwidth limitations prohibit transmission of multiple audio channels, an alternative is to transmit only one or two channels (reference channels, *e.g.* the left and right signals in a traditional stereo recording) and reconstruct the remaining channels at the receiving end. The system proposed there partially provided a solution for reconstructing the channels of a specific recording from the reference channels and was particularly suitable for live concert hall performances. The proposed algorithm was based on information of the acoustics of the specific concert hall and the microphone locations with respect to the orchestra, information that was extracted from the specific multichannel recording.

In this paper the methods for recreating the channels of a multichannel recording proposed in [1] are extended.

Before proceeding to the description of the algorithm proposed here, a brief outline of the previously published analysis is given. The reader is asked to examine Fig. 1, where an example is given of how microphones may be arranged in a recording venue in a multichannel recording. A number of microphones are used to capture several characteristics of the venue, which are then mixed and played back through a multichannel audio system that recreates the spatial realism of the recording venue. By examining the acoustical characteristics of the various microphone recordings, the distinction of microphones is made into reverberant and spot microphones.

Reverberant microphones are the microphones placed far from the sound source, for example C and D in Fig. 1. These microphones are treated separately as one category because they mainly capture reverberant information (that can be reproduced by the surround channels in a multichannel playback system). The recordings captured by these microphones can be synthesized by passing the reference recordings through a linear time-invariant (LTI) filter, designed as in [1]. The challenge of this design is that, for venues such as large concert halls, these filters are of the order of several thousand taps and serious implementation issues arise.

Spot microphones are the microphones that are placed close to the sound source (*e.g.* G in Fig. 1). These microphones introduce an even more challenging situation. Because the source of sound is not a point source but rather distributed such as in an orchestra, the recordings of these microphones depend largely on the instruments that are near the microphone and not so much on the acoustics of the hall. Synthesizing the recordings of these microphones, therefore, involves enhancing certain instruments and diminishing others, which in most cases overlap both in the time and frequency domains. The algorithm described here focuses on this problem and is based on spectral conversion (SC).

## 2 Spectral Conversion

Based on the analysis given in the previous paragraph, the goal is to modify the short-term spectral properties of the reference audio signal in order to recreate the desired one. The short-term spectral properties are extracted by using a short sliding window with overlapping. Each frame is modeled as an autoregressive filter excited by a residual signal. The AR filter coefficients are found by means of linear predictive analysis (LPC, [2]) and the residual signal is the result of inverse filtering the audio signal of the current frame by the AR filter. The LPC coefficients are modified in a way to be described later in this section and
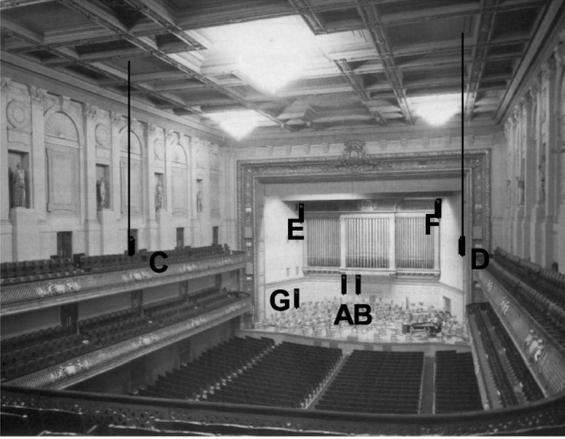
Figure 1: An example of how microphones may be arranged in a recording venue for a multichannel recording.

the residual is filtered with the designed AR filter to produce the desired signal of the current frame. Finally, the desired response is synthesized from the designed frames using overlap-add techniques.

In order to obtain the desired response for each frame, an algorithm is required for converting the LPC coefficients into the desired ones. Although the target coefficients in the application examined can be found by applying the same residual/LP analysis described (assuming that the reference and target waveforms are time-aligned), our intention is to design a mapping function based on the reference and target responses whose parameters will remain constant. The result will be a significant reduction of information as the target response can be reconstructed using the reference signal and this function.

Such a mapping function can be designed by following the approach of voice conversion algorithms [3, 4, 5]. The objective of voice conversion is to modify a speech waveform so that the context remains as is but appears to be spoken from a specific speaker. Although the application is completely different, the approach followed is very suitable for our system. In voice conversion pitch and time-scaling need to be considered, while in the application examined here this is not necessary since the reference and target waveforms come from the same excitation recorded with different microphones and the need is not to modify but to *enhance* the reference waveform. However, in both cases, there is the need to modify the short-term spectral properties of the waveform. The method to do that is briefly described next.

