

Time and Frequency Decomposition of Head-Related Impulse Responses for the Development of Customizable Spatial Audio Models

KENNETH JOHN FALLER II¹, ARMANDO BARRETO¹, NAVARUN GUPTA² and NAPHTALI RISHE³

Electrical and Computer Engineering Department¹ and School of Computing and Information Science³
Florida International University
Miami, FL 33174
USA

kfall001@fiu.edu <http://dsplab.eng.fiu.edu/>

Department of Electrical and Computer Engineering²
University of Bridgeport
Bridgeport, CT 06604
USA

Abstract: - This paper introduces a new approach to the decomposition of measured Head-Related Impulse Responses (HRIRs) based on simultaneous analysis in the time and frequency domains. This approach is computationally less demanding and faster than previous systematic approaches proposed for this purpose. Currently, HRIRs are the most usual representation of Head-Related Transfer Functions (HRTFs), which are, in turn, the basis of many 3D sound spatialization systems used for PC gaming and virtual reality applications, among others. Many of these applications, however, utilize “generic” sets of HRTFs, which may provide only a sub-optimal spatialization experience. The improved HRIR decomposition method will facilitate our first step towards creating easily customizable HRTF representations, as the process of decomposition yields the sets of parameters that can instantiate a simple functional model to be equivalent to the HRTF represented by the corresponding HRIR. Furthermore, the decomposition breaks down the measured HRIR, into multiple delayed and scaled damped sinusoids, which have characteristics (frequency, decaying rate, initial amplitude and latency), that can be associated with anatomical characteristics of the outer ear of the listener.

Key-Words: - Customizable 3D Spatial Audio, Head-Related Impulse Response (HRIR), Head-Related Transfer Function (HRTF), Prony method, Steiglitz-McBride method.

1 Introduction

The evolution of digital signal processing (DSP) technologies and the increased availability of practical platforms that are capable to implement even its advanced algorithms have facilitated the expansion of applications of “Virtual 3D Sound” to many areas. By “Virtual 3D sound” we mean here the ability to produce in a listener the illusion that a certain sound, generated by the DSP system, is emerging from a source located at a specific (virtual) position in the surroundings of the listener. The virtual source position, from which the sound is seemingly emerging, is normally specified in terms of azimuth (θ), elevation (ϕ) and distance (r) coordinates in a spherical coordinate system centered in the head of the listener, and having the orientation shown in Figure 1.

The reason why the subject perceives a sound as if it originated at the virtual location being simulated is that a binaural pair of signals (Left ear signal and Right ear signal) can be created which closely approximates the

signals that would be received at both of the listener’s eardrums if a “real” sound were created at the actual physical location in the listener’s surroundings.

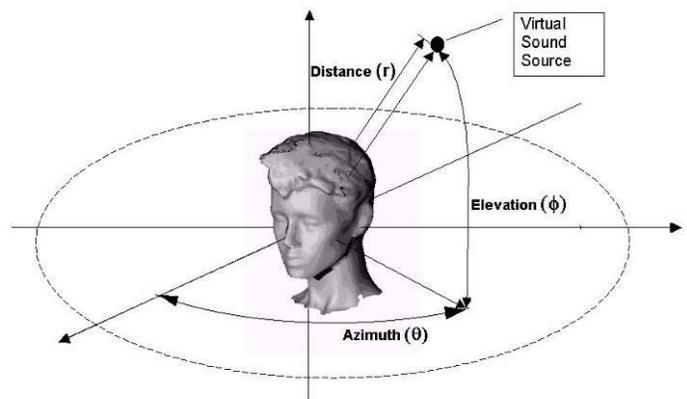


Fig. 1. Diagram of spherical coordinate system.

There are different approaches to the creation of the binaural signal pair. One approach (e.g., Dolby® 5.1 array) uses multiple (more than two) physical sound sources placed all around the listener, at strategically defined positions, to create sounds that combine in the listener’s eardrums, resulting in appropriate binaural acoustic signals into both eardrums. In this approach sounds are actually being generated at positions around the listener, just like they would be in the virtual scenario that is being represented. As such, this setup is not “listener dependent” and a high degree of fidelity in the virtual sound source placement can be archived, regardless of the specific listener using the system. Unfortunately, the multi-speaker approach has clearly a strong dependence on the location of the speakers around the listener, such that there will only be a small spatial region and a defined orientation in which the listener will experience the proper acoustic mixing from the speaker sounds into the correct binaural eardrum signals. Because of this, the multi-speaker approach to 3D virtual sound is applicable only in stationary uses (e.g., home theater).

A fundamentally different approach to the creation of convincing binaural sounds for 3D sound spatialization pursues the creation of the binaural signals digitally, so that they can be delivered directly to the ears of the listener, using stereo headphones. This approach must modify the sound delivered to each of the listener’s ears in the same way the physical environment (e.g., floor, ceiling, walls, listener’s torso, listener’s head, listener’s outer ear) would modify the acoustic signal from its origin at a specific location to each eardrum. Since all the manipulations on the original sound (stored digitally) will be performed by a computer, this approach can be completely portable.

