

# Manifold Learning using Isomap applied to Spatial Audio Personalization

Felipe Grijalva\*, Siome Klein Goldenstein<sup>†</sup>, Dinei Florencio<sup>‡</sup>, and Luiz César Martini\*

\*School of Electrical and Computer Engineering, University of Campinas, Campinas, Brasil

<sup>†</sup>Institute of Computing, University of Campinas, Campinas, Brasil.

<sup>‡</sup>Multimedia, Interaction and Communication Group, Microsoft Research, Redmond, WA, USA

**Abstract**—As augmented reality applications become more important, there is increasing effort in spatial audio research. The term spatial audio or 3D sound refers to techniques where a person’s anatomy is modeled as digital filters. By filtering a sound source with these filters, a listener is capable of perceiving a sound as though it were reproduced at a specific spatial location. In the frequency domain, these filters are known as Head-Related Transfer Functions (HRTFs). A significant problem for the implementation of 3D sound systems is the fact that spectral features of HRTFs differ widely among individuals due to their anatomical differences. Thus, it is necessary to personalize them to guarantee high quality sound perception. With this aim, we introduce a new anthropometric-based method for customizing of HRTFs in the horizontal plane using manifold learning. The method uses Isomap, artificial neural networks (ANN), and a neighborhood-based reconstruction procedure. We first modify Isomap’s graph construction step to emphasize the individuality of HRTFs and perform a customized nonlinear dimensionality reduction of the HTRFs. We then use an ANN to model the nonlinear relationship between anthropometric features and our low-dimensional HRTFs. Finally, we use a neighborhood-based reconstruction approach to reconstruct the HRTF from the estimated low-dimensional version. Simulations show that our approach performs better than PCA (Principal Component Analysis) and confirm that Isomap is capable of discovering the underlying nonlinear relationships of sound perception.

**Keywords**—HRTF; Isomap; Manifold Learning; Spatial Audio

## I. INTRODUCTION

The objective of spatial audio or 3D audio is to simulate a sound source in arbitrary spatial locations. The core components of spatial audio are the so-called Head-Related Impulse Responses (HRIRs) or their frequency-domain representation Head-Related Transfer Functions (HRTFs). HRTFs model the spectral filtering of a sound source caused by the head, pinna (i.e. the outer part of the ear ) and torso before it reaches the eardrum. By filtering a sound source with these filters, a listener is capable of perceiving a sound as though it were reproduced at a specific location in space [1].

Spatial audio has a wide range of applications from hearing aids and entertainment (e.g. home theaters, video games) to virtual reality [2] (e.g. Oculus Rift<sup>TM</sup>, Google Glass<sup>TM</sup>, air traffic controllers [3]). In fact, as virtual reality applications become more important, there is increasing research effort in the spatial audio research. In this sense, several works has proposed the use of spatial audio as natural user interface for sensory substitution and augmented reality prototypes aimed at visually impaired people [4], [5], [6], [7].

It is precisely in this type of application that the project “Vision for the blind: translating 3D Visual Concepts into 3D Auditory Clues” focused through the Microsoft Research/Fapesp cooperation agreement 2012/50468-6. The goal of this project was to construct and validate a complete proof-of-concept assistive device for the blind. This device uses computer vision algorithms to extract high-level 3D information from a Microsoft Kinect Sensor and communicates this information to the visually impaired user using 3D audio to exploit the inherent spatial sense of the auditory system.

With this in mind, the main objective of the M.Sc. dissertation <sup>1</sup> associated to this paper, in the context of the aforementioned project, was to provide the theoretical basis and characteristics of HRTFs as the main components of spatial audio. Besides, we proposed a novel approach for HRTF personalization using the manifold learning technique, Isomap.

This paper is organized as follows. Section II gives an overview of HRTFs and why we need to personalize them. Section III analyze related works using machine learning techniques for HRTF personalization. Section IV introduces our approach for HRTF personalization using Isomap. In Section V, we performed experiments and compared our approach to a PCA-based method.

## II. THE NEED OF HRTF PERSONALIZATION

The sound from an audio source reaches both ears after interacting with the anatomical characteristics of the individual (i.e. head, torso and pinna). The resultant signal contains statics cues (i.e. binaural and monaural cues) that in conjunction with dynamic cues (i.e. produced by head movements) define our three dimensional perception of audio [1]. The static cues are modeled through HRTFs that encode binaural cues, such as the Interaural Time Difference (ITD, i.e. difference in arrival time of a sound between two ears) and the Interaural Level Difference (ILD, i.e. difference in air pressure level of a sound between two ears). Besides, they also encode monaural cues that are mainly produced by the pinna. A pair of complex-valued HRTFs at distance  $r$ , azimuth  $\theta$  and elevation  $\phi$ , for

<sup>1</sup>M.Sc Dissertation directed by Luiz César Martini (FEEC-Unicamp), co-directed with Siome Klein Goldenstein (IC-Unicamp) and in collaboration with Dinei Florencio (Microsoft Research).

