

Estimation of HRTFs on the horizontal plane using physical features

Takanori Nishino ^{a,*}, Naoya Inoue ^b, Kazuya Takeda ^b,
Fumitada Itakura ^c

^a Center for Information Media Studies, Nagoya University, Japan

^b Graduate School of Information Science, Nagoya University, Japan

^c Faculty of Science and Technology, Meijo University, Japan

Received 19 December 2005; received in revised form 12 December 2006; accepted 26 December 2006

Available online 23 February 2007

Abstract

Sound localization can be controlled by using head related transfer functions (HRTFs), which are related to the size of the head, the ears and so on. Since HRTFs are characterized by source directions and subjects, it is necessary to conduct measurements in all directions for all subjects. However, such measurement is expensive and time-consuming. In this paper, we propose a simpler and more useful method that investigates the relationship between HRTFs and physical size by multiple regression analysis. The estimated HRTFs are evaluated by objective and subjective measures. For objective results, the average spectral distortion score is 4.0 dB in a bandwidth ranging from 0 to 8 kHz. Subjective results indicate no significant difference between the measured and the estimated HRTFs in that frequency range. These results support the hypothesis that the proposed method is effective for estimating HRTFs.

© 2007 Elsevier Ltd. All rights reserved.

PACS: 43.66.Dc; 43.66.Pn; 43.66.Qp

Keywords: Head related transfer function; Physical feature; Multiple regression analysis; Principal component analysis; Sound localization

* Corresponding author. Tel.: +81 52 789 4432; fax: +81 52 789 3172.

E-mail address: nishino@media.nagoya-u.ac.jp (T. Nishino).

1. Introduction

In headphone listening, head related transfer functions (HRTFs) are sometimes used to achieve spatial hearing. HRTFs make it possible to localize a sound image at an arbitrary point. An HRTF is an acoustical transfer function between the sound source and an arbitrary point in the ear canal. Since a head, the ears, and such physical parts refract, reflect, and block sound waves, an HRTF has very complex characteristics. The differences in the physical characteristics among listeners result in an HRTF depending on the listener and the direction of the sound source. Since accurate sound localization cannot be obtained with unsuitable HRTFs [1], it is very important to prepare suitable HRTFs. However, this is very difficult to do without measuring the HRTFs of all listeners and sound directions.

Several methods have been proposed to generate suitable HRTFs. Generating HRTFs by the boundary element method was investigated in [2,3]. In those studies, if the detailed shapes of the head and ears could be measured, in theory excellent HRTFs could be obtained. However, their proposed method requires a powerful computer with an enormous memory and a three-dimensional shape-measurement system.

The HRTFs produced by the method of clustering many listeners' HRTFs [4] provided good subjective performance; however, as yet there is no criterion for selecting suitable HRTFs for each listener.

We consider that estimating suitable HRTFs by a simple or an inexpensive measurement would be effective. In this paper, we propose, and objectively and subjectively evaluate, a simpler and more useful estimation method. Since HRTF contains reflections and refractions of sound waves at the ears and head, it can be considered that there is the relationship between HRTFs and such physical features as the size of the head, ears, and so on. In our method, HRTFs are estimated with physical features by multiple regression analysis. This proposed method can provide suitable HRTFs without special measuring equipment.

2. Measurement

HRTFs and physical features were measured for 86 subjects, 71 males and 15 females, who ranged from 17 to 33 years old.

2.1. Measurement of HRTFs

In our experiments, an HRTF is defined as the acoustical transfer function between a sound source and the entrance of the ear canal in the frequency domain. In the time domain, the impulse response between a sound source and the entrance of the ear canal is called a head related impulse response (HRIR).

