

# Robust 3D Robotic Sound Localization Using State-Space HRTF Inversion\*

Fakheredine Keyrouz and Klaus Diepold  
*Technische Universität München*  
80290, Germany  
{keyrouz,kldi}@tum.de

Patrick Dewilde  
*Technical University Delft*  
Delft, South Holland, Netherlands  
dewilde@dimes.tudelft.nl

**Abstract**— We address the problem of robotic real-time binaural hearing using a humanoid head. The sound detection ability for a previously proposed robotic 3D binaural sound localization technique is stabilized and considerably enhanced. The proposed technique is based on using two small microphones placed inside the ear canal of a robot dummy head mounted on a torso. As initially proposed, the detection algorithm is based on a simple correlation approach using a generic set of head related transfer functions, HRTFs, inverted using the fast fourier transform, FFT, method. In the current work, we replace the fourier inversion mechanism with state-space inversion using outer-inner factorization, and we demonstrate the performance through simulation and further in a household environment. Due to this modification, our robotic sound detection set up proves to be more noise-tolerant and able to localize sound sources in a three-dimensional space with a higher precision.

**Index Terms**— robotic binaural localization, HRTF, DFE, PCA, outer-inner factorization.

## I. INTRODUCTION

The human hearing organ has, throughout evolution, developed elaborate mechanisms for segregating and analyzing sound signals from separate sources in such a way that, under daily-life adverse conditions, robustness and low sensitivity with respect to varying external and internal working conditions are intuitively and automatically ensured. This hearing organ is often regarded as a signal conditioner and processor [1], and provides astonishing signal processing facilities. It comprises mechanic, acoustic, hydro-acoustic, and electric components, which, in total, realize a sensitive receiver and high-resolution spectral analyzer.

From the viewpoint of signal processing the underlying physical principles and a too-detailed description of a very complex system, like the ear organ, are of little interest and rather undesired, because computing times are dramatically increased. Many specialized cells in the auditory pathway contribute to the highly complex signal processing, which by far exceeds the performance of modern computers.

Eventually, auditory signal processing starts at the ear pinnae. The external sound field has to couple into the ear canals. The relative positions of the two ear canals and the sound sources lead to a coupling that is strongly dependent on

frequency. In this context, not only the two pinnae but also the whole head and torso have an important functional role, which is best described as a spatial filtering process. This linear filtering is normally quantified in terms of so-called head related transfer functions, HRTFs, or their corresponding head related impulse response, HRIR, which can also be interpreted as the directivity characteristics of the two ears. The interaural level and time differences, ITDs and ILDs, of the two signals entering the two ears are the basis of binaural hearing, which enables a fairly accurate spatial localization of sound sources. The limits to sound localization are determined by the limits of detecting ITDs and IIDs. The minimum audible angle, MAA, for sinusoidal sounds coming from directly in front of the subject is  $1^\circ$  for frequencies below 1000 Hz. Performance worsens around 1500-1800 Hz, due to fact that above 1500 Hz phase differences between the two ears are ambiguous cues for localization, and ILDs are small and do not change much with azimuth [2].

So far, a complete model, i.e. a model that is able to describe all types of binaural phenomena, does not exist. The problem is that neurophysiologists do not completely understand how living organisms localize sounds. The question of how ITDs and ILDs are combined, in the auditory system, to estimate the position of the sound source, has been only partly answered so far. Furthermore, the proposed models of the human ear organ, with all its underlying processing blocks and physical aspects, depend on the specific characteristics of the sources and the environments, and are therefore complex and hard to optimize.

Towards a fast robotic sound detection, using only two microphones, we have proposed a three-dimensional sound localization method [3] based on a very simple correlation algorithm which has access to a catalog of reduced KEMAR HRTFs. Inspired by the important role of the human pinnae to focus and amplify sound, this method demonstrates that a robot can perform sound localization in a three dimensional space using only two microphones, two artificial ears and a HRTF database. The method utilizes the effects of pinnae and torso on the original sound signal in order to localize the sound source in a simple matched filtering process. The 3D position of the source is determined by convolving the ear input signals, collected at the end of the robot ear canals, with the inverses of all possible HRTF pairs. The correct inverse

\*This work is fully supported by the German Research Foundation (DFG) within the collaborative research center SFB453 "High-Fidelity Telepresence and Teleaction".

yields maximum cross-correlation between the convolved left and right signals. Applying appropriate reduction techniques, the length of the impulse response of the HRTFs, was reduced to a hundred or even fewer samples, reducing thus the overall localization time and complexity.