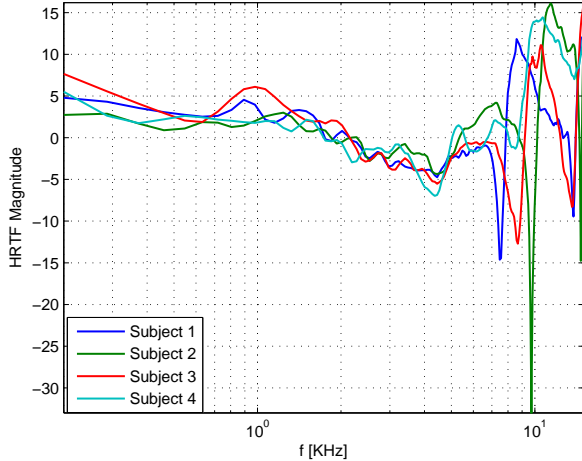


Fig. 1. The left HRTF at a specific location for four different subjects. The variability among subjects, more evident at high frequencies, is due to the monaural cues introduced by the pinna.

left and right ear,  $H_L$  and  $H_R$ , are

$$\begin{aligned} H_L(r, \theta, \phi, f, a) &= \frac{P_L(r, \theta, \phi, f, a)}{P_0(r, f)}, \\ H_R(r, \theta, \phi, f, a) &= \frac{P_R(r, \theta, \phi, f, a)}{P_0(r, f)}, \end{aligned} \quad (1)$$

where,  $P_L$  e  $P_R$  represent the pressure level for each ear and  $P_0$  represents the pressure level in the center of the head with the head absent [8]. The variable  $a$  depends on the subject's anatomy. If  $r > 1$  m, the HRTFs do not depend on distance and are called *far-field HRTFs*. Otherwise, they are called *near-field HRTFs*. Near-field HRTFs are beyond the scope of this study and when we refer to HRTFs, we are referring to far-field HRTFs.

A significant problem for the implementation of 3D sound systems is the fact that spectral features of HRTFs differ widely among individuals [1]. This inter-subject variability is due to their anatomical differences such as the head dimensions and the pinna shape and size. Figure 1 shows the left HRTF at a specific location for four different subjects. The inter-subject variability is specially notable for high frequencies (i.e.  $f > 4$  KHz) where the monaural cues introduced by the pinna are more prominent [1]. For this reason, the pinna is considered as the acoustic fingerprint of a subject. Various studies show a decrease in localization accuracy due to nonindividualized HRTFs [9], [10], often producing front/back reversals (i.e. the listener perceive a sound source at a front location as though it were coming from a back location, or vice versa), poorly sound externalization (i.e. the subject perceives the sound inside her head) and incorrect elevation perception. Thus, it is necessary to personalize HRTFs to guarantee high quality 3D sound perception. The most accurate approach of personalizing HRTFs is through direct measurements. This way, for each spatial location, a loudspeaker reproduces a sound signal which in turn is captured by a microphone

inserted in the subject's ear. Then, the captured microphone signals are processed to extract their corresponding HRTFs [8]. Note that custom HRTF measurement involves expensive apparatus (e.g. an anechoic chamber, low-noise microphones) and it is a complex, time consuming, and not scalable procedure [8]. To avoid HRTF measurements, several theoretical models (spherical head model [11], the snowman model [12]) and numerical methods (boundary element method [13]) have been proposed. Nevertheless, theoretical models are approximations of complicated anatomy and numerical methods are computationally intensive.

On the other hand, since HRTFs are closely related to certain anthropometric parameters, they can therefore be customized from anthropometric measurements [14]. Anthropometric-based regression methods predict the individualized HRTFs of a new subject using a model derived from a baseline database. Usually, some dimensionality reduction is applied to the HRTFs prior to customization. It is this kind of HRTF personalization methods that this work focuses on.

### III. PRIOR WORK

Nishino et al. [15] performed Principal Component Analysis (PCA) on the log magnitude HRTFs in the horizontal plane for each direction and ear separately. Then, linear regression analysis for each direction and ear is applied on a baseline database, using 9 anthropometric parameters as inputs and 5 PCA weights as outputs. For a new subject outside the training database, the PCA weights are predicted from the linear models and then used to reconstruct the log magnitude of HRTFs. Finally, minimum-phase reconstruction [16] estimates the final complex-valued HRTFs.

Due to the inability of linear methods (such as PCA) to represent the complex relationship between HRTF and multiple variables (i.e. direction, frequency and individual), Grindlay et al. [17] introduced a multilinear tensor framework representation for HRTF decomposition. The tensor has 3 modes: frequency mode, direction mode and subject mode. A single linear regression model is used for mapping anthropometric features to a 5 dimension vector representing the subject mode in the tensor. Li et al. [18] employ a similar approach for dimensionality reduction but instead of linear regression, they use an artificial neural network (ANN).

