

# Personalization of Head-Related Transfer Functions with Limited Acoustic Measurements

**Griffin D. Romigh**

Dept. of Elec. & Comp. Eng.  
Carnegie Mellon University  
Pittsburgh, PA 15213  
*gromigh@ece.cmu.edu*

## Abstract

The main goal of this project was to take advantage of the spatial structure inherent within a set of HRTFs in order to provide a mechanism by which an entire set of individualized HRTFs could be estimated from a small set of measurement locations. Several methods were investigated for this purpose including missing feature techniques such as latent variable model decomposition and k-Nearest Neighbor averaging, as well as naïve methods such as spherical harmonic decomposition and linear interpolation. The most successful method investigated was a Linear Minimum Mean Square Error (LMMSE) estimation procedure which showed near complete reconstruction of the 232 HRTFs from a subset of less than 30 randomly distributed locations and promising results for as few as 20 locations. More practical measurement schemes like measurements only taken on the horizontal plane proved to be less successful than even distributions; however they still provided benefit over other rapid HRTF personalization techniques such as derivation from anthropometric measurements.

## 1 Introduction

Human listeners have the ability to determine the location of a sound source in three dimensions. This ability is largely due to robust cues based on the spatial separation of the two ears which causes a sound originating off to one side of the head to arrive at the near ear sooner and with greater amplitude than at the far ear. Thus creating an inter-aural time difference (ITD) and an inter-aural level difference (ILD). While the ITD and ILD dominate sound localization in most regards, these two cues alone fail to distinguish sound source positions which lie on contours of constant relative distance from the two ears; for instance a source directly in front of the listener and one directly above. In reality, these equal distance contours occur at every location in space and are referred to as cones of confusion. Sound source localization within one of these cones of confusion is accomplished using less robust spectral cues caused by the acoustic filtering properties of the listeners head, shoulders, and outer ear.

37 The concepts needed to accurately recreate this complete set of cues have been known for  
 38 some time, and involve the characterization of the acoustic filtering effects of the listeners'  
 39 anatomy. This characterization is accomplished by playing a known source from many locations in  
 40 space and making recordings at the entrance of the listeners' two ear canals. For a given location  
 41 in space the complex ratio of the Fourier transform of the recorded signal to that of the original  
 42 signal is called the Head Related Transfer Function (HRTF), and can be used to recreate the entire  
 43 set of cues needed for sound source localization. To adequately generate a 3-D representation  
 44 anywhere in space however, HRTFs for more than 250 locations in space must be collected, and  
 45 are generally only applicable for the person they were measured on. These two limitations make  
 46 attaining high fidelity spatial audio difficult.

47 Several authors have shown that once an HRTF set is collected, individual filters can be  
 48 represented well by weighted sums of a limited set of spectral basis vectors, and similarly that a  
 49 single frequency value can be represented anywhere in space with a weighted sum of spatial basis  
 50 vectors. These results prove that more compact and efficient representations of entire HRTF  
 51 databases are feasible, but currently these studies provide no advantage over traditional techniques  
 52 for acquiring the HRTF representation due to the fact that they are derived from HRTF sets  
 53 acquired using traditional methods. The goal of the current project is to use the knowledge gained  
 54 concerning low order representations of HRTF sets along with signal processing and machine  
 55 learning techniques to investigate methods for generating adequate representations of an entire set  
 56 of HRTFs from a limited set of easily attainable acoustic measurements.

## 57 2 Initial Method Evaluations

59 The HRTFs used in this study were gathered from the publically available CIPIC  
 60 database [1] which consists of 200 sample Head Related Impulse Responses (HRIRs)  
 61 captured at a sampling frequency of 44.1 kHz from 45 subjects at 1250 azimuth and  
 62 elevation locations. A 200 point DTF was taken of each HRIR and the magnitude was kept  
 63 as the reference HRTF signal. HRTF magnitudes for a small subset of 232 locations were  
 64 kept from the available 1250 which corresponded to roughly 15 degrees of angular  
 65 resolution in both azimuth and elevation.

66 For all of the techniques below a "one-out" training and testing procedure was used  
 67 and results were averaged over the same five subjects for each condition. Relative  
 68 performance was based on an **Average Spectral Distortion(SD)** measure presented in [3] as  
 69 defined below in equation 1.

$$70 \quad \overline{SD} = \frac{1}{L} \sum_{l=1}^L \sqrt{\frac{1}{N} \sum_{k=1}^N \left( 20 \log \left( \frac{|\hat{H}_k|}{|H_k|} \right) \right)^2} \quad (1)$$

71 Above,  $H$  and  $\hat{H}$  are the reference and modeled HRTFs at each of  $L$  locations respectively, all  
 72 consisting of  $N$  frequency bins. It represents the mean RMS error of the log-magnitude  
 73 spectra over an entire set of HRTFs.