HRIRs were measured in a reverberant room at reverberant times of 150 ms as shown in Fig. 1. Impulse responses were measured by using the time stretched pulse (TSP) [5] with the loudspeaker (BOSE ACOUSTIMASS) attached to the arched traverse. TSP duration was 1 s. Microphones (SONY ECM-77B) were placed at the entrances of both ear canals. In our measurements, the ear canals were not blocked completely by the microphones. The microphone diaphragms faced toward outside, thus putting the microphone diaphragms in the same position as the entrance of each ear canal.

The subjects were sat on a turntable that can be rotated at intervals of 1° with an accuracy of 0.3° .

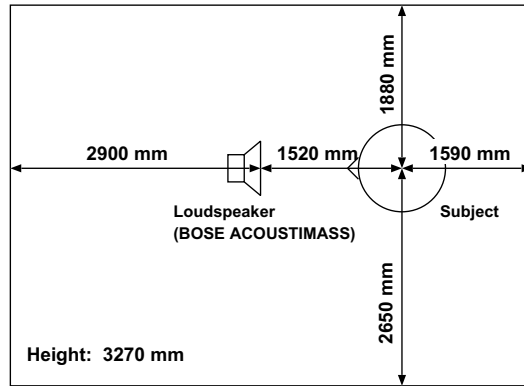


Fig. 1. Measurement environment.

HRIRs were measured for 72 azimuths at intervals of 5° . The distance from the sound source was 1520 mm, and the sampling frequency was 48 kHz. The azimuths corresponded to the following directions: front, 0° ; left, 90° ; back, 180° ; and right, 270° (Fig. 2). The sound source, the bitragion, and the pronasale were all located on a horizontal plane.

The HRIRs used in these experiments can be downloaded at the following Web site: <http://www.ciair.coe.nagoya-u.ac.jp/db/>.

2.2. Measurement of physical features

We measured nine physical features determined to be the data points when the KEMAR manikin was designed (Fig. 3) [6]: (1) ear length; (2) ear breadth; (3) concha length; (4) concha breadth; (5) protrusion; (6) bitragion diameter; (7) arc distance among the bitragion and the pronasale; (8) arc distance among the bitragion and the opistocranion; and (9) arc distance among the pronasale, the vertex, and the opistocranion. Table 1 shows the measurement results.

3. Method

The HRTFs were related to physical features by a multiple regression analysis [7]. Since an HRIR $h[t]$ is very complex, we divided each HRIR into two parts: an initial delay τ and a main response $h[t - \tau]$. Therefore, the estimation procedure is performed for the

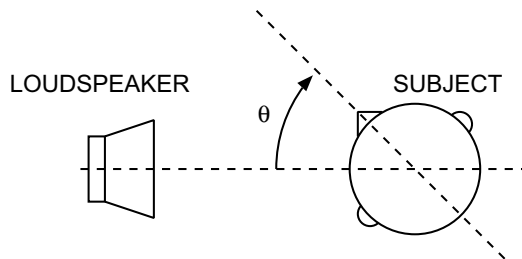


Fig. 2. Definition of angle.

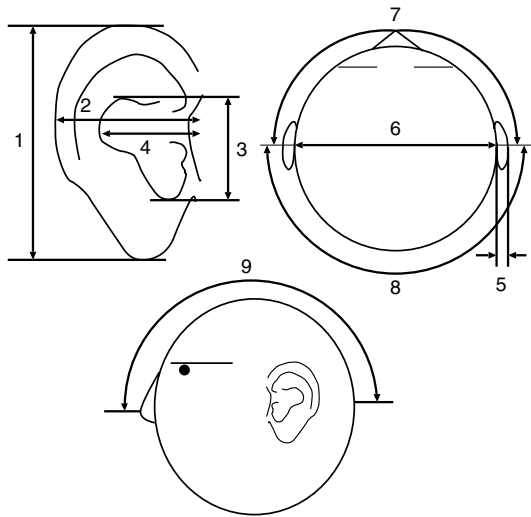


Fig. 3. Measured points of physical feature.