The DSP-based creation of the left and right binaural pair of signals involves the use of special filters, characterized by their impulse responses, which are known as Head-Related Impulse Responses (HRIRs). The transfer functions of these filters are known as Head-Related Transfer Functions (HRTFs). Therefore, every position and each ear will have a specific HRIR. Convolution of a sound signal with the two HRIRs corresponding to a specific source position results in a binaural sound (left channel, right channel) that, when played to a listener through stereo headphones will cause a perception similar to that of a sound emanating from the source location in question, specified by azimuth, elevation and distance (Fig. 1).

Since each person’s head and torso are different, the creation of highly convincing binaural sounds requires the convolution of digital sounds with the pair of HRIRs estimated from each individual listener. However, the

determination of the “individual” HRIR pairs corresponding to varied positions around a specific subject requires the use of specialized and expensive equipment and the involvement of trained personnel, which makes it unaffordable to most users of 3D sound systems. Instead, many applications of spatial audio systems make use of “generic” HRIR pairs obtained from a mannequin of “average anatomical dimensions” (e.g., MIT’s measurements of a KEMAR Dummy-Head Microphone) or using a limited number of subjects to represent the general population (e.g., the CIPIC Database [1]). This type of “generic” HRIRs provides an approximate sense of source locations in many users, but does not have as high spatialization fidelity as individual HRIRs [2].

Our goal is the development of “customizable” HRIRs obtained from a generic dynamic model that could be instantiated differently for each particular listener, by taking into account the relevant physical measurements of the intended listener, in order to still provide a high-fidelity spatialization. Unfortunately, the currently used representation of HRIRs as long (e.g., 128, 256, 512) collections of values obtained as the response to impulse-equivalent functions, such as Golay codes, cannot be altered in any simple way that would factor in the geometrical characteristics of the intended listener. Therefore, we believe that the first step towards customizable HRTFs is the substitution of their current representation in terms of large sets of HRIR sample values with an equivalent functional model requiring the instantiation of a much smaller number of parameters related to the geometry of each intended listener.

2 Methodology

2.1 Simplified Structural Pinna Model

Brown and Duda [3] have proposed that a “structural” model for binaural sound synthesis should “cascade” the effects (e.g., diffraction, inter-aural delay, etc.) of the listener’s head with the local monaural effects of the geometry of the pinna or outer ear. Previous work by Algazi et al. [4] has already yielded a functional model for the listener’s head that can be customized according to 3 simple anatomical measurements. Therefore, the objective of our work is to establish a reduced-parameter pinna model which could be instantiated on the basis of geometrical measurements from the intended listener.

We have previously proposed a pinna model in which the transformation of sound traveling to the eardrum takes place by superposition of a number of reflections of the incoming sound in the ear, which are also affected by the effect of the pinna cavities, such as the concha, acting as resonators [5]. The basic formulation of this model is shown as a block diagram in Figure 2.

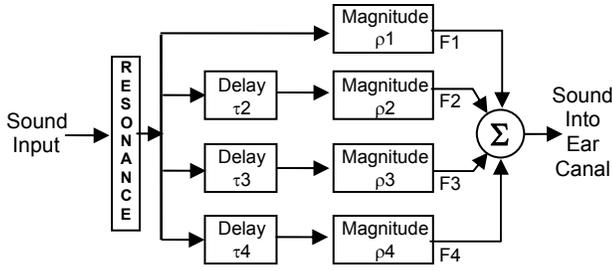


Fig. 2. Block diagram of pinna model (for 4 components).

In this model the parallel paths represent the multiple bounces of the incoming sound wave on the geometrical structures of the pinna. Different trajectories will have different lengths, modeled by different delays, τ_i . Similarly, the loss of energy in each reflection of sound is modeled by a magnitude factor, ρ_i , in each of the paths of the block diagram. If the resonator block in the diagram can be represented by two parameters, such as the angle and the radius of its poles in the Z-plane, the instantiation of a model such as the one shown in Figure 2 would require the definition of only 9 parameters, all of which can reasonably be expected to relate to geometrical measurements of the outer ear of the subject. This approach, therefore, satisfies the requirements of a reduced-parameter pinna model, which could be “cascaded” with Algazi’s functional head model to represent a given HRTF.

Since the parameters of our model cannot be measured directly; our first goal is to convert HRIR sequences measured by specialized equipment (e.g., AuSIM’s HeadZap HRTF measurement system) into specific instances of our model. Once the values for our model are known for multiple source positions and for a large enough number of experimental subjects, from whom anatomical measurements are also available, empirical rules will be developed to assign parameter values from anatomical measurements. Such rules could then be used to assign a custom set of parameters to the model for spatial audio generation for any subject if his/her geometric measurements are known. This paper describes a new approach for the definition of the model parameters from the decomposition of a measured HRIR for one of the ears of a given subject, i.e., the conversion of a “traditional” measured HRIR sequence of values to a smaller set of model parameters.