Moreover, nonlinear techniques have been applied to both dimensionality reduction of HRTFs (e.g. Isomap, Locally Linear Embedding) and to regression of HRTFs based on anthropometric features (e.g. Support Vector Regression [19], ANNs [20], [18]). In [21], Duraiswami et al. present an exploratory study on learning the nonlinear manifold structure in vertical plane HRTFs using Locally Linear Embedding (LLE). They also propose a new method for HRTF interpolation and a new distance metric between two HRTFs based on the geodesic distance on the learned manifold.

Kapralos et al. [22], [23] conducted a comparative study from a quantitative point of view between PCA, Isomap and LLE for HRTF dimensionality reduction, finding that Isomap and LLE perform better than PCA in subjective experiments.

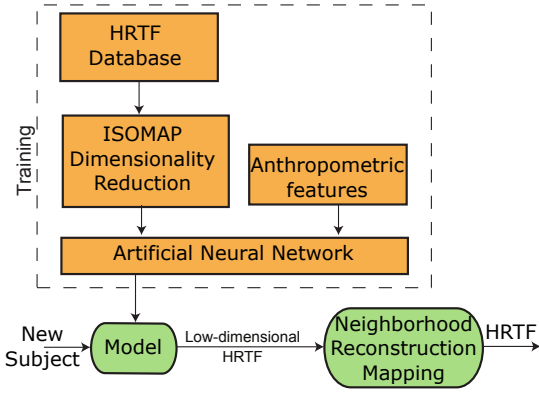


Fig. 2. HRTF Personalization using Isomap.

As in [20], [18], we employ an ANN for regression to predict the HRTFs for a new subject based on his anthropometric parameters. Unlike this prior work, we use nonlinear reduction technique, Isomap, to construct a manifold structure in horizontal plane HRTFs..

Our work is inspired on the successful results by Duraiswami et al [21] and Kapralos et al [22], [23] using LLE and Isomap for HRTF interpolation and dimensionality reduction. Their findings support the idea suggested by Seung et al [24] that nonlinear manifold techniques are crucial for understanding how perception arises from the dynamics of neural networks in the brain. However, neither of them addresses the customization of HRTFs as we do.

As in previous work [15], [20], we use the minimum phase approximations for HRTFs, a minimum-phase function cascaded with a pure delay [16]. In practice, the pure delay is the ITD (Inteaural Time Difference) and it is commonly cascaded in either the left or right HRTF of each left-right HRTF pair [15]. Calculation of ITD is beyond the scope of this paper. Several studies address the ITD calculation based on anthropometric parameters, notably in [25]. Here, we focus only on the spectral features of HRTFs magnitude and, unless otherwise stated, when we refer to HRTF we are referring to its magnitude.

#### IV. HRTF PERSONALIZATION USING ISOMAP

Figure 2 summarizes our HRTF personalization method. First, we reduce the HRTF dimensionality using Isomap. Then, we train an ANN with anthropometric parameters as inputs and the low-dimensional HRTFs as output. For each new subject with known anthropometric features, the ANN model predicts the low-dimensional HRTF representation. Finally, we use neighbor reconstruction mapping to recover the high-dimensional HRTFs from the low-dimensional space.

##### A. Dimensionality Reduction using Isomap.

In general, dimensionality reduction algorithms provide a method for taking a dataset represented in a  $D \times N$  matrix  $\mathbf{X}$  consisting of  $N$  sample vectors  $\mathbf{x}_i$ , i.e.  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^D$  and calculating a corresponding low-dimensional representation in a  $d \times N$  matrix  $\mathbf{Y} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\} \subset \mathbb{R}^d$ , where

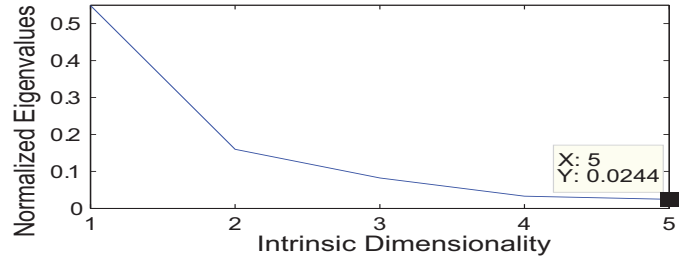


Fig. 3. Intrinsic Dimensionality Estimation.

$d < D$ . Here, consider all HRTFs in the horizontal plane as points in the  $D$  high-dimensional space.

Isomap is a nonlinear dimensionality reduction algorithm, first introduced in [26]. The first step in the Isomap algorithm is to construct a graph  $G(V, E)$  on the input data set  $\mathbf{X}$ . Each sample  $\mathbf{x}_i \in \mathbf{X}$  is represented by a node  $v_i \in V$ , and two nodes  $v_i$  and  $v_j$  are connected by an edge  $(v_i, v_j) \in E$  with length  $d_{\mathbf{X}}(\mathbf{x}_i, \mathbf{x}_j)$  if  $\mathbf{x}_i$  is one of the  $K$  nearest neighbor of  $\mathbf{x}_j$ . The edge length  $d_{\mathbf{X}}(\mathbf{x}_i, \mathbf{x}_j)$  is given by the Euclidean distance between  $\mathbf{x}_i$  and  $\mathbf{x}_j$  [26], [27].