Table 1
Results of measuring physical features in mm

Measured point	Average	Maximum	Minimum	Standard deviation
1	65.6	82	55	4.9
2	32.0	40	22	3.4
3	18.4	30	10	3.5
4	18.4	24	13	2.2
5	22.1	30	15	3.4
6	148.0	181	113	13.4
7	307.3	371	270	19.1
8	227.6	275	195	16.3
9	424.9	482	360	21.1

magnitude response $|H(f)|$ of the main response and for the initial delay. Once multiple regression models have been calculated, only the physical features are measured to estimate the HRTFs. In other words, the listeners can obtain their own HRTFs by only measuring the physical features.

3.1. Estimation method for the magnitude response

The following describes the method for estimating the magnitude response.

First, a principal component analysis (PCA) is applied, which is represented as the product of the basis functions by the weights. Each HRTF magnitude response is subtracted from the mean at the angle θ

$$20\log_{10}|\widetilde{H}_{\theta,k}(f_i)| = 20\log_{10}|H_{\theta,k}(f_i)| - \frac{1}{M} \sum_{m=1}^M 20\log_{10}|H_{\theta,m}(f_i)|, \tag{1}$$

where $|H_{\theta,k}(f_i)|$ is the magnitude response of the k th subject's HRTF, and f_i is the frequency; $H_{\theta,k}(f_i)$ is calculated using the 512-point discrete Fourier transform.

The covariance matrix of $|\tilde{H}_{\theta,k}|$ is given by

$$s_{ij} = \frac{1}{M} \sum_{m=1}^M \{20 \log_{10} |\tilde{H}_{\theta,m}(f_i)| \times 20 \log_{10} |\tilde{H}_{\theta,m}(f_j)|\}, \quad (2)$$

where s_{ij} is an element of the covariance matrix \mathbf{S} ; i denotes the column index, and j the row. The covariance matrix \mathbf{S} is $I \times I$ matrix. For the analysis, variables i , j , f_i , and f_j are changed with each frequency.

The eigen-vector matrix \mathbf{C} is calculated using the eigen-vector of the covariance matrix \mathbf{S} . The eigen-vector matrix \mathbf{C} is also $I \times I$ matrix. The PCA weights $w_{\theta,k}$ are given by

$$\begin{aligned} w_{\theta,k} &= \mathbf{C}^{-1} \tilde{\mathbf{H}}_{\theta,k}, \\ w_{\theta,k} &= (w_{\theta,k}[1], \dots, w_{\theta,k}[N], \dots, w_{\theta,k}[I])^T, \\ \tilde{\mathbf{H}}_{\theta,k} &= (20 \log_{10} |\tilde{H}_{\theta,k}(f_1)|, \dots, 20 \log_{10} |\tilde{H}_{\theta,k}(f_I)|)^T, \end{aligned} \quad (3)$$

where T means transposition. In our experiments, the $w_{\theta,k}[n]$ ($n = 1, \dots, N, \dots, I$) are the criterion variables. N is the number of principle components for estimating magnitude responses.

Second, the physical features are corresponded to the weights by multiple regression analysis, and these weights can be estimated by a regression model. Here, x_{lk} ($l = 1, \dots, 9$) are the k th subject's physical features. The n th PCA weight at the angle θ is given by

$$w_{\theta,k}[n] = \alpha_{0n} + \sum_{l=1}^9 \alpha_{ln} x_{lk} + \varepsilon_n, \quad (4)$$

where α_{ln} is the multiple regression coefficient and ε_n denotes the estimation error. These parameters are calculated using every n .