2.2 HRIR decomposition for the determination of pinna model parameters

According to the model shown in Figure 2 a measured HRIR sequence will be conformed by the superposition of

several damped sinusoids (since the resonator will provide that kind of signal as response to an impulse), appearing scaled by a magnitude factor ρ_i and delayed by a latency τ_i . ($\tau_1 = 0$). Therefore the magnitude factors and delays needed for the model will become apparent if the original measured HRIR sequence is decomposed into a number of scaled (ρ_i) and delayed (τ_i) damped sinusoids.

Previously [5-7] this process of HRIR decomposition into damped sinusoids has been attempted by sequential application of second-order Prony or Steiglitz-McBride (STMCB) signal modeling algorithms to consecutive windows defined on the measured HRIR sequence. The aim of that sequential process was to always restrict the analysis to a partial window of data where only one damped sinusoid (second-order approximation) is expected to be present. So, an initial window is defined from the beginning of the HRIR sequence to a point where the second damped sinusoid is estimated to start, i.e., τ_2 . The amplitude of this first estimated sinusoid is considered to be the value of the magnitude factor ρ_1 . The sinusoid estimated for the initial interval (F1) is then extrapolated to the end of the HRIR sequence and it is subtracted from it, to remove the influence of this first damped sinusoid from the rest of the measured HRIR. At this point the residue obtained is analyzed in the same manner as the original HRIR, except that the origin of analysis is re-established at time τ_2 . The next stage of decomposition will only use the window of data between τ_2 and the point in which the onset of the third damped sinusoid component is estimated to occur, τ_3 . The damped sinusoid estimated in the second stage of decomposition (F2) will also be extrapolated and subtracted from the complete extent of the HRIR remnant still being analyzed. The amplitude of this second estimate will be considered to be ρ_2 . To begin the third stage of decomposition the origin of analysis will be re-established at τ_3 , and the process can be repeated through as many decomposition stages as damped sinusoids are sought. An example of results from the process is shown in Figure 3.

The goal of the process is to obtain a set of damped sinusoids which, when added, form a “reconstructed” HRIR that is a good approximation of the original, measured HRIR sequence. The goodness of fit of the reconstructed HRIR is assessed by means of Equations (1) and (2), where MS means mean square value:

$$\text{Error} = \text{Original HRIR} - \text{Reconstructed HRIR}, \quad (1)$$

$$\text{Fit} = [1 - \{\text{MS}(\text{Error})/\text{MS}(\text{Original HRIR})\}]. \quad (2)$$

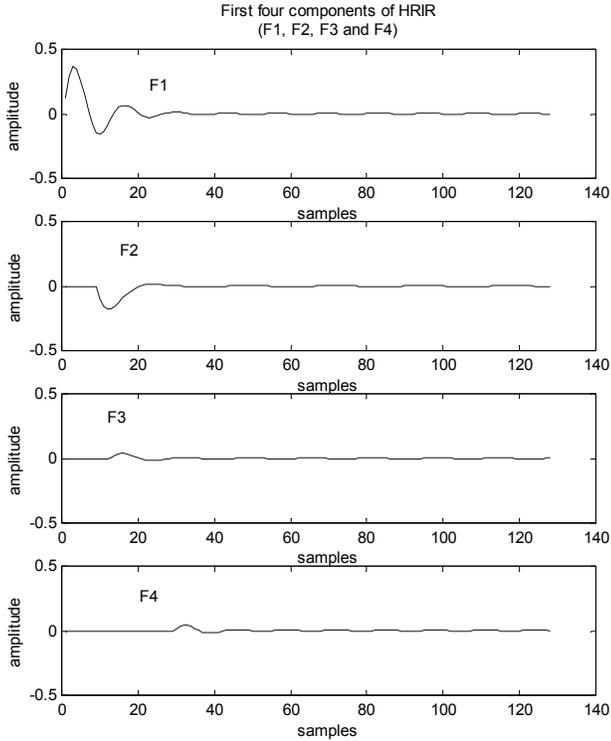


Fig. 3. Four damped sinusoidal components obtained from a measured HRIR.