The proposed approach for sound localization depends on finding the inverse filters of the whole reduced HRTF dataset, and saving for later use in localization. The inverse filter was directly made available by simply exchanging the values of the numerator and the denominator. However, getting the inverse filter is not as easy as it looks. The problem arises when the inverted filter is an unstable one. This is the case with the inverted HRTF filters, using the FFT method, especially that all HRTFs include a linear-phase component, i.e. pure delay, which is vital for maintaining the correct inter-aural time difference. The challenge is then to take the unstable filter and modify it to get a stable filter without disrupting either the magnitude or phase response or both. This can be difficult, if not impossible, to do while retaining causality. Most of the times, the resulting filter is non-causal but a stable and an excellent approximation of the original unstable filter. Non-causality is a small price to pay if magnitude and phase information are critical to performance which is the case in our situation. The following section will be mostly motivated by the inversion problem, and will cover important theoretical grounds on outer-inner factorization which forms the basis for the inversion algorithm we have adopted in this work.

## II. TRANSFER FUNCTION INVERSION

The stable inversion of transfer functions is of valuable importance for sound synthesis and channel equalization, not only to compensate from deficiencies of the transduction chain (amplifiers, loudspeakers, headphones), but also to reproduce a spatially coherent sound field. Direct methods to invert finite impulse response, FIR, transfer functions, like the head related transfer functions, HRTFs, may give undesired unstable results. A very efficient method, however, which handles this problem, and ensures stability, by simply translating the unstable inverse into an anti-causal yet bounded inverse, is the state-space inversion method using the celebrated outer-inner factorization. According to this method, the inner factor captures the part of the transfer function that causes the instability in the inverse, while the outer part can be straightforwardly inverted.

For rational time-invariant single-input single-output systems, the outer-inner factorization is a factorization of an analytical (causal) transfer function  $T(z)$  into the product of an inner and an outer transfer function,

$$T(z) = T_o(z)V(z). \quad (1)$$

The inner factor  $V(z)$  has its poles outside the unit disc and has modulus 1 on the unit circle, whereas the outer factor  $T_o(z)$  and its inverse are analytical in the open unit disc. Before deriving these inner and outer factors, we will recapitulate the steps involved in inverting a given transfer

function within the state-space domain, as it proves helpful in understanding the outer-inner factorization method.

Let  $T(z) = D + Bz(I - Az)^{-1}C$  be a minimal state-space representation of a given HRTF, with  $D$  being square and non-singular. This transfer function can also be represented by the state-space equations,

$$\begin{aligned} x_{k+1} &= x_k A + u_k B \\ y_k &= x_k C + u_k D \end{aligned} \quad (2)$$

The direct way of inverting the transfer function  $T(z)$  results in an unstable system, knowing that  $T(z)$  represents a non minimum-phase transfer function. This unstable inverse could be, in state-space, directly derived. We first take the second part of (2) and solve for  $u_k$ :

$$u_k = y_k D^{-1} - x_k C D^{-1}. \quad (3)$$

Inserting this in the first part of (2)

$$x_{k+1} = x_k A + B(y_k D^{-1} - x_k C D^{-1}), \quad (4)$$

leads to the inversion in the state-space, which can be written as

$$\begin{aligned} x_{k+1} &= (A - C D^{-1} B)x_k + D^{-1} B y_k \\ u_k &= -C D^{-1} x_k + D^{-1} y_k \end{aligned} \quad (5)$$

We denote the quadruple  $\{\tilde{A}, \tilde{B}, \tilde{C}, \tilde{D}\}$ , where  $\tilde{A} = A - C D^{-1} B$ ,  $\tilde{B} = D^{-1} B$ ,  $\tilde{C} = -C D^{-1}$ , and  $\tilde{D} = D^{-1}$ , as the state-space realization of the unstable inverse transfer function,  $T^{-1}(z)$ , corresponding to the inverse HRTF. The state-space realization  $\tilde{A}$  has its zeros outside the unit circle and, therefore, drives the system unstable. To ensure stability, we implement the outer-inner factorization theorem stated below.