To minimize ε_n by the method of least squares, the following calculations are performed:

$$E = \sum_{m=1}^M \varepsilon_n^2 = \sum_{m=1}^M \left(w_{\theta,m}[n] - \alpha_{0n} - \sum_{l=1}^9 \alpha_{ln} x_{lm} \right)^2. \quad (5)$$

If the result of partial differentiation is zero, the minimization of ε_n is achieved

$$\frac{\partial E}{\partial \alpha_{ln}} = 0 \quad (l = 0, \dots, 9). \quad (6)$$

Eq. (6) gives

$$\begin{pmatrix} \sum_{m=1}^M 1 & \sum_{m=1}^M x_{1m} & \cdots & \sum_{m=1}^M x_{9m} \\ \sum_{m=1}^M x_{1m} & \sum_{m=1}^M x_{1m}^2 & \cdots & \sum_{m=1}^M x_{1m} x_{9m} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{m=1}^M x_{9m} & \sum_{m=1}^M x_{1m} x_{9m} & \cdots & \sum_{m=1}^M x_{9m}^2 \end{pmatrix} \times \begin{pmatrix} \alpha_{0n} \\ \alpha_{1n} \\ \vdots \\ \alpha_{9n} \end{pmatrix} = \begin{pmatrix} \sum_{m=1}^M w_{\theta,m}[n] \\ \sum_{m=1}^M w_{\theta,m}[n] x_{1m} \\ \vdots \\ \sum_{m=1}^M w_{\theta,m}[n] x_{9m} \end{pmatrix}. \quad (7)$$

In Eq. (7), the following notations are applied:

$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \cdots & x_{91} \\ 1 & x_{12} & \cdots & x_{92} \\ \vdots & \vdots & \ddots & \vdots \\ 1 & x_{1M} & \cdots & x_{9M} \end{pmatrix},$$

$$\mathbf{W} = (w_{\theta,1}[n], w_{\theta,2}[n], \dots, w_{\theta,M}[n])^T,$$

$$\mathbf{A} = (\alpha_{0n}, \alpha_{1n}, \dots, \alpha_{9n})^T.$$

The resultant equation is

$$(\mathbf{X}^T \mathbf{X}) \mathbf{A} = \mathbf{X}^T \mathbf{W}. \quad (8)$$

Both sides of Eq. (8) multiplied by $(\mathbf{X}^T \mathbf{X})^{-1}$ give

$$\mathbf{A} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{W}. \quad (9)$$

The estimated weight $\hat{w}_{\theta,k}[n]$ is calculated using the multiple regression coefficients and the physical sizes

$$\hat{w}_{\theta,k}[n] = \alpha_{0n} + \sum_{l=1}^9 \alpha_{ln} x_{lk}. \quad (10)$$

Finally, estimated magnitude responses are obtained from the product of the basis function and the estimated weight:

$$\begin{aligned} \hat{\mathbf{H}}_{\theta,k} &= \mathbf{C} \hat{\mathbf{w}}_{\theta,k} + \bar{\mathbf{H}}_{\theta,k}, \\ \hat{\mathbf{H}}_{\theta,k} &= (20 \log_{10} |\hat{H}_{\theta,k}(f_1)|, \dots, 20 \log_{10} |\hat{H}_{\theta,k}(f_I)|)^T, \\ \hat{\mathbf{w}}_{\theta,k} &= (\hat{w}_{\theta,k}[1], \dots, \hat{w}_{\theta,k}[N], \dots, \hat{w}_{\theta,k}[I])^T, \\ \bar{\mathbf{H}}_{\theta,k} &= \left(\frac{1}{M} \sum_{m=1}^M 20 \log_{10} |H_{\theta,m}(f_1)|, \dots, \frac{1}{M} \sum_{m=1}^M 20 \log_{10} |H_{\theta,m}(f_I)| \right)^T, \end{aligned} \quad (11)$$

where T means transposition, and $\hat{w}_{\theta,k}[n] (n = N + 1 \dots I)$ are 0.

The estimated HRTF of the main response is reconstructed as a minimum phase [8].