Assuming that a sequence $[\boldsymbol{x}_1\boldsymbol{x}_2\cdots\boldsymbol{x}_n]$ of reference spectral vectors (*e.g.* line spectral frequencies (LSF's), cepstral coefficients, *etc.* ) is given, as well as the corresponding sequence of target spectral vectors $[\boldsymbol{y}_1\boldsymbol{y}_2\cdots\boldsymbol{y}_n]$, a function $\mathcal{F}(\cdot)$ can be designed which, when applied to vector $\boldsymbol{x}_i$, produces a vector close in some sense to vector $\boldsymbol{y}_i$. Many algorithms have been described for designing this function (see [3, 4, 5, 6] and the references therein). Here the algorithms based on vector quantization (VQ, [3]) and Gaussian mixture models (GMM, [4, 5]) were implemented and compared.

## 2.1   Spectral Conversion based on VQ

Under this approach, the spectral vectors of the reference and target signals (training data) are vector quantized re-

spectively using the well-known modified K-means clustering algorithm (see for example [7] for details). Then, a histogram is created indicating the correspondences between the reference and target centroids. Finally, the function $\mathcal{F}$ is defined as the linear combination of the target centroids using the designed histogram as a weighting function. It is important to mention that in this case the spectral vectors were chosen to be the cepstral coefficients so that the distance measure used in clustering is the truncated cepstral distance.

## 2.2   Spectral Conversion based on GMM

In this case, the assumption made is that the sequence of spectral vectors $x$ has a probability density function (pdf) that can be modeled as a mixture of $m$ multivariate Gaussian pdf's. Thus, the pdf of $x$, p($x$), can be written as

$$\mathrm{p}(x) = \sum_{i=1}^{m} \alpha_i \mathrm{N}(x, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i) \qquad (1)$$

where, $\mathrm{N}(x, \boldsymbol{\mu}, \boldsymbol{\Sigma})$ is the normal multivariate distribution with mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\Sigma}$ and the weights $\alpha_i$ are non-negative under the constraint that $\sum_{i=1}^{m} \alpha_i = 1$ ($\alpha_i$ is the prior probability of class $\mathcal{C}_i$). The parameters of the GMM, *i.e.* the mean vectors, covariance matrices and priors, can be estimated using the expectation maximization (EM) algorithm [8].

As already mentioned, the function $\mathcal{F}$ is designed such that the spectral vectors $\boldsymbol{y}_i$ and $\mathcal{F}(\boldsymbol{x}_i)$ are close in some sense. In [4], the function $\mathcal{F}$ is designed such that the error

$$\mathcal{E} = \sum_{k=1}^{n} \|\boldsymbol{y}_k - \mathcal{F}(\boldsymbol{x}_k)\|^2 \qquad (2)$$

is minimized. Since this method is based on least-squares estimation, it will be denoted as the LSE method. This problem becomes possible to solve under the constraint that $\mathcal{F}$ is piecewise linear, *i.e.*

$$\mathcal{F}(\boldsymbol{x}_k) = \sum_{i=1}^{m} \mathrm{P}(\mathcal{C}_i|\boldsymbol{x}_k) \left[ \boldsymbol{v}_i + \boldsymbol{\Gamma}_i \boldsymbol{\Sigma}_i^{-1} (\boldsymbol{x}_k - \boldsymbol{\mu}_i) \right] \qquad (3)$$

where the conditional probability that a given vector $\boldsymbol{x}_k$ belongs to class $\mathcal{C}_i$, $\mathrm{P}(\mathcal{C}_i|\boldsymbol{x}_k)$ can be computed by applying Bayes' theorem

$$\mathrm{P}(\mathcal{C}_i|\boldsymbol{x}_k) = \frac{\alpha_i \mathrm{N}(\boldsymbol{x}_k, \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)}{\sum_{j=1}^{m} \alpha_j \mathrm{N}(\boldsymbol{x}_k, \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)} \qquad (4)$$

The unknown parameters ($\boldsymbol{v}_i$ and $\boldsymbol{\Gamma}_i$, $i = 1, \cdots, m$) can be found by minimizing (2) which reduces to solving a typical least-squares equation.