### A. Inner-Outer Factorization

Given a minimal state-space realization of the transfer function  $T(z)$ , we would like to find the factors  $V(z)$  and  $T_o(z)$ , as in (1), where  $V(z)$  is unitary and  $T_o(z)$  is an outer function, that is to say it is minimum phase, and hence  $T_o^{-1}(z)$  is bounded but not causal.

Equation (1) can be expressed in a state-space form:

$$D + Bz(I - Az)^{-1}C = \quad (6)$$

$$[D + B_o z(I - Az)^{-1}C_o] [D_v + B_v z(I - A_v z)^{-1}C_v] \quad (7)$$

where  $\{A_v, B_v, C_v, D_v\}$  is a realization for  $V(z)$ , and  $\{A, B, C_o, D_o\}$  is a realization for  $T_o(z)$ .

Expansion of the quadratic term in (6) and equating members leads to

$$\begin{bmatrix} A & C \\ B & D \end{bmatrix} \begin{bmatrix} Y & 0 \\ 0 & I \end{bmatrix} = \begin{bmatrix} Y & C_o \\ 0 & D_o \end{bmatrix} \begin{bmatrix} A_v & C_v \\ B_v & D_v \end{bmatrix} \quad (8)$$

where the diagonal matrix  $Y$  satisfies the Lyapunov-Stein equation

$$Y = C C_v^* + A Y A_v^* \quad (9)$$

To get the inner and outer factors we are looking for, (8) must be solvable, and the kernel and maximality requirements on  $Y$  and  $D_o$  must indeed produce an outer factor  $T_o(z)$ .

*Theorem:* Let  $W$  be a unitary matrix, and  $Y$  be a uniformly bounded matrix which satisfies the following equality,

$$\begin{bmatrix} A & C \\ B & D \end{bmatrix} \begin{bmatrix} Y & 0 \\ 0 & I \end{bmatrix} = \begin{bmatrix} Y & C_0 \\ 0 & D_0 \end{bmatrix} W \quad (10)$$

such that  $Y$  has a maximal dimension and  $\ker(Y) = 0$ . Let

$$W = \begin{bmatrix} A_v & C_v \\ B_v & D_v \end{bmatrix}. \quad (11)$$

Then  $\{A_v, B_v, C_v, D_v\}$  is an isometric realization for the sought inner factor  $V(z)$ , and  $\{A, B, C_0, D_0\}$  is a realization for the outer factor  $T_o(z)$ .

The proof of the above theorem is detailed in [4]. We will recapitulate here its main steps.

Since the  $Y$  sought is such that  $\ker(Y) = 0$ , using  $RQ$ -factorization or SVD, we can always express it as

$$Y = V \begin{bmatrix} \sigma \\ 0 \end{bmatrix}, \quad (12)$$

in which  $\sigma$  is square non-singular and  $V$  is unitary.

We now let  $\Delta$  be a block Schur eigenspace decomposition for the term  $A - CD^{-1}B$  in (5),

$$\Delta = A - CD^{-1}B = V^* \begin{bmatrix} \delta_{11} & \delta_{21} \\ \delta_{21} & \delta_{22} \end{bmatrix} V, \quad (13)$$

where  $\delta_{11}$  collects the eigenvalues of  $\Delta$  which are strictly outside the unit circle, thus, causing instability.

Given the state-space realization  $\{A, B, C, D\}$ , and assuming that  $D$  is square invertible, we can write down the following schur factorization

$$\begin{bmatrix} A & C \\ B & D \end{bmatrix} = \begin{bmatrix} I & CD^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} \Delta & 0 \\ 0 & D \end{bmatrix} \begin{bmatrix} I & 0 \\ D^{-1}B & I \end{bmatrix}, \quad (14)$$

where  $\Delta = A - CD^{-1}B$  is the Schur complement of  $D$ . We can now write (10) as

$$\begin{bmatrix} \Delta & 0 \\ B & D \end{bmatrix} \begin{bmatrix} Y & 0 \\ 0 & I \end{bmatrix} W^* = \begin{bmatrix} I & CD^{-1} \\ 0 & I \end{bmatrix} \begin{bmatrix} Y & C_0 \\ 0 & D_0 \end{bmatrix}. \quad (15)$$