3.2. Estimation method for initial delay

In the case of initial delay estimation, the individual initial delay is decided first. An initial delay τ , which is a part prior to the arrival time, does not change the magnitude response of the impulse response. It can be evaluated by calculating the correlation coefficient between the magnitude response of the measured HRIR $h[t]$ and the trimmed HRIR $h[t - \tau]$. The correlation coefficient r is defined with

$$\begin{aligned} r &= \frac{\text{Cov}(\mathbf{H}, \mathbf{H}_\tau)}{\sqrt{\text{Var}(\mathbf{H}) \text{Var}(\mathbf{H}_\tau)}}, \\ \mathbf{H} &= (|H(f_1)|, |H(f_2)|, \dots, |H(f_{257})|)^T, \\ \mathbf{H}_\tau &= (|H_\tau(f_1)|, |H_\tau(f_2)|, \dots, |H_\tau(f_{257})|)^T, \end{aligned} \quad (12)$$

where $|H(f)|$ is the magnitude response of the measured HRIR $h[t]$, and $|H_\tau(f)|$ is the magnitude response of the trimmed HRIR $h[t - \tau]$. The trimmed HRIR $h[t - \tau]$ is the same length as the measured HRIR $h[t]$ due to zero padding. Since the magnitude responses are calculated with the 512-point discrete Fourier transform, f_1 is 0 kHz and f_{257} is 24 kHz.

The correlation coefficient is 1 when τ is shorter than the correct initial delay. However, a smaller coefficient is obtained when the impulse response divides at an incorrect time [9]. In our method, the maximum τ that satisfies Eq. (13) is considered to be the initial delay [10]

$$r > 0.997. \quad (13)$$

The threshold 0.997 was determined experimentally. Since this threshold depends on the measurement condition, such as the background noise level, the measurement error and so on, the score should be changed every experiment.

These initial delays are corresponded to the physical features by multiple regression analysis for every azimuth:

$$\tau_\theta[k] = \beta_{0\theta} + \sum_{l=1}^9 \beta_{l\theta} x_{lk} + \varepsilon_\theta, \quad (14)$$

where $\beta_{l\theta}$ denotes the regression coefficient, x_{lk} is physical features, and ε_θ represents the estimation error decided by the method of least squares

$$\hat{\tau}_\theta[k] = \beta_{0\theta} + \sum_{l=1}^9 \beta_{l\theta} x_{lk}. \quad (15)$$

4. Experiments

The estimated HRTFs were evaluated both objectively and subjectively. Since Kistler and Wightman reported that an HRTF can be constructed with five principal components [8], we used five principal components related to the magnitude response in our experiments ($N = 5$). Furthermore, we used 82 subject's data for analysis and four subject's for estimation.

4.1. Experimental conditions

4.1.1. Objective tests

The estimation performance was evaluated using a spectral distortion (SD) score given by

$$\text{SD} = \sqrt{\frac{1}{I} \sum_{i=1}^I \left(20 \log_{10} \frac{|H(f_i)|}{|\hat{H}(f_i)|} \right)^2} \quad [\text{dB}], \quad (16)$$

where $|H(f_i)|$ denotes the magnitude response of the measured HRTF, $|\hat{H}(f_i)|$ is that of the estimated HRTF, and f_i is the frequency. The estimated HRTF is more similar to the measured HRTF when a small SD is obtained.

The errors of the initial delay were evaluated by

$$\text{error}_\theta = |\tau_\theta[k] - \hat{\tau}_\theta[k]|, \quad (17)$$

where $\tau_\theta[k]$ denotes the measured initial delay and $\hat{\tau}_\theta[k]$ is the estimated initial delay.

4.1.2. Subjective tests

We performed two subjective tests to evaluate the estimation. One of them was an evaluation of the measured interaural time delay and the estimated interaural time delay; the other, a sound localization test with the measured HRTFs and the estimated HRTFs [11,12].

