

Purely Auditory Online-Adaptation of Auditory-Motor Maps

Tobias Rodemann, Kalina Karova, Frank Joublin, Christian Goerick

2007

Preprint:

This is an accepted article published in IEEE-RSJ International Conference on Intelligent Robot and Systems (IROS 2007). The final authenticated version is available online at: https://doi.org/[DOI not available]

Purely Auditory Online-Adaptation of Auditory-Motor Maps

Tobias Rodemann, Kalina Karova, Frank Joublin, and Christian Goerick

Abstract—We present a system for an online-adaptation of auditory-motor maps that doesn't require a special set-up or dedicated robot movements and can therefore work during the normal operation of the robot. Our approach is based purely on auditory cues and motor position feedback for estimating the correct sound source position. The system can learn the correct auditory-motor map within 1–2 hours, starting from a random initialization, in a room with an active radio as the main sound source.

I. INTRODUCTION

In robotics there are a number of implicit or explicit sensory-motor maps which contain the information needed to map sensory inputs to motor outputs. One of these is the audio-motor map that relates sound localization cues to source position/bearing (or the motor command to focus on the sound source). Most sound localization systems use a set of two binaural cues for estimating the position of a sound source. The first cue is the Interaural Time Difference (ITD) which depends basically on the distance between the microphones and, for a humanoid robot with two ears, therefore on the size and form of the robot's head. The second cue is the Interaural Intensity Difference (IID), which depends on size, form, material and density of the robot's head. In the following we restrict ourselves to a binaural system, but the same approach could also be used for arrays of microphones. For more information see e.g. [1], [2], [3]. In order to estimate the sound source position from measured audio cues it is necessary to know the relation between cues and positions - the audio motor map. For ITD this relation is often, as in the case of most microphone arrays, relatively easy to compute. In other cases, e.g. when microphones are mounted on a humanoid robot's head, the relation becomes somewhat more complicated. The IID cue shows an even more complex dependence on sound source position. The IID-motor map also changes frequently, e.g. when hardware modifications are made. Therefore, when using IID or to improve ITD, the audio-motor map has to be relearned. This calibration process normally requires that a number of sound files is played from several position (ideally from different 3D positions, but normally only the azimuth position is varied). During this time the robot has to stand still and no other sound sources should be active. In realistic settings this will require at least one hour in a dedicated operation

mode. This calibration procedure would need to be repeated whenever any relevant aspect has changed. Among those aspects are microphones (type or position), recording hardware, robot body (especially the head), sound preprocessing, cue computation, and environmental conditions like echoes or background noise (if those can't be fully compensated by the preprocessing). Since it is difficult to estimate just how well the current map is adapted one either has to recalibrate frequently or risk a severe performance degradation due to badly tuned audio-motor maps. The obvious solution therefore is to do a continuous, unsupervised, online adaptation of the audio-motor map during normal operation of the robot so that the system always operates at or close to the optimum. A similar approach has been shown already for visual saccading systems [4]. In this paper we will describe a method for online adaptation of audio-motor maps using only auditory and motor information. Since audio localization cues (IID and ITD) are computed continuously by the sound localization system (e.g. [1], [5], [6]), the main issue in online adaptation is to get a good position estimation without relying on the potentially decalibrated auditory motor map. Nakashima et al. [7], [8] presented a system where vision is used to provide the position information. However, this approach requires a specific set-up, since it is a-priori very difficult to find the matching visual source to an audio source. In [7] the problem was solved by explicitly marking the audio source in the environment (a red marker on the loudspeaker) beforehand to ensure correct visual position estimation. Another approach was used by [3], where a face tracker was used to identify the position of a human speaker. This approach assumes that there is just one human in the environment and that all sounds are generated by this person.