The second step in Isomap involves computation of the shortest paths between all nodes in  $G$ . Distances are stored pairwise in a matrix  $\mathbf{D}_G$ . The distance matrix  $\mathbf{D}_G$  represents geodesic distances between all samples on the manifold [27]. Because these distances are Euclidean, Isomap makes the same assumption of local linearity as LLE [27].

The third and final step is to construct the  $d$ -dimensional embedding calculating the eigenvectors of  $\tau(\mathbf{D}_G)$ , where  $\tau(\mathbf{D}) = -\mathbf{H}\mathbf{S}\mathbf{H}/2$  and  $S_{ij} = D_{ij}^2$  ( $\mathbf{S}$  is the matrix of squared distances) and  $H_{ij} = \delta_{ij} - 1/N$ . Recall that  $N$  is the number of sample points and  $\delta$  is the Kronecker delta function. Finally, let  $\lambda_p$  be the  $p^{\text{th}}$  eigenvalue (in decreasing order) of the matrix  $\tau(\mathbf{D}_G)$ , and  $v_p^i$  be the  $i^{\text{th}}$  component of the  $p^{\text{th}}$  eigenvector. Then set the  $p^{\text{th}}$  component of the  $d$ -dimensional coordinate vector  $\mathbf{y}_i$  equal to  $\sqrt{\lambda_p} v_p^i$  [26].

Isomap first step is the construction of a graph. The simplest approach is to select, for each data point, a fixed number of nearest neighbors,  $K$ , as measured by Euclidean distance. Other criteria, however, can also be used to choose neighbors, and in general, neighborhood selection in Isomap presents an opportunity to incorporate a priori knowledge [28].

We know that some correlation exists due to left-right symmetry of HRTFs at frequencies below 5.5 KHz [29]. Moreover, to emphasize the individuality of HRTFs across directions, Nishino et al. [15] perform PCA reduction separately for each direction and ear. Here, instead of applying Isomap separately for each direction and ear, we propose construct the graph taking into account this knowledge.

One of our contributions is our graph  $G$  construction procedure. Consider again the high-dimensional dataset in a  $D \times N$  matrix  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^D$  formed by  $N$  HRTFs of two ears of  $P$  subjects at  $M$  azimuths in the horizontal plane (i.e.  $N = 2 \cdot P \cdot M$ ).

We connect each datapoint  $\mathbf{x}_i$  to  $K = 2P + 1$  neighbors and we set its edge lengths to  $s_{ij} d_{\mathbf{X}}(\mathbf{x}_i, \mathbf{x}_j)$ , where  $s_{ij}$  is a scale factor, according to the following rules:

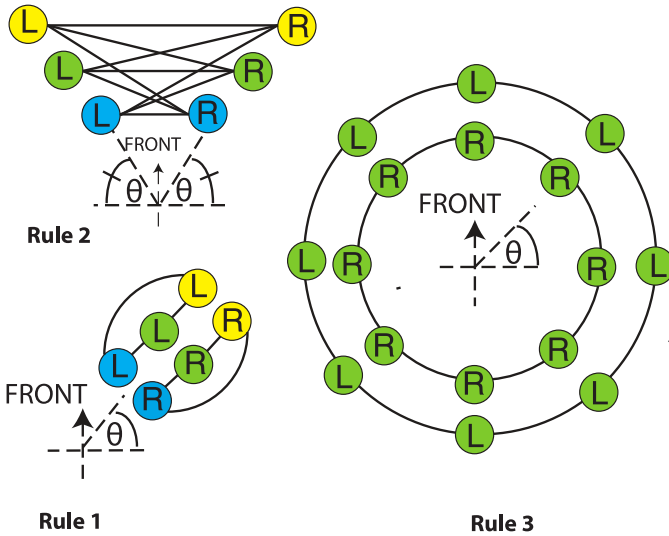


Fig. 4. Example of graph construction procedure for  $P = 3$  subjects. Each color represents a subject. L=Left, R=Right and  $\theta$  is the azimuth angle.

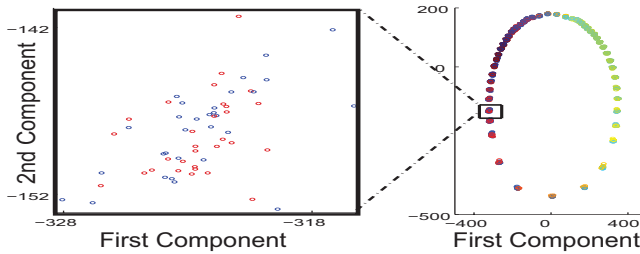


Fig. 6. Variability inside clusters due to inter-subject differences. Red and blue points represent HRTFs at symmetric azimuths of left and right ears respectively

**Rule 1.** If  $\mathbf{x}_i$  and  $\mathbf{x}_j$  represent HRTFs of the same azimuth and ear but different subject, then connect them and set  $s_{ij} = 1/100$  in order to emphasize the individuality of HRTFs across directions.