In the evaluation of the initial delay, stimuli were white noise of 0.5 s duration, and the sampling frequency was 48 kHz. This test was performed with the method of paired comparison, where each subject answered a difference of two stimuli. One of stimuli had a measured interaural time delay, the other an estimated interaural time delay. Fig. 4 shows an example of the stimuli. There was a 0.5 s gap of silence between both stimuli. If subjects could not perceive the difference, the measured interaural time delay and the estimated interaural time delay were deemed to be the same. In this test, target azimuths were from 0° to 180° at intervals of 30° on the horizontal plane. The stimuli were transduced by headphones (ETYMOTIC RESEARCH, ER-4B), and each subject was presented with eight stimuli at each azimuth. The subjects were four males whose HRTFs were estimated by the multiple regression models.

In the sound localization test, stimuli comprised white noise of 1 s duration, and each subject was presented with four stimuli in each direction. In this test, target azimuths were from 0° to 330° at intervals of 30° on the horizontal plane. There was a 4 s gap of silence between each stimulus. The stimuli were transduced by headphones (ETYMOTIC RESEARCH, ER-4B). The subjects for these tests were four males whose HRTFs were estimated by the multiple regression models.

The subjects in both subjective tests were the same.

4.2. Results

4.2.1. Objective results

Fig. 5 shows the left ear SD scores by azimuth under two different bandwidth conditions. One bandwidth ranges from 0 to 8 kHz ($I = 86$), the other from 0 to 24 kHz ($I = 257$). Since SD scores were obtained by a cross-validation method, there were 83 experimental conditions. The average SD scores for all azimuths were 4.0 dB for the bandwidth from 0 to 8 kHz and 6.2 dB from 0 to 24 kHz. Low and high SD scores were respec-

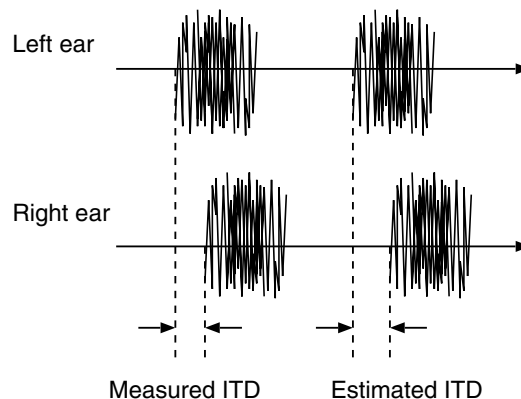


Fig. 4. Stimulus for the evaluation of the initial delay.

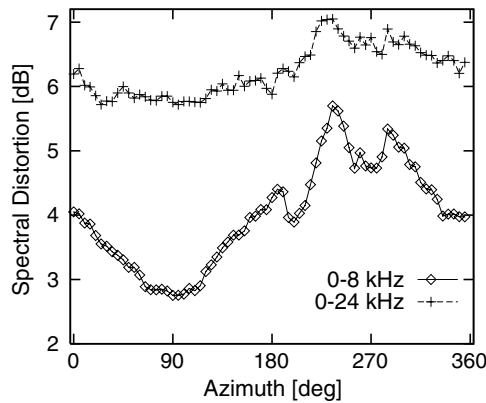


Fig. 5. Spectral distortion score of the estimated HRTFs (left ear). This figure was originally published in [12].

tively obtained at the angles close to the left and right ears. Since the HRTFs on the right side are formed by the reflection and refraction of sound waves, the HRTFs on the right side are more complicated than those on the left; therefore, the estimation performances are inferior. However, the influence on sound localization is slight because the sound pressure level is low at these angles.

It was found in results of previous research about interpolating the HRTF on the horizontal plane [10] that the rough sound localization was achieved when the SD score was 5.7 dB. Therefore, we assume that the performance of the estimated HRTFs will be the similar to the measured HRTFs.

The average error of the initial delay was 0.030 ms for the bandwidth from 0 to 24 kHz in the case of estimation using the physical features.