A different solution for function $\mathcal{F}$ is given when a different function compared to (2) is minimized [5]. In the mean-squared sense, the optimal choice for the function $\mathcal{F}$ is

$$\begin{aligned} \mathcal{F}(\boldsymbol{x}_k) &= \mathrm{E}(y|\boldsymbol{x}_k) \qquad (5) \\ &= \sum_{i=1}^{m} \mathrm{P}(\mathcal{C}_i|\boldsymbol{x}_k) \left[ \boldsymbol{\mu}_i^y + \boldsymbol{\Sigma}_i^{yx} \boldsymbol{\Sigma}_i^{xx-1} (\boldsymbol{x}_k - \boldsymbol{\mu}_i^x) \right] \end{aligned}$$

where $\mathrm{E}(\cdot)$ denotes the expectation operator and now

$$\mathrm{P}(\mathcal{C}_i|\boldsymbol{x}_k) = \frac{\alpha_i \mathrm{N}(\boldsymbol{x}_k, \boldsymbol{\mu}_i^x, \boldsymbol{\Sigma}_i^{xx})}{\sum_{j=1}^{m} \alpha_j \mathrm{N}(\boldsymbol{x}_k, \boldsymbol{\mu}_j^x, \boldsymbol{\Sigma}_j^{xx})} \qquad (6)$$

If the source and target vectors are concatenated, creating a new sequence $z = [x^T y^T]^T$, (where $^T$ denotes transposition) then all the required parameters in the above equations can be found by estimating the GMM parameters of $z$. Then,

$$\boldsymbol{\Sigma}_i = \left[ \begin{array}{cc} \boldsymbol{\Sigma}_i^{xx} & \boldsymbol{\Sigma}_i^{xy} \\ \boldsymbol{\Sigma}_i^{yx} & \boldsymbol{\Sigma}_i^{yy} \end{array} \right], \boldsymbol{\mu}_i = \left[ \begin{array}{c} \boldsymbol{\mu}_i^x \\ \boldsymbol{\mu}_i^y \end{array} \right] \qquad (7)$$

Once again, these parameters are estimated by the EM algorithm. Since this method estimates the desired function based on the joint density of $x$ and $y$, it will be referred to as the JDE method.

## 2.3 Subband Processing

Audio signals contain more information than speech signals. The sampling rate for audio signals is usually 44.1 or 48 kHz compared to 16 kHz for speech. Moreover, high acoustical quality for audio is essential. For these reasons, the decision to follow an analysis in sub-bands seems natural. Instead of warping the frequency spectrum using the Bark scale as is usual in speech analysis, the frequency spectrum was divided in sub-bands and each one was treated separately under the analysis of the previous paragraphs. Perfect reconstruction filter banks, based on wavelets [9], provide a solution with acceptable computational complexity as well as the appropriate, for audio signals, octave frequency division. The choice of filter bank was not a subject of investigation but steep transition is a desirable property.

## 2.4 Residual Processing for Percussive Sounds

The SC methods described earlier will not produce the desired result in all cases. One such case of particular importance is the case of percussive drum-like sounds. It is usual in multichannel recordings to place a microphone close to the tympani as drum-like sounds are considered perceptually important in recreating the acoustical environment of the recording venue. For percussive sounds, a similar model to the residual/LP model described here can be used [10], but for the enhancement purposes investigated in this paper, the emphasis is given to the residual instead of the LP parameters. It is proposed to extract the residual of an instance of the particular percussive instrument from the recording of the microphone that captures this instrument and then recreate this channel from the reference channel by simply substituting the residual of all instances of this instrument with the extracted residual. This residual corresponds to the interaction between the exciter and the resonating body of the instrument and lasts until the structure reaches a steady vibration so it is independent of the frequencies and decays of the harmonics of the instrument at a given time (after the instrument has reached a steady vibration) and it can be used for synthesizing different sounds by using an appropriate LP filter. This method was successfully tested and more details are given in the next section. The drawback of this approach is that a robust algorithm is required for identifying the particular instrument instances in the reference recording. A possible improvement of this method would be to extract all instances of the instrument from the target response and use some clustering technique for choosing the residual that is more appropriate in the resynthesis stage. The reason is that the residual/LP model introduces modeling error which is larger in the spectral valleys of the AR spectrum, thus better results would be obtained by using a residual which corresponds to an AR