It has been observed that proper recognition of the boundaries for the windows to be analyzed (τ_2 , τ_3 , etc.) was critical to the achievement of high fit values. Originally, these break points were determined in each stage of decomposition by tentatively widening the window of analysis, finding a tentative damped sinusoid and calculating the mean square error (MSE) between the tentative approximation and the values of the HRIR under analysis, within the tentative window. Typically the MSE value found would decrease as the window was widened, until it would reach the onset of the next damped sinusoidal present, where the MSE would spike up. It was later found [6, 7] that an exhaustive search which tried all probable widths (range of 2 to 10 sampling intervals of 10.4 μ s each) for all the sequential windows and selected the parameters from the combination of widths that resulted in the larger overall fit (between the complete reconstructed HRIR and the original measured HRIR) provided better average fits for a database of 14 subjects including a total of 2016 HRIRs (2 ears, 72 source positions per subject). However, the exhaustive search approach is extremely computationally intensive, even with just the 5 windows processed in those studies. In fact, the tree-diagram needed to track all possible width combinations of 5 sequential windows has

$9 \times 9 \times 9 \times 9 \times 9 = 59,049$ leaf nodes and the addition of any subsequent windows with this approach will multiply the number of leaf nodes by 9, per additional window. To truly select the best of all possible alternatives, all the branches of the tree need to be explored and the reconstructed HRIR defined at each leaf node compared with the measured HRIR to assess its fit. It became clear that increasing the number of windows of analysis (which may be necessary to model late components in the HRIRs) would be impractical using the exhaustive search method. This has prompted us to develop a new, faster method of HRIR decomposition into sequential damped sinusoids.

2.3 Pole Approximation of Damped Sinusoids

The goal of this new method is to avoid having to pre-set the width of each sequential window of data analyzed in each subsequent stage of the decomposition process. The need to isolate small windows of data was connected to the assumption that such windows could be defined so that they would only contain a single damped sinusoid and not a superposition of several of them. Under that premise the Prony or STMCB second order algorithms were applied in each window, seeking to approximate a single damped sinusoid. In general, a single damped sinusoidal component sequence will be represented by a conjugate pair of poles within the unit circle and a zero at the origin of the Z-plane (Figure 4) [8]. Hence, a damped sinusoid in the Z-domain can be described with the following general equation:

$$X(z) = \frac{k \cdot z}{(z - p_1)(z - p_2)} \quad (3)$$

where k is a scalar and p_1 and p_2 are complex poles. According to Equation 3, if the scalar k and the poles are known then, using the inverse Z-transform, it is possible to find the time domain representation of a damped sinusoid. Based on these considerations, the new approach to the decomposition of HRIRs into scaled and delayed damped sinusoids will use the complete remnant of the HRIR available for processing at each decomposition stage (instead of a bounded window), searching for multiple damped sinusoids by application of a higher-order STMCB approximation for the whole remnant. The method will then isolate individual damped sinusoids by their complex conjugate pole signatures in the Z-domain, pursuing as many alternative outcomes as complex pairs can be identified for the specific decomposition stage in question.

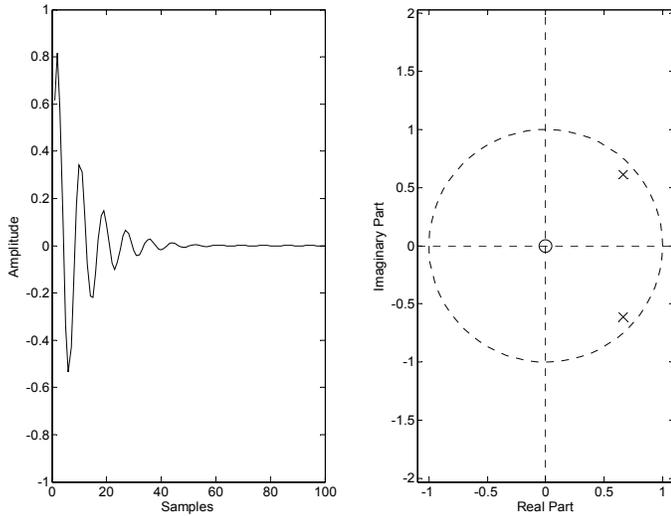


Fig. 4: Time domain and Zero-Pole plot of a single damped sinusoidal.

This also results in a tree-search approach. However, the branching factor of this search tree starts at the amount of damped sinusoids being extracted from the whole HRIR but decreases by one in every subsequent stage of the decomposition, which makes the number of leaf nodes much smaller than for the previous algorithm. For example, if 5 damped sinusoids will be extracted, only $5 \times 4 \times 3 \times 2 \times 1 = 5! = 120$ leaf nodes will exist.

The details of the process are explained with a simulated example in the paragraphs below.

The example addresses the decomposition of a synthetic signal created by summing three delayed damped sinusoids using this new “pole approximation” method. The sinusoids were created using the following equation:

$$x_i(n) = e^{d_i * n} \cdot \sin(\omega_d \cdot \pi \cdot n) \quad (4)$$

where N is the length of the signal, $n = 1, \dots, N$, d_i is the negative damping factor and ω_d is the digital frequency. Once the three sinusoids (x_1 , x_2 and x_3) are created, the desired delays (τ_2 and τ_3) are applied to the last two sinusoids respectively, resulting in x_{2s} and x_{3s} . Finally, the sinusoids are then summed point-to-point to produce the test signal (x). In this example $N=100$, $\tau_2=3$, $\tau_3=6$, $\omega_d=0.711$, $d_1=-0.1$, $d_2=-0.125$ and $d_3=-0.15$. The three signals (x_1 , x_{2s} and x_{3s}) and the resulting signal (x) are shown in Figure 5.