4.2.2. Subjective results

Fig. 6 shows the acceptance rate of the estimated interaural time delay. Because subjects answered whether there was the difference or not, the acceptance rate corresponds to the

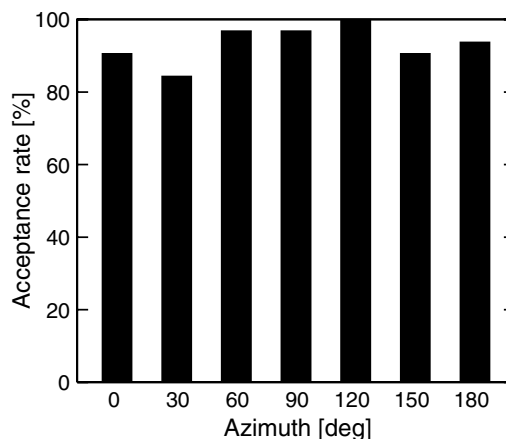


Fig. 6. Acceptance rate of the estimated ITD.

number of answers that subject could not perceived the difference. Significant tests (*T*-test) were performed for these results with the levels of significant α being 0.05. There was no significant difference for each azimuth except between 30° and 120°. The previous research [13] reported that there was a linear relationship between the lateral displacement and the interaural time delay. Their results suggested that the estimated interaural time delay could not cause the lateral displacement. Therefore, we considered that the estimated interaural time delay does not cause any lateralization error for the 30° intervals sound localization tests.

Fig. 7a–d shows the sound localization results. In these figures, circles correspond to the number of answers, which lie along the solid diagonal line. The front–back confusion, on the other hand, is represented on the dashed-and-dotted lines. Table 2 shows the correct and front–back confusion rates. In our experiment, a correct answer is one in which the presented direction and the perceived direction are the same. Front–back confusion is when the presented stimulus is perceived in a direction symmetrical to the bithrion diameter.

Significant tests (*T*-test) were performed for correct rates, with the levels of significant α being 0.05.

In the case of a comparison between the measured and the estimated HRTFs, there was no significant difference when the bandwidth ranged from 0 to 8 kHz, whereas there was

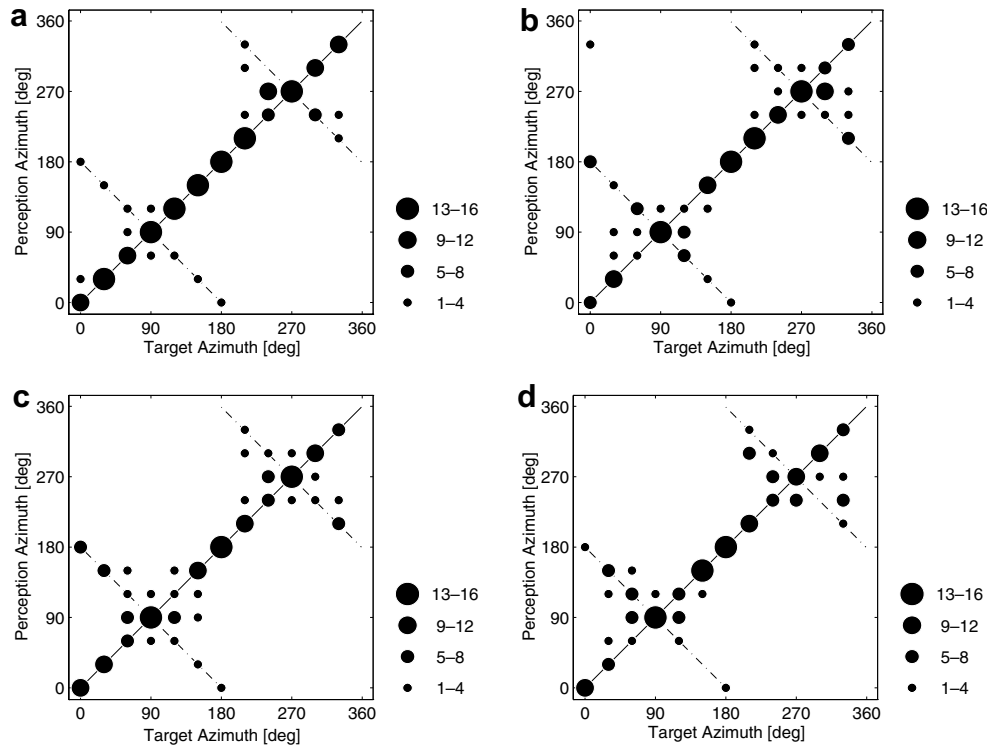


Fig. 7. Results of sound localization test using white noise. This figure was originally published in [12].