Taking the first block column of the right part of the equation, we find

$$\begin{bmatrix} \Delta & 0 \\ B & D \end{bmatrix} \begin{bmatrix} Y & 0 \\ 0 & I \end{bmatrix} \begin{bmatrix} W_{21}^* \\ W_{22}^* \end{bmatrix} = \begin{bmatrix} Y \\ 0 \end{bmatrix}, \quad (16)$$

where  $W_{21}^* = A_v^*$  and  $W_{22}^* = C_v^*$ .

From the last equation we can write,  $\Delta Y = YW_{21}^{-*}$ , or  $W_{21}^{-*} = Y^+ \Delta Y$ , where

$$Y^+ = [\sigma^{-1} \quad 0] V^*, \quad (17)$$

thus

$$W_{21}^{-*} = [\sigma^{-1} \quad 0] V^* \Delta V \begin{bmatrix} \sigma \\ 0 \end{bmatrix} \quad (18)$$

$$= \sigma^{-1} \delta_{11} \sigma. \quad (19)$$

It is important, at this point, to observe that since  $W_{21} = A_v = \sigma^* \delta_{11}^{-*} \sigma^{-*}$ , the matrix,  $\delta_{11}^{-*}$ , must have its eigenvalues

strictly inside the unit disc, achieving thus the sought system stability.

Multiplying (16) with  $Y^+$  we get

$$\begin{bmatrix} \sigma^{-1} \delta_{11} \sigma & 0 \\ BYD & 0 \end{bmatrix} \begin{bmatrix} W_{21}^* \\ W_{22}^* \end{bmatrix} = \begin{bmatrix} I \\ 0 \end{bmatrix}. \quad (20)$$

Introducing

$$\beta = D^{-1} B V \begin{bmatrix} \delta_{11}^{-1} \\ 0 \end{bmatrix},$$

we find

$$\begin{bmatrix} W_{21}^* \\ W_{22}^* \end{bmatrix} = \begin{bmatrix} \sigma^{-1} \delta_{11}^{-1} \sigma & 0 \\ -\beta \sigma D^{-1} & 0 \end{bmatrix} \begin{bmatrix} I \\ 0 \end{bmatrix}. \quad (21)$$

Finally, setting  $M = \sigma^{-*} \sigma^{-1}$ , we obtain the Lyapunov-Stein equation

$$M = \beta^* \beta + \delta_{11}^{-*} M \delta_{11}^{-1}. \quad (22)$$

Knowing that  $\delta_{11}^{-1}$  has its eigenvalues strictly inside the unit disc of the complex plane, (22) must have a unique solution, and both  $W_{21} = A_v$  and  $W_{22} = C_v$  have also unique solutions,

$$\begin{bmatrix} W_{21}^* & W_{22}^* \end{bmatrix} = [\sigma^* \delta_{11}^{-*} \sigma^{-*} \quad -\sigma^* \beta^*]. \quad (23)$$

Once we find both  $A_v$  and  $C_v$  realizations corresponding to the inner-factor  $V(z)$ , we can proceed to compute  $Y$  in (12) and use it in (10) to solve a system of four equations for the remaining inner-factor realizations  $C_v$  and  $D_v$  as well as the outer-factor,  $T_o(z)$ , realizations  $C_0$  and  $D_0$ . Finally, the sought stable inverse,  $T^{-1}(z) = T_o(z)^{-1} V^*(z)$ , is calculated, in which  $T_o(z)^{-1}$  is singular and causal, and  $V^*(z)$  is anticausal and stable.

It should be noted that the outer-inner factorization computed using the Lyapunov-Stein equation, ensures a linear solution adequate for a fast real-time implementation, compared with other methods [5] which solve the same outer-inner factorization problem quadratically using the celebrated Ricatti equation, and needlessly condition a problem which is already well-conditioned.

After we have recollected the essential theory we should use to invert the HRTFs, we will direct our attention to the novel robotic sound detection method, which will serve, in this case, as a test scenario for the above mentioned inversion technique.