| Band Nr. | Frequency Range | | LPC Order | GMM Centroids |
|---|---|---|---|---|
| | Low (kHz) | High (kHz) | | |
| 1 | 0.0000 | 0.1723 | 4 | 4 |
| 2 | 0.1723 | 0.3446 | 4 | 4 |
| 3 | 0.3446 | 0.6891 | 8 | 8 |
| 4 | 0.6891 | 1.3782 | 16 | 16 |
| 5 | 1.3782 | 2.7563 | 32 | 16 |
| 6 | 2.7563 | 5.5125 | 32 | 16 |
| 7 | 5.5125 | 11.0250 | 32 | 16 |
| 8 | 11.0250 | 22.0500 | 32 | 16 |

Table 1: Parameters for chorus microphone example.

filter as close as possible to the resynthesis AR filter. However, this approach would again require identifying all the instances of the instrument.

## 3 Implementation Details

The three spectral conversion methods outlined in Section 2 were implemented and tested using a multichannel recording, as described in the first section of this paper. The objective was to recreate the channel that mainly captured the chorus of the orchestra (residual processing for percussive sound resynthesis is also considered at the last paragraph of this section). Acoustically, therefore, the emphasis was on the male and female voices. At the same time, it was clear that some instruments, inaudible in the target recording but particularly audible in the reference recording, needed to be attenuated. A database of about 10,000 spectral vectors for each band was carefully created (so that only parts of the recording where the chorus is present are used) with the choice of spectral vectors being the cepstral coefficients. Parts of the chorus recording were selected so that there were no segments of silence included. Results were evaluated through informal listening tests and through objective performance criteria. The SC methods were shown to provide promising enhancement results. Formal listening tests are currently underway and will be available in the near future. The experimental conditions are given in Table 1. The number of octave bands used was 8, a choice that gives particular emphasis on the frequency band 0-5 kHz and at the same time does not impose excessive computational demands. The frequency range 0-5 kHz is particularly important for the specific case of chorus recording resynthesis since this is the frequency range where the human voice is mostly concentrated. For even better results the entire frequency range 0-20 kHz must be considered. The order of the LPC filter varied depending on the frequency detail of each band and for the same reason the number of centroids for each band was different.

In Table 2, the average quadratic cepstral distance (averaged over all vectors and all 8 bands) is given for each method, for the training data as well as for the data used for testing (9 sec. of music from the same recording). The cepstral distance is normalized with the average quadratic distance between the reference and the target waveforms (*i.e.* without any conversion of the LPC parameters). The improvement is large for both the GMM-based algorithms, with the LSE algorithm being slightly better, for both the training and testing data. The VQ-based algorithm, in contrast, produced a deterioration in performance which was audible as well. This can be explained based on the fact that the GMM-based methods result in a conversion
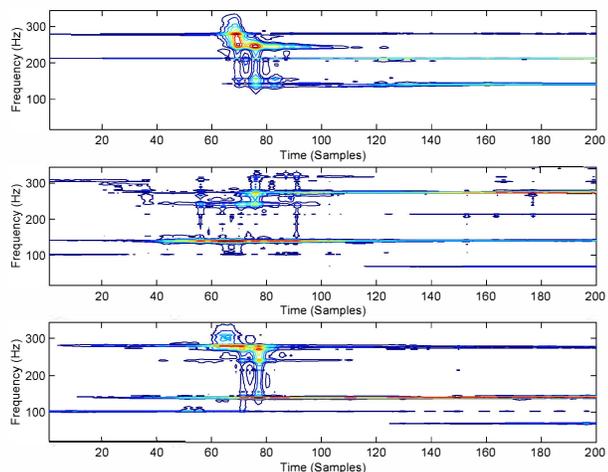
Figure 2: Choi-Williams distribution of the desired (top), reference (middle) and synthesized (bottom) waveforms at the time points during a tympani strike (samples 60-80).

| SC Method | Cepstral Distance | | Centroids per Band |
| --- | --- | --- | --- |
| | Train | Test | |
| LSE | 0.6451 | 0.7144 | Table 1 |
| JDE | 0.6629 | 0.7445 | Table 1 |
| VQ | 1.2903 | 1.3338 | 1024 |

Table 2: Normalized distances for LSE-, JDE- and VQ-based methods.

function which is continuous with respect to the spectral vectors. The VQ-based method, on the other hand, produces audible artifacts introduced by spectral discontinuities because the conversion is based on a limited number of existing spectral vectors. This is the reason why a large number of centroids was used for the VQ-based algorithm as seen in Table 2 compared to the number of centroids used for the GMM-based algorithms. However, the results were still unacceptable both from the objective and subjective perspectives.