It seems obvious to use vision feedback to provide the necessary position information for the online-adaptation. But in an uncontrolled real-world environment it is extremely difficult to know which visual target corresponds to which auditory one. Even with sophisticated audio and visual processing modules (providing access to e.g. speaker and object identity) it is hard to imagine how to find the correct visual target, especially given the limited viewing angle of most camera systems which means that most auditory sources will be outside the camera's field of view. Therefore we believe that at least for the bootstrapping of the audio-motor map, a purely auditory method for position estimation has to be used.

The ability to adapt sensory-motor maps is well-known in biology. A specific class of experience-dependent plasticity is known from experiments with barn owls (see e.g.[9]). Using ear plugs the relation between localization cues and sound

T. Rodemann, F. Joublin, C. and Goerick are with the Honda Research Institute Europe, Carl-Legien Strasse 30, 63073 Offenbach, Germany, Tobias.Rodemann@honda-ri.de

Kalina Karova is with the Institut für Regelungstheorie und Robotik, Technische Universität Darmstadt, Germany, kalina.karova@rtr.tu-darmstadt.de

position could be altered. Animals were capable of adapting to this new situation. It was recently found [10] that even mature ferrets are capable of readapting their audio-motor map after severe externally induced modifications and that this readaptation works even without vision, relying purely on auditory and motor inputs.

II. SYSTEM ARCHITECTURE

In this section we will outline the basic architecture of the sound processing system (see Fig. 1). We will focus on the major items only, for a more detailed description please consult [1], [2]. The hardware setup is a humanoid (styrofoam) head with silicon pinnae mounted on top of a pan/tilt element. As sound sources we used loudspeakers or normal radios. The latter setup is depicted in Fig. 2.



Fig. 1. The complete system graph. In addition to the normal sound localization architecture (as described in [1]), there are also three modules for online adaptation (grey box): one calculates the current localization error to update the adaptation width σ , the next one estimates the position of the sound source and the last one updates the audio-motor map which is used for the main sound localization system.



Fig. 2. The humanoid head mounted on a pan/tilt element in front of a radio used as a sound source.

After sound acquisition from a stereo microphone system we employ a Gammatone Filterbank (GFB) with Equivalent Rectangular Bandwidth (ERB) [11] to split the signal into separate frequency bands (frequency channels). We then extract a number of binaural cues: the Interaural Envelope Difference (IED), the Interaural Intensity Difference (IID), and Interaural Time Difference (ITD). Cues are measured only at signal onsets to reduce the effect of echoes [2]. IED and ITD are both based on the difference between consecutive zero-crossings [12] of the left and right microphone signal. While ITD acts directly on the GFB signal, IED operates on the signal envelope after an additional high-pass filtering. IID is computed as the difference between the left and right envelope signal divided by the maximum of left and right. In order to reduce stationary noise we employ binaurally synchronized spectral subtraction (see [1]). Our online test system uses 40 frequency channels between 50 and 4000 Hz (due to the limited range of the speaker in the radio), while the simulation test system uses 100 frequency channels between 100 Hz and 10 kHz. The choice of frequency channels is not relevant for the adaptation algorithm. During ego-motion of the head we interrupt sound localization due to motor noise.



Fig. 3. Example audio motor maps for ITD, IID, and IED (left column). In the right column single channels are visualized. The maps were generated via an offline calibration of one of our robot heads. The cue vector (box) represents the information that can be obtained from a single sound source position for an ideal (broadband) sound signal.

III. ONLINE ADAPTATION

In this section we will outline the basic principles of the online adaptation approach. We will first discuss how to estimate the position of the sound source and how to adapt the audio motor map and then describe some extensions and criteria for evaluating the quality of the audio-motor map.