### III. HRTFS REDUCTION TECHNIQUES

Our aim is to construct a three-dimensional binaural sound source localization system using only two microphones and a set of generic HRTF measurements. We use these measurements and develop a low-complexity model, based on simple correlation, for estimating the azimuth and the elevation for a sound wave impinging on a robot artificial pinna.

Experiments have shown that measured HRTFs can undergo a great deal of distortion (i.e. smoothing, reduction, etc.) and still be relatively effective at generating spatialized sound [1]. This implies that the reduced HRTF still contain

all the necessary descriptors of localization cues and is able to uniquely represent the transfer of sound from a particular point in the 3D space. We can take advantage of this fact, to greatly simplify the task of sound source localization, by using approximations of HRTFs, shortening the length of each impulse response and consequently reducing the overall localization processing time.

We use the Knowles Electronics Manikin for Acoustic Research, KEMAR, dataset which collects 710 impulse responses of measured HRTF filters. The 512 samples of each HRTF-measurement can directly be considered to be the coefficients of a FIR representation of the filter. However, for real-time processing, FIR filters of this order are computationally expensive. That's to say, using the 512 samples slows down the localization process and does not offer memory savings. Therefore, we investigate two techniques for reducing the length of the HRTF, namely, diffuse-field equalization, DFE, and principal component analysis, PCA. The original KEMAR HRTFs containing the 512 coefficients of the FIR filter shall be denoted by  $H_{512}^{FIR}$ . Using the reduced datasets, we presented in [3] a novel approach to localize sound sources using only two microphones in a real environment. This new method is considerably enhanced, using state-space inversion, as will be discussed later in the scope of the current paper.

#### A. Diffuse-Field Equalization

We aim to shorten the length of the original filters in order to reduce the computational burden for convolution, while preserving the main characteristics of the measured impulse responses. Towards this end, we adopt the algorithm proposed in [6] for a diffuse-field equalization (DFE). In DFE, a reference spectrum is derived by power-averaging all HRTFs from each ear and taking the square root of this average spectrum. Diffuse-field equalized HRTFs are obtained by de-convolving the original by the diffuse-field reference HRTF of that ear. This leads to the fact that the factors that are not incident-angle dependent, such as the ear canal resonance, are removed. The DFE is performed according to the following four steps: 1) Remove the initial time delay from the beginning of the measured impulse responses, which typically has a duration of about 10-15 samples. 2) Remove features from modeling that are independent of the incident angle [6]. 3) Smooth the magnitude response using a critical-band auditory smoothing technique [7]. 4) Construct a minimum-phase filter, ensuring thus stability for the final filter.

Thus, we shorten the length of the FIR representation of the original KEMAR HRTFs,  $H_{512}^{FIR}$ , from 512 to 128 coefficients. The resulting DFE HRTF database is denoted by  $H_{128}^{FIR}$ .

#### B. Principal Component Analysis

In order to examine to which extent the HRTF can be reduced while still preserving the characteristic information which makes it unique, we further reduce the previously

derived diffused-field HRTF dataset,  $H_{128}^{FIR}$ , by adopting Principal Component Analysis, PCA.

The first step in PCA is the computation of a frequency covariance matrix. These covariances provide a measure of similarity across the KEMAR's 712 HRTFs for each pair of frequencies. The covariance matrix  $S$  for a given pair  $(i, j)$  of frequencies is given by:

$$S_{i,j} = \frac{1}{N} \sum_k H_{k,i} \cdot H_{k,j}, \quad i, j = 1, 2, \dots, P; k = 1, 2, \dots, N \quad (24)$$

where  $N$  is the total number of HRTFs (712 in this case),  $P$  is the total number of frequency samples (512 in this case since HRTF equalization has been applied), and  $H_{k,i}$  is the magnitude at the  $i$ th frequency of the  $k$ th HRTF. The basis vector  $BF_q$ , is then derived from the  $q$ th eigenvector of the covariance matrix  $S$  that corresponds to the  $q$ th largest eigenvalue. The HRTF can then be modeled as a linear combination of several weighted basis functions. For a given HRTF, the weights representing the contribution of each basis function to that HRTF are given by:

$$w_i(\theta, \phi) = BF^T \times H_k(\theta, \phi) \quad i = 1, 2, \dots, 512 \quad (25)$$