**Rule 2.** Let  $\theta_i$  and  $\theta_j$  be azimuths of HRTFs represented by  $\mathbf{x}_i$  and  $\mathbf{x}_j$  respectively. Regardless of the subject, if  $\mathbf{x}_i$  and  $\mathbf{x}_j$  represent HRTFs of opposite ears and  $\theta_j$  is the mirror horizontal azimuth of  $\theta_i$  (i.e.  $\theta_j = 360 - \theta_i$ ), then connect them and set  $s_{ij} = 1/100$  in order to take advantage of left-right symmetry.

**Rule 3.** Let  $\theta_i$  and  $\theta_j$  be azimuths of HRTFs of the same subject represented by  $\mathbf{x}_i$  and  $\mathbf{x}_j$  respectively. If  $\theta_j$  is the nearest azimuth greater than  $\theta_i$  or if  $\theta_j$  is the nearest azimuth less than  $\theta_i$ , then connect  $\mathbf{x}_i$  and  $\mathbf{x}_j$  and set  $s_{ij} = 1$ .

In order to clarify how the above mentioned rules were applied, Figure 4 shows an illustrative example.

Before applying Isomap, we first need to select the number of neighbors,  $K$ , and the intrinsic dimensionality,  $d$ . Due to our proposed graph construction explained above, the number of neighbors is set to  $K = 2P + 1$ , where  $P$  is the number of subjects on the dataset  $\mathbf{X}$ . The intrinsic dimensionality was estimated analyzing the residual variance. Figure 3 shows the normalized eigenvalues (in decreasing order) calculated over the complete dataset  $\mathbf{X}$ . Since eigenvalues give the variance in each dimension, when they are lower than a threshold, little

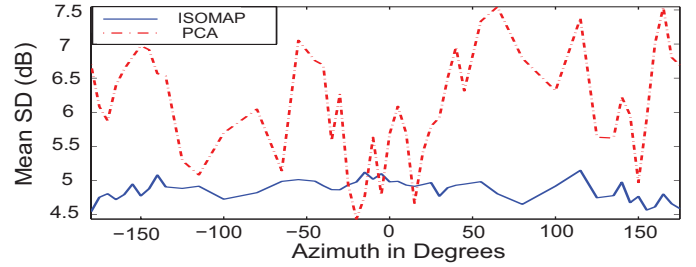


Fig. 7. Mean Spectral Distortion as a function of azimuth.

is gained from adding a dimension [27]. Using 0.025 for the threshold we find the intrinsic dimensionality  $d = 4$  [30].

Unlike previous works [21], we apply Isomap only once, over the entire dataset – a single procedure for the HRTFs of all subjects, ears and directions taking into account our proposed neighborhood selection. Figure 5 shows the Isomap manifold calculated for all directions and ears of 30 individuals (i.e.  $P = 30$ , so  $K = 2P + 1 = 61$  neighbors) from CIPIC database [14] in the horizontal plane, where the color represent the azimuth angle.

In Figure 5a, we plot the first embedded component of Isomap as a function of azimuth in order to highlight the symmetric properties of HRTFs. In Figure 5b and 5c, the manifold embedded in two and three dimensions show the variability of HRTFs across directions. Note that for each direction there are small clusters of reduced HRTFs. The variability inside these clusters is due to inter-subject differences (see Figure 6). Figure 5b illustrates that clusters are not uniformly distributed – the large gaps between some clusters is due to the HRTF non-uniform sampling in CIPIC database.

### B. Regression using an Artificial Neural Network

ANN is a system inspired by human brain capable of approximating nonlinear functions of their inputs. Since the relationship between HRTFs and anthropometric parameters is very complex, it is difficult to express them with linear functions. Here, we apply a back propagation ANN with sigmoid activation function in the hidden layer and a linear activation function in the output layer. The inputs are  $s$  anthropometric parameters, the azimuth angle in the horizontal plane and the ear (Left=1, Right=-1). The outputs are the coordinates of the HRTFs in the low-dimensional space obtained in Section IV-A. In order to determine the number of hidden nodes, we varied it from 5 to 30 and selected 20 hidden nodes that produced the lowest mean squared error. Note that our approach requires training only one ANN for all directions and ears. After the regression model is learned, the individual HRTF on the low-dimensional space for a new subject can be predicted by his anthropometric parameter measurements.

### C. Neighborhood Reconstruction Mapping

Unlike PCA and similar linear reduction methods, Isomap produce a low-dimensional embedding

$$\mathbf{Y}_{d \times N} = \{\mathbf{y}_1, \dots, \mathbf{y}_N\} \in \mathbb{R}^d$$