### 74 2.1 Naïve Methods

76 In this context the term naïve is chosen to represent methods which do not require  
 77 training data. Their predictions are based entirely on the available HRTF measurements for  
 78 the test subject and can be thought of as a form of interpolation.

79 The simplest of these methods is a  $k$ -nearest neighbor linear interpolation, where  $k$ -  
 80 nearest neighbors refer to the  $k$  measured HRTFs whose locations are closest in absolute  
 81 angular distance to the location of the HRTF being predicted. The  $k$  HRTFs are then scaled  
 82 and summed to produce the predicted HRTF where the scaling factors are defined as the  
 83 ratio of the angular distance from that particular measurement location to the location being  
 84 predicted divided by the sum of all  $k$  angular distances as seen in equation 2.

$$85 \quad \hat{H}_m = \frac{1}{D} \sum_{i=1}^k d_i H_i \quad \text{where} \quad D = \sum_{i=1}^k d_i \quad (2)$$

86 The second naïve method which was investigated was spherical harmonic

87 decomposition. While the k-NN interpolation relied exclusively on local information,  
 88 spherical harmonic decomposition uses information from all known measurement locations  
 89 to determine the predicted HRTF. Spherical harmonics are continuous basis functions of  
 90 angular position, and can be used to estimate the underlying function at positions other than  
 91 those that were sampled. The HRTFs at the  $L$  measured locations are used to estimate  
 92 spherical harmonic coefficients,  $C_{pm}$ , which represent the underlying continuous HRTF  
 93 function in spherical harmonics as shown in equation 3[2].

$$94 \quad C_{pm}^f = \sum_{l=1}^L H^f(l) Y_{pm}^{*f}(l) \quad (3)$$

95 The spherical harmonic coefficients are then used to estimate the continuous HRTF  
 96 representation which can then be sampled at the missing locations as seen in equation 4.

$$97 \quad \hat{H}^f(l) = \sum_{p=0}^P \sum_{m=-p}^p C_{pm}^f Y_{pm}^f(l) \quad (4)$$

98 In this context  $P$  is the order of the spherical harmonic expansion and can be thought of as a  
 99 measure of spatial variation (i.e. a higher order is capable of capturing more rapid spatial  
 100 changes in the underlying function)[2].

## 101 102 **2.2 Missing Feature Methods**

103 In general, the term missing feature applies to any problem where portions of the  
 104 data are unknown or corrupted and need to be filled with predicted or average values. In this  
 105 way, an HRTF set lacking measurements at certain locations can be thought of as a missing  
 106 feature problem.

107 A basic solution for a missing feature problem is the standard k-Nearest Neighbor  
 108 (kNN) approach. As in the above naïve approach k of the “closest” known samples are  
 109 averaged to predict the unknown values, however in this implementation the HRTFs  
 110 available in the training set are evaluated and the k training HRTFs with the lowest spectral  
 111 distortion at a certain location are averaged to obtain the prediction. In this way equation 2  
 112 holds for this formulation as well with all  $d_i$ 's equal to  $1/k$  and the  $H_i$ 's are training HRTFs  
 113 for that location. One possible outcome of this method if k is chosen to be one and the same  
 114 training HRTF is chosen at every frequency and location corresponds to the scenario where a  
 115 complete HRTF from a single training subject is used for the missing locations. This  
 116 condition is reported as the “BestFit” training subject for use as a baseline performance  
 117 condition since it represents the best non-individualized HRTF set available.

118 A more recent approach to a missing feature problem was proposed by Smaragdis et  
 119 al. [4] who showed that a time sample of a spectrogram can be treated as a histogram  
 120 generated by repeatedly drawing from a mixture multinomial distribution with time  
 121 dependant mixture weights and source specific frequency multinomial bases. Since under  
 122 previous assumptions an HRTF is essentially a space indexed spectrogram, the HRTF of a  
 123 certain subject  $s$  at a given frequency  $f$  can be modeled as a histogram of  $N_d$  repeated draws  
 124 from a mixture multinomial distribution with *subject* dependant mixture weights and a set of  
 125 *spatial* multinomial bases as shown in equations 5 and 6.

$$126 \quad P_s^f(l) = \sum_{z=1}^Z P_s^f(z) P^f(l|z) \quad (5)$$

$$127 \quad \hat{H}_s^f(l) = N_d P_s^f(l) \quad (6)$$

128 The subject independent spatial bases can be trained using the EM algorithm from sample  
 129 HRTFs. The expectation step at each location consisted of computing the *a posteriori*  
 130 probability of the bases as in equation 7.

$$131 \quad P_s^f(z|l) = \frac{P_s^f(z) P^f(l|z)}{\sum_{z'=1}^Z P_s^f(z') P^f(l|z')} \quad (7)$$