In this example the goal is to decompose x into three damped sinusoids. Therefore, the process starts by applying a sixth-order STMCB approximation process to the complete x sequence. The results from the sixth-order STMCB approximation will have the pole structure shown

in Figure 6. As seen in Figure 6, x will result in 2 conjugate pairs. These pairs will be used to compute two separate impulse responses. One of which should approximate x_1 accurately.

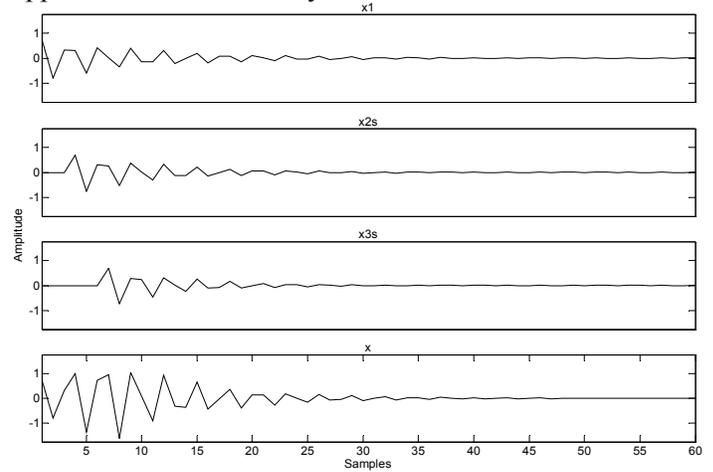


Fig.5: Plot of the three damped sinusoids (x_1 , x_{2s} with delay τ_2 and x_{3s} with delay τ_3) and the sum of them (x).

The damped sinusoidal impulse responses associated with the conjugate pole pairs shown in Figure 6 will be investigated as candidates to represent the first sinusoidal present in x (i.e., there will be up to three branches at the initial node of this search tree, if all the poles were complex). The investigation of each of these alternatives involves its subtraction from x to define a residue sequence, as shown in Figure 7, which will then be thresholded. The threshold level used for this segmentation was set at 25% of the signal peak, in this synthetic example. A slightly different threshold to process real HRIRs was found as described in the following sections.

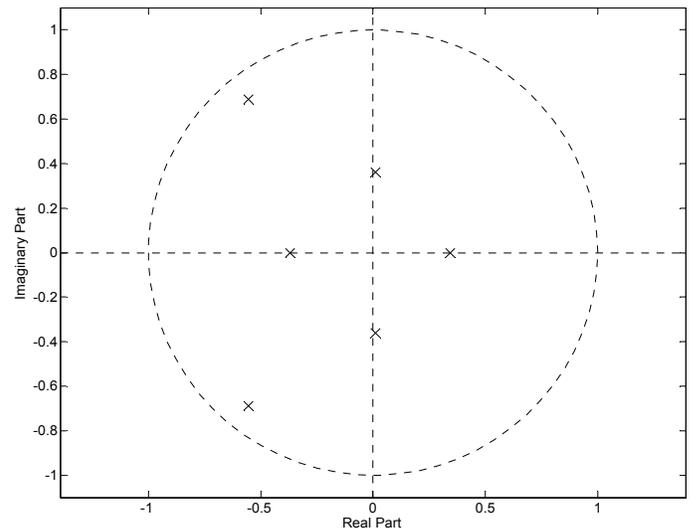


Fig.6: Poles obtained from the sixth-order STMCB approximation of the complete sequence x .

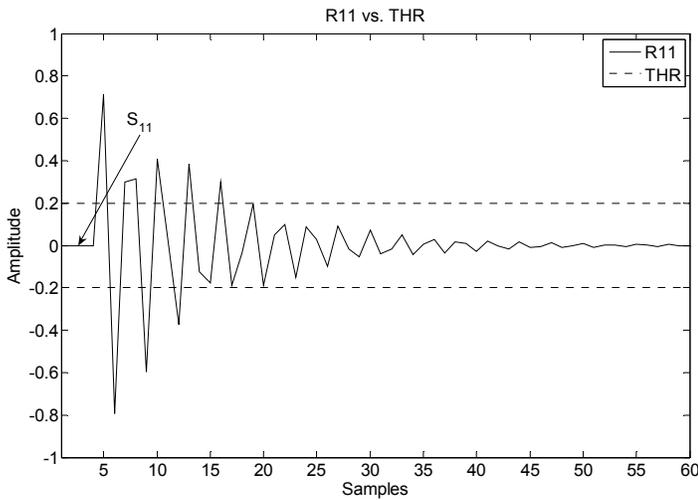


Fig.7: Plot of R11 with threshold lines (THR).