Note that if the terms are rearranged, the HRTF magnitude vector is equal to a weighted sum of the basis vectors:

$$H_k(\theta, \phi) = BF \times w_i(\theta, \phi) \quad i = 1, 2, \dots, 512 \quad (26)$$

However, this equality holds if and only if  $q = p$ , or the maximum possible number of eigenvectors and basis vectors is retained. In practice,  $q \ll p$ , here we take  $q = 5, 10, 15, 20, 50$ .

Thinking about the accuracy, we found that it has reached 91.337% of the total variance by taking the first 20 basis functions only. Taking more basis functions can reduce the error between the measured and estimated HRTFs to some extent, but more calculation time and storing space are required. Once the reduced basis functions are selected, the weighting matrix is then calculated. These two matrices (correlation matrix and weighting matrix) are stored e.g. in a DSP memory, for a fast computation of reduced HRTFs. A thorough description of the PCA technique in modeling HRTFs is available in [8]. We shall denote the PCA-reduced HRTFs by  $H_m^{FIR}$ , where every HRTF has a length of  $m$  samples, and for every value of  $m$ , we have a truncated HRTF dataset.

### IV. ENHANCED SOUND SOURCE LOCALIZATION TECHNIQUE

We now recollect in detail the localization method which was suggested in [3] and adjust it in a such a way that robustness and better performance are achieved. The main idea in the general algorithm was to first minimize the HRTFs and remove all redundancy. The resulting minimized HRTFs are then used for localizing sound sources in the same way the full HRTFs would be used. The algorithm relies on a straight-forward matched filtering concept.



We assume that we have received the left and right signals of a sound source from a certain direction. The received signal by each ear is therefore the original signal filtered with the HRTF that corresponds to the given ear and direction.

Match Filtering the received signals through the correct HRTF should give back the original mono signal of the sound source. Although the system has no information about what the sound source is, the result of filtering the left received signal by the correct inverse left HRTF should be identical to the right received signal filtered by the correct inverse right HRTF.

In order to determine the direction from which the sound is arriving, the two signals must be filtered by the inverse of all of the HRTFs. The inverse HRTFs that result in a pair of signals that closely resemble each other should correspond to the direction of the sound source. This is determined using a simple correlation function. The direction of the sound source is assumed to be the HRTF pair with the highest correlation. This method is illustrated in Figure 1.

We have previously assumed that the reduced HRTFs, i.e. using DFE and PCA, have already been calculated and their fast fourier-transforms, FFTs, are saved. Once the audio samples are received to the left and right inputs, they must also be transformed using FFT. Then, the transformed left signal is divided (or multiplied by a pre-calculated inverse) by each of the left HRTFs. Similarly, the transformed right signal is divided by each of the right HRTFs. Finally, the correlation of each pair from the left and right is calculated. After the correlations are found, the direction that corresponds to the maximum correlation value is taken to be the direction from which the sound is arriving.

In the current work, instead of applying FFT to the reduced HRTFs and then inverting them, which is a crude unstable approximation of the real inverse, we directly invert them using the above mentioned state-space inversion method, and we save them for fast convolution with the incoming signals in the time-domain. The impact of this adjustment onto the existing system is considerable, and will be thoroughly discussed in the following simulation and experimental part.

## V. SIMULATION AND EXPERIMENTAL RESULTS

The simulation test consisted of having a broadband sound signal filtered out by the effect of the 512-sample HRTF at a certain azimuth and elevation. Thus, the test signal was virtually synthesized using the original HRTF set. For the test signal synthesis, a total of 100 random HRTFs were used. Since each of the HRTFs represents a unique 3D position around the robot's head, a total of 100 different random source locations in the 3D space were simulated. In order to insure rapid localization of multiple sources, small parts of the filtered left and right signal is considered (about 350msec). These left and right signal parts are then correlated with the available 710 reduced and state-space inverted HRTFs. Basically, the correlation should yield a maximum value when the saved HRTF ratio corresponds to the location from which the simulated sound source is originating. Therefore, we base