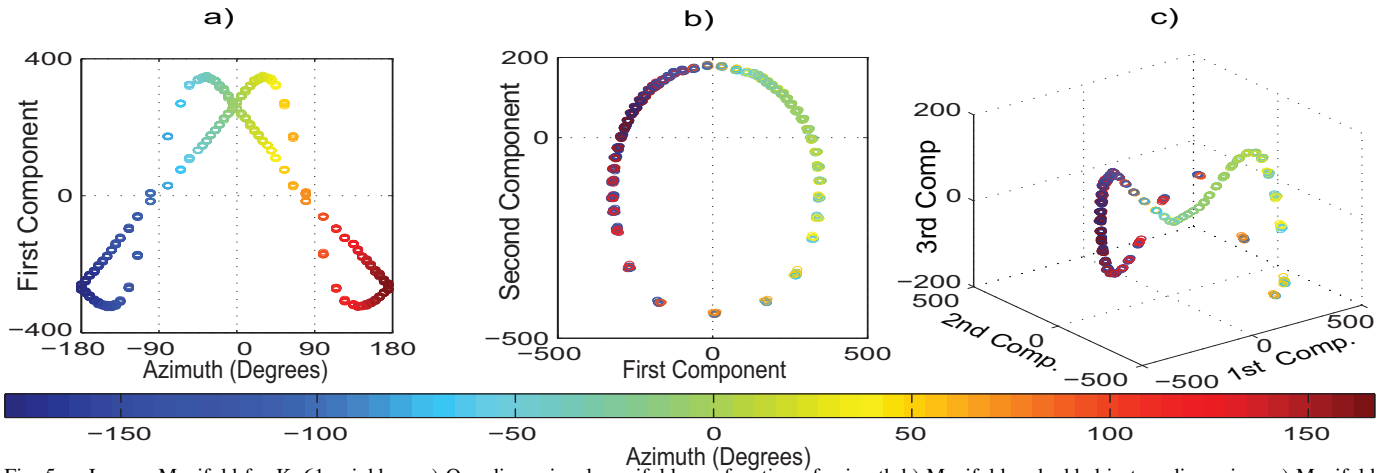


Fig. 5. Isomap Manifold for  $K=61$  neighbors a) One-dimensional manifold as a function of azimuth b) Manifold embedded in two dimensions c) Manifold embedded in three dimensions

from the samples in  $\mathbf{X}$  without generating an explicit map [27]. As we are interested in reconstructing an HRTF in the high-dimensional space from the low-dimensional HRTF predicted by the ANN, we need to project a low-dimensional point  $\mathbf{y}$  back into the original space. Since Isomap assumes that a sample and its neighbors are locally linear, we can perform the mapping using a linear combination of a sample's  $K$  neighbors [27], and the reconstructed HRTF,  $\hat{H}$ ,

$$\hat{H} = \sum_i^K w_i \mathbf{x}_i \quad (2)$$

to calculate the weights  $w_i$ , we follow Brown et al. [27], and choose  $w_i$  to be the inverse Euclidean distance between the sample and the neighbor  $i$  in the low-dimensional space.

## V. SIMULATIONS

We use the publicly available CIPIC database [14] which contains head related impulse responses (HRIRs) measured for 45 subjects at 1250 directions (25 azimuths and 50 elevations in interaural coordinate system). We employ 50 azimuth directions per subject and ear corresponding to horizontal plane. Each HRIR is 200 samples long (roughly 4.5 ms at 44.1 KHz sampling rate and 16 bit resolution). Each HRIR was transformed into an HRTF by a 512-point FFT. To reduce the effects of error due to nonlinearity introduced by equipments used to measure HRIRs, HRTFs were filtered to preserve frequencies between 200 Hz and 15 kHz, leaving 172 frequencies in each HRTF magnitude. We use only subjects that has the complete anthropometric parameters (i.e. 35 subjects). Performance was evaluated using a K-fold cross-validation approach. We split the HRTF dataset into 7 folds of 5 subjects each (6 folds for training and 1 fold for testing). Because the number of subjects for training each fold is  $P = 30$ , then according to our neighborhood selection proposed, the number of neighbors for Isomap is set to  $K = 2P + 1 = 61$

The CIPIC database also contains anthropometric measurements. We selected 8 anthropometric parameters for regression in accordance to [31]: head width, head depth, neck width, shoulder width, *cavum concha* height, *cavum concha* width,

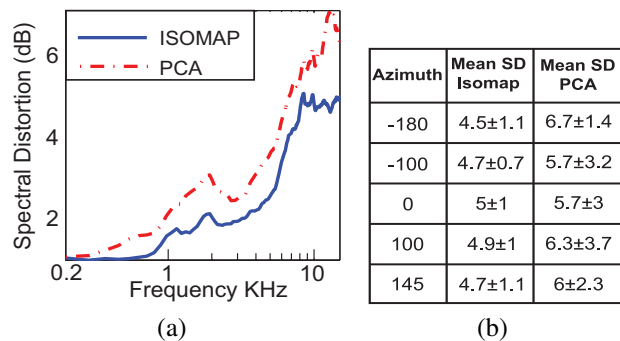


Fig. 8. a) Spectral Distortion b) Confidence Interval

pinna height, and pinna width. As explained in Section IV-B, the azimuth angle, the L/R ear, and the eight anthropometric parameters are the inputs for the ANN and the outputs are the low-dimensional HRTFs reduced using Isomap. We used Matlab Neural Network Toolbox 8.0.