The time at which the residual surpasses this threshold will be considered the onset of the next damped sinusoidal, i.e., the estimate of τ_2 . As in the previous method, the decomposition process will continue on to a second stage after re-establishing the origin of analysis at the estimated τ_2 . The assumption made in every subsequent decomposition stage is that there should be one less damped sinusoidal present in the new remnant (since one has just been removed in the previous stage). As such, a fourth-order STMCB approximation will be applied in the second decomposition stage, yielding 4 poles, which will then be used to synthesize up to two candidates for the second damped sinusoid extracted from x . The same pattern of steps will be applied through all subsequent stages of the decomposition, until the stage in which a second-order STMCB approximation will be applied to the last remnant to identify the last damped sinusoid.

After M stages of decomposition there will be $M!$ leaf nodes in the search tree, each representing a set of M delayed and scaled damped sinusoids that, when added together, form candidate approximations to the original signal x . The fit of each of those $M!$ candidate approximations with respect to x will be evaluated (Equations 1 and 2) and the candidate with the highest fit will be selected as the final decomposition of x . In our example, the winning candidate approximation had a 99.99% fit with the original x , and the individual damped sinusoids obtained through each stage of decomposition also matched x_1 , x_2 s and x_3 s very closely. Figure 8 shows the steps in detail for this 3-component example.

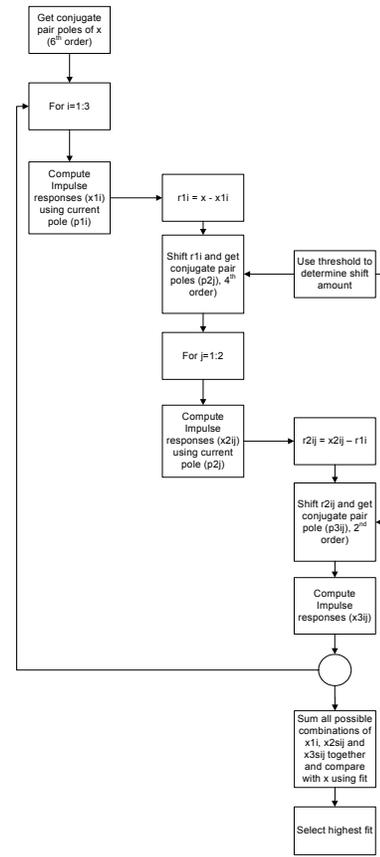


Fig. 8: Flow chart for decomposition.

3 Decomposition of Measured HRIRs

The method described in Section 2.3 was applied to the decomposition of 14 actual HRIRs, recorded from 14 subjects using the AuSIM HeadZap system at Florida International University. The goal in each case was to obtain $M = 5$ damped sinusoidal components. Therefore, the order of the first STMCB approximation process was 10. The procedure was identical as the one explained for the decomposition of the synthetic sequence x , in Section 2.3, with the exception that an empirically-defined slightly different threshold level was applied to each reduced remnant of the HRIR.

The empirical determination of the best threshold level to use in decomposing actual HRIR signals was performed by sweeping through thresholds from 0.005 to 0.4 in increments of 0.005 for HRIRs corresponding to sound sources at $\pm 90^\circ$ azimuth (i.e., directly lateral from the ear measured) and elevations from -36° to 54° at increments of 18° . The resulting fits for all these thresholds for each of the 14 subjects was recorded. Average fits for each threshold were calculated for all the elevations. For example, Figure 9 is the plot of threshold vs. average fit for elevation $\phi = -36^\circ$. As can be seen in this plot, there is a curvature which has a maximum at a threshold value of

about 0.18. This was also apparent in the other plots and, as a result, 0.18 (i.e., 18% of the signal peak) was selected.

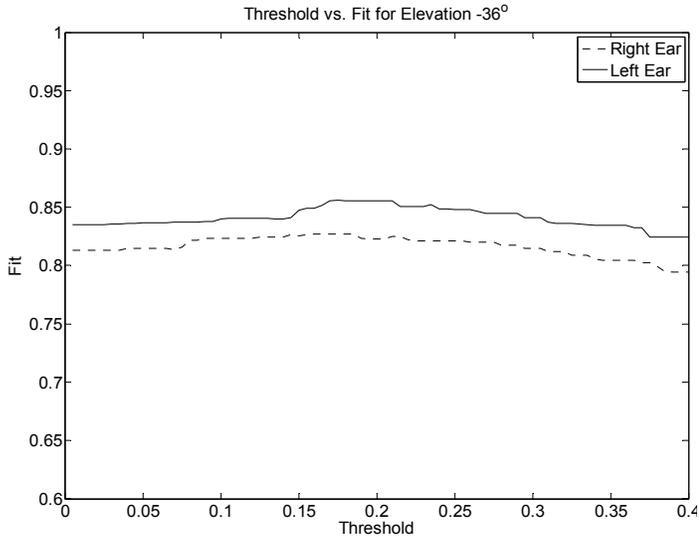


Fig.9: Threshold versus fit for elevation -36° .