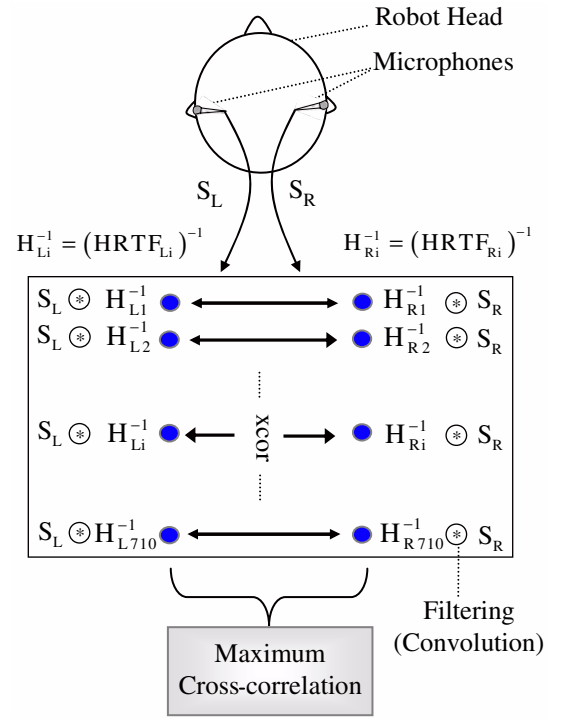


Fig. 1. Sound localization algorithm as proposed in [3]

our localization on the obtained maximum correlation factor. The reduction techniques, namely DFE and PCA, were used to create two different reduced models of the original HRTFs. The performance of each of these models under the state-space inversion method is illustrated in Figure 2. The circled line and the star sign in the Figure show the outer-inner factorization state-space inversion percentage of correct localization versus the length of the HRTF in samples. For comparison, the plus line and the multiplication sign in the figure refer to the previous FFT method performance [3]. Using the diffuse-field equalized 128 samples HRTF set,  $H_{128}^{FIR}$ , the state-space inversion simulated percentage of correct localization is 99%, this means that out of 100 locations, 99 were detected by applying the state-space inversion algorithm compared to 96% for the previous method using FFT inversion. Using the PCA-reduced set,  $H_m^{FIR}$ , the state-space inversion localization percentage falls between 55% to 97% compared to 42% to 91% for the previous FFT method, with the HRTF being within 10 to 45 samples, i.e.  $10 \leq m \leq 45$ . It should be noted that, while using 35 PCA-reduced HRTFs, all of the falsely localized angles fall within the close neighborhood of the simulated sound source locations.

In our household experimental setup, 100 binaural recordings from different directions were obtained using a dummy head and torso with two artificial ears in a reverberant room. The microphones were placed at a distance of 26 mm away from the ear's opening. The recorded sound signals, also containing external and electronic noise, were used as inputs to our state-space inversion algorithm. Using 35 samples of

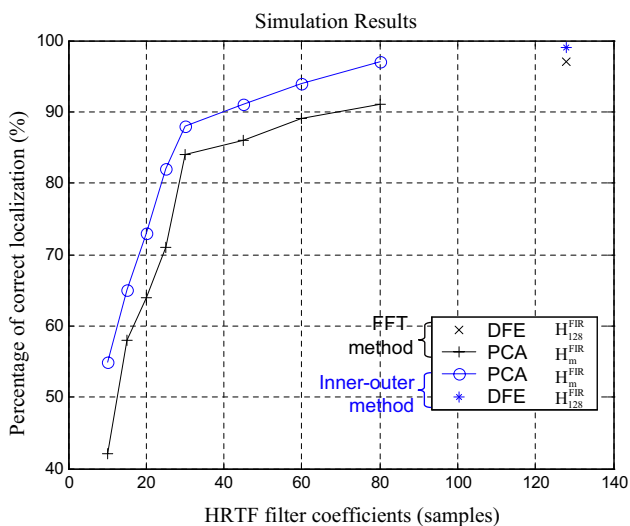


Fig. 2. Simulated percentage of correct localization using SCA, compared with the method in [3].