We implemented a PCA-based customization, for comparison, with seven principal components (90% of variance). We used a similar ANN structure for the regression model and K-fold cross-validations for testing. We used Matlab Dimensionality Reduction Toolbox [30] for both PCA and Isomap.

We choose the mean spectral distortion as an error metric,

$$SD_M = \sqrt{\frac{1}{N_f} \sum_{f_k} \left( 20 \log_{10} \frac{|H(f_k)|}{|\hat{H}(f_k)|} \right)^2} \quad (3)$$

where  $H$  and  $\hat{H}$  represent the measured and reconstructed HRTF respectively and  $N_f$  is the number of frequency points. The reconstructed HRTF,  $\hat{H}$ , was calculated using Equation 2.

As can be seen in Figure 7, our approach performs better than PCA. The confidence interval ( $\pm 2\sigma$ , 95%) shows that our method has less variability than PCA (see Figure 8b). Moreover, our approach achieves better performance even with less dimensions than PCA. As in other studies [15], error increases at high frequencies due to complex scattering caused by pinna (Figure 8a) but in our approach it stays roughly below 5dB.

## VI. CONCLUSIONS

In this paper, we presented the problem of spatial audio and the need of HRTF personalization to ensure high 3D audio quality. Moreover, we introduced a new method for customizing HRTFs in the horizontal plane based on anthropometric measurements. Unlike previous works, we keep the multi-factor nature of HRTFs (i.e. frequency, direction and subject) by performing dimensionality reduction once on the entire HRTF dataset for all subjects, directions and ears in the horizontal plane. Besides using Isomap as a nonlinear dimensionality reduction technique, we introduce a brand-new graph construction technique that incorporates important prior information about the HRTFs that aims to exploit the correlations existent among HRTFs. The results show that incorporating prior knowledge in the neighborhood selection in Isomap can lead to a better manifold representation, and we can conclude that Isomap is a promising reduction technique for HRTFs analysis and synthesis.

We are currently extending our approach to estimate HRTFs beyond just the horizontal plane.

## PUBLICATIONS

The research presented here produced the following scientific paper and the first author was awarded a FAPESP M.Sc. grant 13/21349-1:

Felipe Grijalva, Luiz Martini, Siome Goldenstein, and Dinei Florencio. Anthropometric based customization of Head-Related Transfer Functions using Isomap in the horizontal plane. In *2014 IEEE Int. Conf. on Acoustics, Speech, and Signal Processing (ICASSP)*, Florence, Italy, May 2014.

## ACKNOWLEDGMENT

This work was supported by Microsoft-FAPESP grant 2012/50468-6, FAPESP M.Sc. grant 13/21349-1, FAPESP Ph.D. grant 14/14630-9, CNPq 308882/2013-0, CNPq: 454082/2014-2, and CAPES, and it is part of a project that was approved by Unicamp's IRB CAAE 15641313.7.0000.5404.