The final fit achieved by the new “pole approximation” method using the threshold at 18% of the signal peak was recorded and compared to the fit achieved using the previous, exhaustive search method. The results for each ear are displayed in Table I.

TABLE I: HRIR DECOMPOSITION RESULTS

METHOD:	Average Fit – Left Ear	Average Fit – Right Ear
Exhaustive, variable window width (old)	97.57%	97.57%
Pole approximation w/ Threshold (new)	91.03%	89.96%

Figures 10 to 12 show the highest, average and lowest fits for HRIRs using the “new” method, respectively.

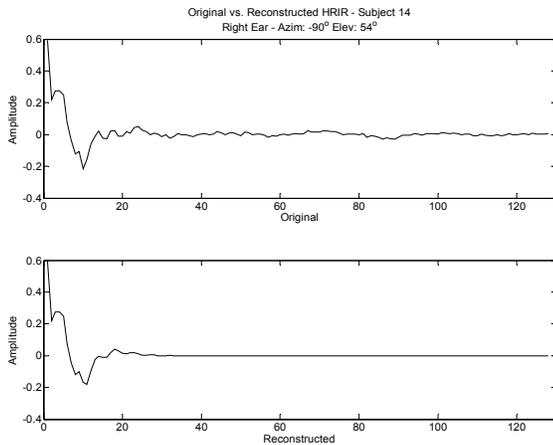


Fig.10: Original (top) vs. reconstructed HRIRs for the right ear of subject 14 for azimuth -90° and elevation 54° .

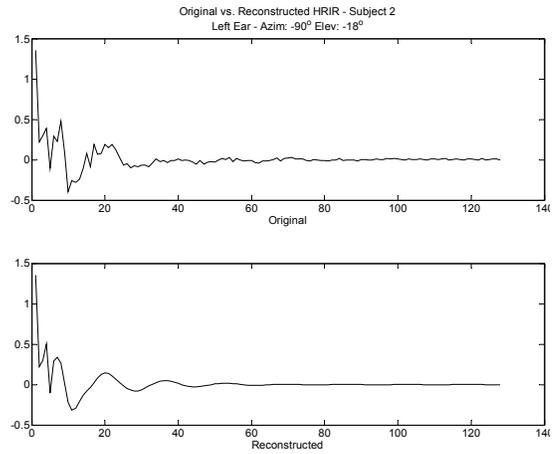


Fig.11: Original (top) vs. reconstructed HRIRs for the left ear of subject 2 for azimuth -90° and elevation -18° .

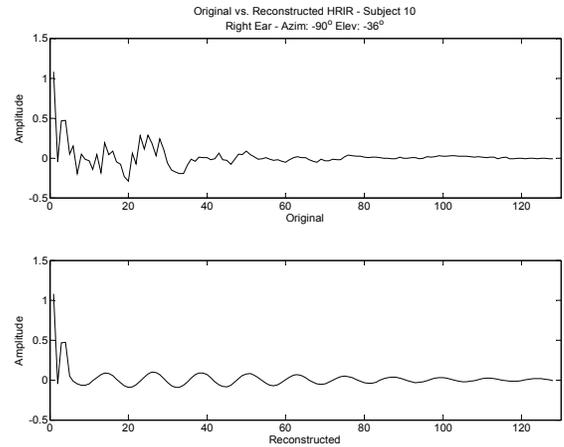


Fig.12: Original (top) vs. reconstructed HRIRs for the right ear of subject 10 for azimuth -90° and elevation -36° .

In our results the best fits were obtained for positive elevations (average 93%) and for zero elevation (average 92%). Negative elevations yielded slightly lower fits (average 86%). This might be due to the potential for more numerous longer-latency reflections when the sound is bounced in the multiple pinna structures above the ear canal (when the source is placed at negative elevations). It is possible there might have been HRIR components beyond the fifth sinusoidal in those cases, which were not addressed by the decomposition process.

4 Conclusions

The results shown in Table I indicate that the “old” method achieved a slightly higher average fit, but exhibited several drawbacks. First, the average calculation time was found to be about 100 times longer for the “old” method when a test set of 14 HRIRs were decomposed by both approaches (429 s to 4.2 s). Secondly, when the delay

is small (less than 5 samples wide), the second-order STMCB sequential method alone tends to inaccurately reconstruct the signal. As an example, using x from Section 2.3, the “new” and “old” methods were used to decompose x to obtain x_1 . The results are displayed in Figure 13. The top plot shows the original sinusoid (x_1), the middle plot shows the results of the “new” method (x_{1rn}) and the bottom plot shows the results of the “old” method (x_{1ro}). Clearly, the “new” method captures more of the time domain features of the original signal x_1 . Additionally, the accuracy of the reconstructed signal is closer for the x_{1rn} in the frequency domain as well. In Figure 14, the complex conjugate poles and the zeros for x_{1rn} lie within a close proximity of x_1 's poles and zeros whereas the poles and zeros for x_{1ro} fall along the real axis. This indicates that the “new” method was able to retain more of the spectral features of x_1 as well.