The algorithm described in Section 2 considering the special case of percussive sound resynthesis was tested as well. Fig. 2 shows the time-frequency evolution of a tympani instance using the Choi-Williams distribution [11], a distribution that achieves the high resolution needed in such cases of impulsive nature. Fig. 2 clearly demonstrates the improvement in drum-like sound resynthesis. The impulsiveness of the signal at around samples 60-80 is observed in the desired response and verified in the synthesized waveform. The attack part is clearly enhanced, significantly adding naturalness in the audio signal, as our informal listening tests clearly demonstrated. Formal listening tests demonstrating the perceptual benefits of this method will be described in a future publication.

## 4 Conclusions

Multichannel audio resynthesis is a new and important application that allows transmission of only one or two channels of multichannel audio and resynthesis of the remaining channels at the receiving end. Spectral conversion algorithms that have been used successfully for voice conversion can be adopted for the task of multichannel audio resynthesis quite favorably. Three of the most common spectral conversion methods have been compared and our objective results, in accordance with our informal listening tests, have indicated that GMM-based spectral conversion can produce extremely successful results. Residual signal enhancement was also found to be essential for the special case of percussive sound resynthesis. Our current research is focused on audio quality improvement for the proposed methods, conducting formal listening tests as well as extensions of this research for the purpose of remastering existing monophonic and stereophonic recordings for multichannel rendering.

## References

[1] A. Mouchtaris, Z. Zhu, and C. Kyriakakis, "High-quality multichannel audio over the internet," in *Conf. Record of the Thirty-Third Assilomar Conf. Signals, Systems and Computers*, Pacific Grove, CA, October 1999, vol. 1, pp. 347–351.

[2] S. Haykin, *Adaptive Filter Theory*, Prentice Hall, 1996.

[3] M. Abe, S. Nakamura, K. Shikano, and H. Kuwabara, "Voice conversion through vector quantization," in *IEEE Proc. Int. Conf. Acoustics, Speech and Signal Processing (ICASSP)*, New York, NY, April 1988, pp. 655–658.

[4] Y. Stylianou, O. Cappe, and E. Moulines, "Statistical methods for voice quality transformation," in *Proceedings of EUROSPEECH*, September 1995, pp. 1029–1032.

[5] A. Kain and M. W. Macon, "Spectral voice conversion for text-to-speech synthesis," in *IEEE Proc. Int. Conf. Acoustics, Speech and Signal Processing (ICASSP)*, Seattle, WA, May 1998, pp. 285–289.

[6] G. Baudoin and Y. Stylianou, "On the transformation of the speech spectrum for voice conversion," in *IEEE Proc. Int. Conf. Spoken Language Processing (ICSLP)*, Philadephia, PA, October 1996, pp. 1405–1408.

[7] L. Rabiner and B.-H. Juang, *Fundamentals of Speech Recognition*, Prentice Hall, Englewood Cliffs, NJ, 1993.

[8] D. A. Reynolds and R. C. Rose, "Robust text-independent speaker identification using Gaussian mixture speaker models," *IEEE Trans. Speech and Audio Processing*, vol. 3, no. 1, pp. 72–83, January 1995.

[9] G. Strang and T. Nguyen, *Wavelets and Filter Banks*, Wellesley-Cambridge, 1996.

[10] J. Laroche and J.-L. Meillier, "Multichannel excitation/filter modeling of percussive sounds with application to the piano," *IEEE Trans. Speech and Audio Processing*, vol. 2, pp. 329–344, 1994.

[11] H.-I. Choi and J. Williams, "Improved time-frequency representation of multicomponent signals using exponential kernels," *IEEE Trans. Acoust., Speech, and Signal Process.*, vol. 37, no. 6, pp. 862–871, June 1989.