132 The maximization step then consists of updating the bases and mixture weights as in  
 133 equations 8 & 9.

134

$$P_s^f(z) = \frac{\sum_{l=1}^L P_s^f(z|l)H_s^f(l)}{\sum_{l'=1}^L \sum_{l=1}^L P_s^f(z|l)H_s^f(l)} \quad (8)$$

135

$$P^f(l|z) = \frac{\sum_{s=1}^S P_s^f(z|l)H_s^f(l)}{\sum_{l'=1}^L \sum_{s=1}^S P_s^f(z|l')H_s^f(l')} \quad (9)$$

136

137

138

139

The unknown HRTFs for the test subject can then be estimated via the above EM algorithm where equation 9 is skipped and the bases used are those obtained from the previous training. In this instance the  $L$  locations being used for the update will only consist of the measured HRTF locations.

140

141

### 2.3 Minimum Mean Square Error Linear Estimation

142

143

144

145

146

147

While both forms of kNN described earlier are technically also linear estimations, a more common way of determining the scaling coefficients is to derive a closed form solution for them which minimizes some objective function. A frequent choice is the minimum mean square error solution which minimizes the average square difference between the actual and predicted values of the function being estimated. For known locations  $k$ , and missing locations  $m$ , the linear minimum mean square error solution is given in equation 10.

148

$$\tilde{H}_m = \mu_m + C_{mk} C_{kk}^{-1} (\tilde{H}_k - \mu_k) \quad (10)$$

149

150

151

Here  $C_{mk}$  is the cross covariance of the training HRTFs at a given frequency at missing and those at known locations;  $C_{kk}$  is the auto covariance at the known locations, and  $\mu_{(\cdot)}$  represent the corresponding means.

152

153

## 3 Initial Method Comparison

154

155

156

157

158

To save space, results for the variations of individual methods such as the various values of  $k$  in the  $k$ -NN approach are left out of the figures and discussions, and their best representative variation are used in all cross method comparisons. The final values for relevant methods and parameters are summarized in Table I.

TABLE I  
FINAL PARAMETER VALUES

Method	Parameter	Value
Naïve k-NN	$k$	4
Sph. Harm.	$P$	7
k-NN	$k$	3
LVM	$Z$	30

159

160

161

162

163

164

165

166

167

The average spectral distortions over all of the predicted HRTFs for 5 subjects are plotted in Figure 1 versus the number of measure locations used. Spherical harmonic decomposition (SHD) is the only method which failed to reach the baseline best fit training subject condition. This large amount of error is likely due to the fact that SHD has an inherent tradeoff between reconstruction accuracy and necessary spatial sampling rate. In other words higher order modeling (larger  $P$  parameter) will yield better overall reconstruction of a function, however, the higher the order of the model the more spatial samples are needed to accurately calculate the spherical harmonic coefficients.

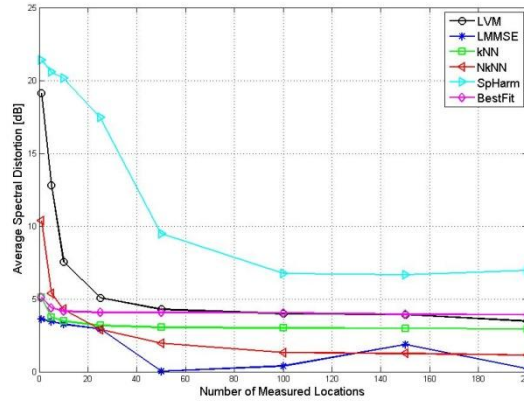


Fig. 1. Comparison of average spectral distortion for tested prediction methods as a function of the number of measurement locations.

168 K-NN proved to be one of the best performing methods for a low number of measurement  
 169 locations, but failed to improve for high numbers of measurements. LVM and Naïve kNN also  
 170 showed minimal improvement as the number of locations was increased, but Naïve kNN provided  
 171 fairly strong results for when greater than 100 measurement locations were used. Linear minimum  
 172 mean squared error estimation provided a clear improvement over all other methods and  
 173 approached zero spectral distortion when around 50 measurements were used.

174

#### 175 4 LMMSE Model Tuning

176 While LMMSE seemed to outperform other methods, it still suffered from discontinuities  
 177 in resulting spectra for low numbers of sources as can be seen in the top row of Figure 2 and some  
 178 degree of over fitting when more than 75 locations were used as shown by the small bump in  
 179 Figure 1.

180 To help avoid over fitting of the training HRTFs in the LMMSE model, a Forward Model  
 181 Selection algorithm was implemented based on the Akaike Information Criterion (AIC). AIC is a  
 182 cost function which seeks to find the minimize model error while simultaneously penalizing  
 183 models of increasing order (i.e. AIC will be lower for a model with less parameters if the two have  
 184 similar amounts of error). AIC is calculated as in Equation 11 and indicates a better model with  
 185 lower AIC scores [5].

$$186 \quad AIC = 2k + n[\ln(RSS)] \quad (11)$$

187 Above, *RSS* refers to the sum of squared residuals (errors) , *k* is the number of measured locations,  
 188 and *n* is the number of training samples.