Table 2

The correct and front–back confusion rates of sound localization

Bandwidth	HRTFs	Correct [%]	Front/back [%]
0–24 kHz	Measured	77.1	11.5
	Estimated	59.9	17.7
0–8 kHz	Measured	60.4	19.8
	Estimated	58.9	15.1

one when the bandwidth ranged from 0 to 24 kHz. In the case of a comparison between the bandwidths, there was a significant difference in the measured HRTFs, while there was none for the estimated HRTFs.

These results suggest that HRTFs can be estimated using the physical features for a bandwidth from 0 to 8 kHz, and the frequency bands above 8 kHz also contribute sound localization. Therefore, it is necessary to estimate the HRTFs in the high frequency band to achieve excellent sound localization.

5. Conclusions

In this paper, we described a method of using physical features to estimate HRTFs. We evaluated the performance of estimated HRTFs with spectral distortion and sound localization, with results indicating that good performance was obtained with no significant difference between the measured and estimated HRTFs with respect to perception when the bandwidth ranged from 0 to 8 kHz. This supports the hypothesis that it is possible to estimate suitable HRTFs for all listeners by using physical features. Future work will involve improving the estimation method to attain good performance in the high frequency band.

References

- [1] Wenzel EM, Arruda M, Kistler DJ, Wightman FL. Localization using nonindividualized head-related transfer functions. *J Acoust Soc Am* 1993;94(1):111–23.
- [2] Otani M, Ise S. Numerical calculation of the head-related transfer functions by using the boundary element method. In: *Proc of WESTPRAC VII*, 2000. p. 305–8.
- [3] Katz B. Boundary element method calculation of individual head-related transfer function. I. Rigid model calculation. *J Acoust Soc Am* 2001;110(5):2440–8.
- [4] Shimada S, Hayashi N, Hayashi S. A clustering method for sound localization transfer functions. *J Audio Eng Soc* 1994;42(7/8):577–84.
- [5] Aoshima N. Computer-generated pulse signal applied for sound measurement. *J Acoust Soc Am* 1981;69(5):1484–8.
- [6] Burkhard MD, Sachs RM. Anthropometric manikin for acoustic research. *J Acoust Soc Am* 1975;58(1):214–22.
- [7] Nishino T, Nakai Y, Takeda K, Itakura F. Estimating head related transfer function using multiple regression analysis. *IEICE Trans A* 2001;J84-A(3):260–8.
- [8] Kistler DJ, Wightman FL. A model of head-related transfer functions based on principal components analysis and minimum-phase reconstruction. *J Acoust Soc Am* 1992;91(3):1637–47.
- [9] Jenkins GM, Watts DG. *Spectral analysis and its applications*. San Francisco: Holden-Day; 1968.
- [10] Nishino T, Kajita S, Takeda K, Itakura F. Interpolation of the head related transfer function on the horizontal plane. *J Acoust Soc Jpn* 1999;55(2):91–9.

- [11] Inoue N, Nishino T, Itou K, Takeda K. HRTF modeling using physical features. In: *Proc Forum Acusticum* 2005. p. L199–202.
- [12] Inoue N, Kimura T, Nishino T, Itou K, Takeda K. Evaluation of HRTFs estimated using physical features. *Acoust Sci Tech* 2005;26(5):453–5.
- [13] Toole FE, Sayers BM. Lateralization judgements and the nature of binaural acoustic images. *J Acoust Soc Am* 1965;37(2):319–24.