PCA-reduced HRTFs, 76% of the estimated azimuth and elevation angles turned out to be exactly at the intended correct location, compared with the 84% obtained from the theoretical simulation. The other falsely located 24% were found to be in the vicinity of the target angles. Due to external reverberation, and internal equipment noise, and due to the differences between the dummy manikin model used in the experiment and the KEMAR model used to obtain the HRTF dataset, 24% of the angles are found at the vicinity of the true location and not exactly at the target location where they originated from. More experimental results are shown in Figure 3.

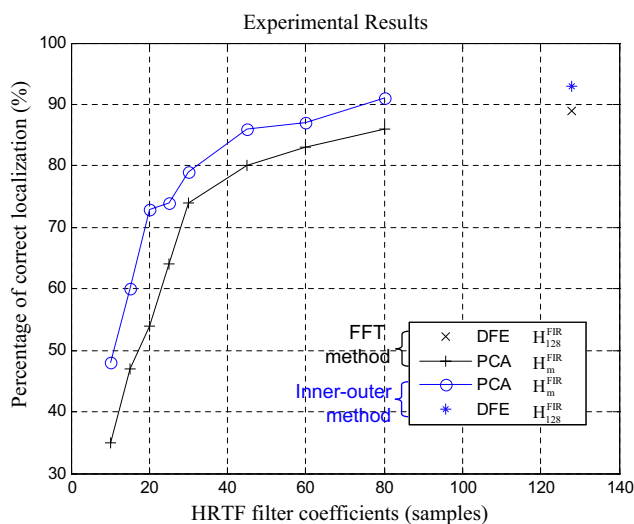


Fig. 3. Experimental percentage of correct localization using SCA, compared with the method in [3].

## VI. CONCLUSION AND FUTURE WORK

We have addressed the HRTF state-space inversion problem, using outer-inner factorization, and its application in binaural robotic sound source localization. Compared to the previous method using FFT inversion, the state-space inversion algorithm is able to increase the accuracy of the three-dimensional sound localization using only two microphones placed inside the ear canal of a robot dummy head mounted on a torso. In contrast to existing 3D sound source localization methods using microphone arrays, this two-microphone technique is based on a simple correlation approach which suggests a cost-effective implementation for robot platforms and allows for a fast localization of multiple sources. Using the presented scheme, two major venues for future work are to be considered, on one hand, we will extend the experimental setup to encompass robotic concurrent sound source localization, and, on the other hand, we will tackle the problems of headphone and loudspeaker equalization as well as crosstalk cancelation.

## REFERENCES

- [1] J. Blauert, "An introduction to binaural technology," in *Binaural and Spatial Hearing*, R. Gilkey, T. Anderson, Eds., Lawrence Erlbaum, USA-Hilldale NJ, 1997, pp. 593–609.
- [2] B. C. J. Moore, *An introduction to the psychology of hearing*. London, United Kingdom: Elsevier Press, 2004, pp. 233–240.
- [3] F. Keyrouz, Y. Naous, and K. Diepold, "A new method for binaural 3d localization based on hrtfs," in *proceedings of IEEE Intl. Conf. on Acoustics, Speech, and Signal Processing (ICASSP)*, Toulouse, France, May 2006, pp. 341–344.
- [4] P. Dewilde and A. J. van der Veen, "Inner-outer factorization and the inversion of locally finite systems of equations," *Linear Algebra and Its Applications*, vol. 313, pp. 53–100, 2000.
- [5] K. Yamada and K. Watanabe, "Inner-outer factorization for the discrete-time strictly proper systems," in *proceedings of IEEE 35th Conf. on Decision and Control*, 1996, pp. 1491–1492.
- [6] H. Moeller, "Fundamentals of binaural technology," *Appl. Acoust.*, vol. 36, no. 3-4, pp. 171–218, 1992.
- [7] J. Mackenzie, J. Huopaniemi, V. Vlimki, and I. Kale, "Low-order modeling of head-related transfer functions using balanced model truncation," *IEEE Singal Processing Letters*, vol. 4, no. 2, pp. 39–41, 1997.
- [8] D. Kistler and F. Wightman, "A model of head-related transfer functions based on principal components analysis and minimum-phase reconstruction," *J. Acoust. Soc. Amer.*, vol. 91, no. 3, pp. 1637–1647, 1992.