189 The Forward Model Selection Algorithm is a technique used to find approximate best  
 190 models. Finding the best model with *N* measurement locations becomes intractable by brute force  
 191 methods after *N* is greater than 3 or 4, so the Forward Model Selection Algorithm approximates  
 192 this by finding the best *N* location model which includes the *N*-1 locations from the best *N*-1  
 193 location model. Using this strategy a good solution can be found in a far less number of  
 194 computations.

195 When the Forward Model Selection Algorithm was implemented using AIC as its  
 196 objective function for determining which of the *k* possible known locations should be used for  
 197 prediction, the large spectral discontinuities for low measurement numbers went away as can be  
 198 seen in the median plane plots in Figure 2. With this additional step, adequate prediction fell from  
 199 around 50 measurement locations to below 30.

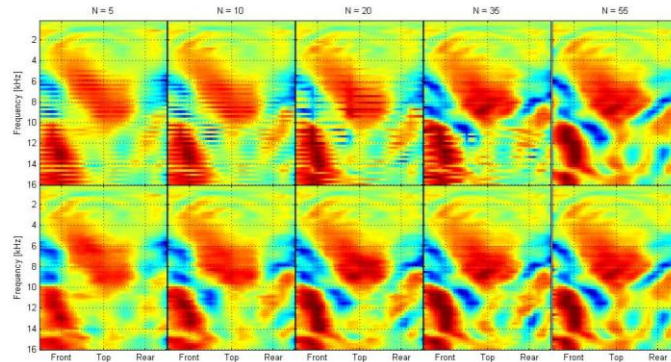


Fig. 2. Predicted HRTFs from the median plane using the LMMSE technique (TOP ROW) and LMMSE technique with forward model selection (BOTTOM ROW).

## 200 5 Practical Measurement Locations

201 While a general decrease in the number of measurement locations needed to accurately predict  
 202 an individualized HRTF is useful, the most useful reductions in measurement locations would be  
 203 ones which restricted measurement locations to those on single directional planes, such as the  
 204 horizontal plane, or hemi-planes, such as the front half of the horizontal and median planes.  
 205 Results for the LMMSE method for several of these setups are displayed in Figure 3. It can be  
 206 seen that in general roughly equally distributed measurement locations out perform the same  
 207 number of locations restricted plane locations, however the two conditions which feature sources  
 208 on the median plain fair better than their equally distributed counterparts.

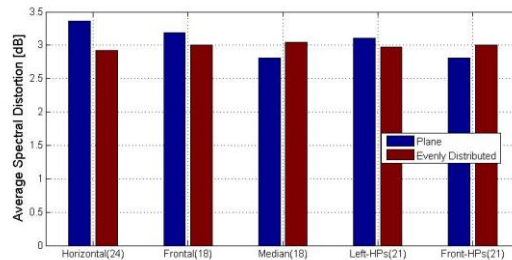


Fig. 3. Average spectral distortion for LMMSE method for practical measurement locations along principle planes and hemi-planes, number of measurement locations indicated in parentheses.

## 209 6 Conclusion

210 Several of the techniques developed show promise for providing personalization of a set of  
 211 HRTFs. More so than the rest, Linear Minimum Mean Square Error estimation proved to be very  
 212 reliable for setups including as little as 12% of the original measurement locations. Locations on  
 213 the median plane also show a higher than average performance versus more distributed locations,  
 214 which indicates highly practical implementations with low numbers of measurement locations are  
 215 also feasible.

## 216 References

- 217 [1] Algazi, V. R., and R. O. Duda. "THE CIPIC HRTF DATABASE." *IEEE Workshop on Applications of Signal*  
 218 *Processing to Audio and Acoustics* (2001).  
 219 [2] Evans, Micheal J. "Analyzing Head-related Transfer Function Measurements Using Surface Spherical  
 220 Harmonics." *J. Acoust. Soc. Am.* 104.2400 (1998).  
 221 [3] Grindlay, Graham, and M. A.O Vasilescu. "A MULTILINEAR APPROACH TO HRTF  
 222 PERSONALIZATION." *ICASSP* (2007).  
 223 [4] Smaragdis, Paris, Bhiksha Raj, and Madhusudana Shashanka. "Missing Data Imputation for Time-Frequency  
 224 Representations of Audio Signals." *JOURNAL OF SIGNAL PROCESSING SYSTEMS* (2010).  
 225 [5] H. Bozdogan, Akaike's Information Criterion and Recent Developments in Information Complexity, *Journal*  
 226 *of Mathematical Psychology*, 44 (2000).