## REFERENCES

- [1] D. R. Begault, *3D Sound for Virtual Reality and Multimedia*. Cambridge: AP Professional, 1994.
- [2] A. Härmä, J. Jakka, M. Tikander, M. Karjalainen, T. Lokki, J. Hiipakka, and G. Lorho, "Augmented Reality Audio for Mobile and Wearable Appliances," *J. of the Audio Eng. Soc.*, vol. 52, no. 6, 2004.
- [3] R. Azuma, M. Daily, and J. Krozel, "Advanced human-computer interfaces for air traffic management and simulation," in *Flight Simulation Technologies Conf.* Reston, Virginia: American Institute of Aeronautics and Astronautics, Jul. 1996.
- [4] B. F. G. Katz, S. Kammoun, G. Parseihian, O. Gutierrez, A. Brilhault, M. Auvray, P. Truillet, M. Denis, S. Thorpe, and C. Jouffrais, "NAVIG: augmented reality guidance system for the visually impaired," *Virtual Reality*, vol. 16, no. 4, pp. 253–269, Jun. 2012.
- [5] S. Shoval, I. Ulrich, and J. Borenstein, "NavBelt and the Guide-Cane [obstacle-avoidance systems for the blind and visually impaired]," *Robotics & Automation Magazine, IEEE*, vol. 10, no. 1, 2003.
- [6] A. Lécuyer, P. Mobuchon, C. Mégard, J. Perret, C. Andriot, and J.-P. Colinot, "HOMERE: a multimodal system for visually impaired people to explore virtual environments," in *Virtual Reality, 2003*. IEEE Comput. Soc, 2003, pp. 251–258.
- [7] J. M. Loomis, R. G. Golledge, and R. L. Klatzky, "Navigation system for the blind: Auditory display modes and guidance," *Presence: Teleoperators and Virtual Environments*, vol. 7, no. 2, 1998.
- [8] H. Möller, "Fundamentals of binaural technology," *Applied Acoustics*, vol. 36, no. 3–4, pp. 171–218, 1992.
- [9] E. M. Wenzel, M. Arruda, D. J. Kistler, and F. L. Wightman, "Localization using nonindividualized head-related transfer functions," *J. of the Acoustical Soc. of America*, vol. 94, no. 1, pp. 111–123, Jul. 1993.
- [10] H. Fisher and S. Freedman, "The role of the pinna in auditory localization," *J. of Auditory research*, 1968.
- [11] P. Morse, *Theoretical acoustics*. McGraw-Hill, New York, USA, 1986.
- [12] V. R. Algazi, R. O. Duda, R. Duraiswami, N. A. Gumerov, and Z. Tang, "Approximating the head-related transfer function using simple geometric models of the head and torso," *J. of the Acoustical Soc. of America*, vol. 112, no. 5, p. 2053, Oct. 2002.
- [13] M. Otani and S. Ise, "Fast calculation system specialized for head-related transfer function based on boundary element method," *J. of the Acoustical Soc. of America*, vol. 119, no. 5, p. 2589, May 2006.
- [14] V. Algazi, R. Duda, D. Thompson, and C. Avendano, "The CIPIC HRTF database," in *Workshop on the Applications of Signal Processing to Audio and Acoustics*. New Platz, NY, USA: IEEE, 2001, pp. 99–102.
- [15] T. Nishino, K. Iida, N. Inoue, K. Takeda, and F. Itakura, "Estimation of HRTFs on the horizontal plane using physical features," *Applied Acoustics*, vol. 68, no. 8, pp. 897–908, 2007.
- [16] D. J. Kistler and F. L. Wightman, "A model of head-related transfer functions based on principal components analysis and minimum-phase reconstruction," *J. of the Acoust. Soc. of America*, vol. 91, no. 3, 1992.
- [17] G. Grindlay and M. A. O. Vasilescu, "A Multilinear (Tensor) Framework for HRTF Analysis and Synthesis," in *IEEE Int. Conf. on Acoust., Speech and Signal Processing (ICASSP)*, vol. 1, 2007.
- [18] L. Li and Q. Huang, "HRTF Personalization Modeling based on RBF Neural Network," in *IEEE Int. Conf. on Acoust., Speech and Signal Processing (ICASSP)*. IEEE, 2013, pp. 3707–3710.
- [19] Q. Huang and Y. Fang, "Modeling personalized head-related impulse response using support vector regression," *J. of Shanghai University*, vol. 13, no. 6, pp. 428–432, 2009.
- [20] H. Hu, L. Zhou, H. Ma, and Z. Wu, "HRTF personalization based on artificial neural network in individual virtual auditory space," *Applied Acoustics*, vol. 69, no. 2, pp. 163–172, Feb. 2008.
- [21] R. Duraiswami and V. Raykar, "The Manifolds of Spatial Hearing," in *IEEE Int. Conf. on Acoust., Speech and Signal Processing (ICASSP)*, vol. 3. IEEE, 2005, pp. 285–288.
- [22] B. Kapralos and N. Mekuz, "Application of dimensionality reduction techniques to HRTFs for interactive virtual environments," in *Int. Conf. on Advances in computer entertainment technology*. ACM, 2007.
- [23] B. Kapralos, N. Mekuz, A. Kopinska, and S. Khattak, "Dimensionality reduced HRTFs: a comparative study," in *Int. Conf. in Advances on Computer Entertainment Technology*. New York, New York, USA: ACM, Dec. 2008, p. 59.
- [24] H. Seung and D. Lee, "The manifold ways of perception," *Science*, vol. 290, pp. 2268–2269, 2000.
- [25] V. R. Algazi, C. Avendano, and R. O. Duda, "Estimation of a Spherical-Head Model from Anthropometry," *J. of the Audio Eng. Soc.*, vol. 49, no. 6, pp. 472–479, Jun. 2001.
- [26] J. B. Tenenbaum, V. de Silva, and J. C. Langford, "A global geometric framework for nonlinear dimensionality reduction," *Science*, vol. 290, pp. 2319–23, Dec. 2000.
- [27] W. M. Brown, S. Martin, S. N. Pollock, E. A. Coutsias, and J.-P. Watson, "Algorithmic dimensionality reduction for molecular structure analysis," *J. of chemical physics*, vol. 129, no. 6, p. 064118, Aug. 2008.
- [28] K. S. Lawrence and T. R. Sam, "Think globally, fit locally: Unsupervised learning of nonlinear manifolds," *J. of Machine Learning Research*, vol. 4, pp. 119–155, 2002.
- [29] B. Xie, X. Zhong, D. Rao, and Z. Liang, "Head-related transfer function database and its analyses," *Science in China Series G: Physics, Mechanics and Astronomy*, vol. 50, no. 3, pp. 267–280, Jun. 2007.
- [30] L. van der Maaten, E. Postma, and J. van den Herik, "Dimensionality reduction: A comparative review," *J. of Machine Learning Research*, vol. 10, pp. 1–41, 2009.
- [31] W. Wahab and D. Gunawan, "Enhanced Individualization of Head-Related Impulse Response Model in Horizontal Plane Based on Multiple Regression Analysis," in *Int. Conf. on Computer Eng. and Applications*, vol. 2. IEEE, 2010, pp. 226–230.