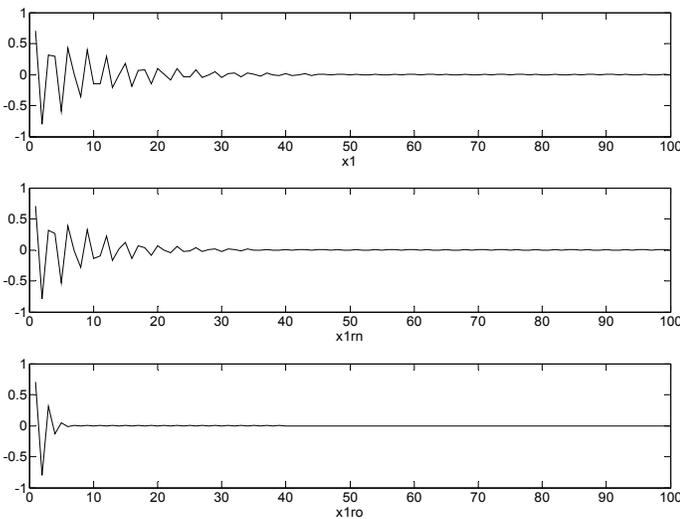


Fig. 13: x_1 (top) vs. x_{1rn} (middle) and x_{1ro} (bottom).

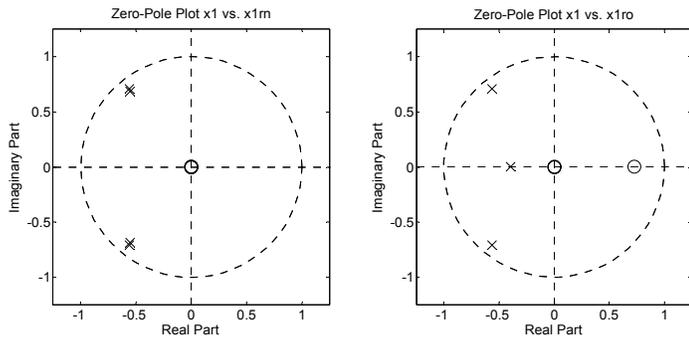


Fig. 14: Zero-pole plot of x_1 vs. x_{1rn} (left) and x_1 vs. x_{1ro} (right).

Similar experiments in which the different signal components were placed at different delays from each other seem to indicate that the new method described in

this paper performs better than the previous method in decomposing test signals when the synthetic delays τ_i were small. According to these observations, it seems that the new time-frequency decomposition approach presented here might be particularly beneficial when the delayed sinusoidal components in the HRIRs are closely packed together (in time), which is expected to be the case for source locations close to the inter-aural axis.

5 Acknowledgements

This work was sponsored by NSF grants IIS-0308155, CNS-0520811, HRD-0317692 and CNS-0426125.

References:

- [1] R. O. Duda, "3-D Audio for HCI," [Online document], [2006 Jan 19], Available at HTTP: http://interface.cipic.ucdavis.edu/CIL_tutorial/3D_home.htm
- [2] E. M. Wenzel, M. Arruda, D. J. Kistler, and F. L. Wightman, "Localization Using Nonindividualized Head-Related Transfer-Functions," *Journal of the Acoustical Society of America*, vol. 94, pp. 111-123, 1993.
- [3] C. P. Brown and R. O. Duda, "A structural model for binaural sound synthesis," *Ieee Transactions on Speech and Audio Processing*, vol. 6, pp. 476-488, 1998.
- [4] V. R. Algazi, C. Avendano, and R. O. Duda, "Estimation of a spherical-head model from anthropometry," *Journal of the Audio Engineering Society*, vol. 49, pp. 472-479, 2001.
- [5] A. Barreto and N. Gupta, "Dynamic Modeling of the Pinna for Audio Spatialization," *WSEAS Transactions on Acoustics and Music*, vol. 1, pp. 77-82, January 2004.
- [6] K. J. Faller II, A. Barreto, N. Gupta, and N. Rische, "Decomposition and Modeling of Head-Related Impulse Responses for Customized Spatial Audio," *WSEAS Transactions on Signal Processing*, vol. 1, pp. 354-361, 2005.
- [7] K. J. Faller II, A. Barreto, N. Gupta, and N. Rische, "Performance Comparison of Two Identification Methods for Analysis of Head Related Impulse Responses," in *Advances in Systems, Computing Sciences and Software Engineering*, T. Sobh and K. Elleithy, Eds. Netherlands: Springer, 2006, pp. 131-136.
- [8] L. P. Charles and H. T. Nagle, *Digital control system analysis and design (3rd ed.)*: Prentice-Hall, Inc., 1995.