In our system the audio-motor map contains the relations between audio (localization) cues like ITD or IID and sound source position. Calibration or adaptation is done by memorizing the cues measured at a known sound source position, while sound localization involves searching for the position in the audio-motor map whose cue vector (see Fig.3) best matches the currently measured cues. Since cue values depend on the frequency channel, there is one cue entry for every frequency channel and every position (for IED, IID, and ITD each). We represent this information in the cue matrix $C(f, \phi)$, where f is the channel index and ϕ the relative angle towards the sound source. Positions are sampled in the range between -90 and +90 degrees in steps of 10 degrees. Individual entries in the map are the mean over several measurements from the same position. Note that there is one $C(f, \phi)$ for IED, IID, and ITD each. For simplicity, however, we treat them as a single one.

A. Basic principles

Our central idea is to use a linear model for estimating the relation between sound source position and ITD. The linear model is chosen, because it is the simplest one with only two parameters that need to be estimated. Model parameters can be estimated by using ITD measurements from two different positions. In this section ITD measurements are averaged over a larger frequency range (between 100 and 650 Hz, since ITD is unambiguous in this range) to increase robustness. The ITD/IID-position relation is visualized in Fig. 4. The two graphs show ITD (top) and IID (bottom) mean values for different positions of the sound source (data taken from an offline calibration). The two graphs look very similar, but there is a crucial difference: ITD values are very close to zero for sound sources in front of the robot (0°) while IID values are normally not. The reason is that ITD basically depends on the distance between the sound source and the two ears. For a symmetric ear design any stimulus in front of the robot will have an equal distance to the two microphones. Therefore virtually all humanoid robots will have an ITD of approximately zero for sources at zero degrees azimuth. IID in contrast depends not only on the position of the microphones but also on the detailed characteristics of the microphones, the external and internal structure and material of the head, and amplification factors of the recording hardware.

We note another important aspect of the mean ITD: it is almost linear between -60 and +60 degrees. At more peripheral angles the curve deviates from the linear course, but in our experiments we found a generally good fit to the linear model.

Theoretically, the influence of source distance on ITD or IID should be minimal (at least beyond 1 m), however, we didn't perform extended tests, and under realistic conditions some performance reduction for sounds from different distances has to be expected.

Taking the observations that zero ITDs are reached at zero degree azimuth and making the linear approximation of the ITD to position ϕ relation:

$$\overline{ITD} = a \cdot (\phi - \phi_0), \tag{1}$$

with *a* as the slope and ϕ_0 as the sound source position, we just need two independent measurements of the sound source to estimate the true position. These two measurement (\overline{ITD}_1 and \overline{ITD}_2) are taken from two different positions (α_{M1} and α_{M2}) of the head (see also Fig. 5). We can then estimate *a* and ϕ_0 by:

$$a = \frac{\overline{ITD}_2 - \overline{ITD}_1}{\alpha_{M2} - \alpha_{M1}}$$
(2)

$$\phi_0 = \frac{\overline{ITD}_1 \cdot \alpha_{M2} - \overline{ITD}_2 \cdot \alpha_{M1}}{\overline{ITD}_1 - \overline{ITD}_2}.$$
 (3)

For this we make the assumption that when the head moves the active source before and after the movement is at the same position. This assumption will not always be true but in the majority of cases. We will later discuss how to optimize the chance of satisfying this criterion. Head movements can be arbitrary, either randomly as for our test scenario, or due to some other task like sound localization or visual exploration. For every pair of positions and assuming that audio cues could be computed before and after the move, the audio motor map can be updated for two different positions:

$$\phi_1 = \alpha_{M1} - \phi_0 \tag{4}$$

$$\phi_2 = \alpha_{M2} - \phi_0. \tag{5}$$

and the cue vector can be moved towards the measured values at positions ϕ_i (*i* = 1,2):

$$C^{k+1}(f,\phi_i) = C^{k+1}(f,\phi_i) + \beta \cdot (m_i(f) - C^k(f,\phi)).$$
(6)

Here, $m_i(f)$ is the IED, IID, or ITD value measured at relative position ϕ_i in frequency channel f. The learning parameter β determines the degree of adaptation for a single adaptation step k. A good choice is a value of $\beta = 0.3$.



Fig. 4. Mean cue values of ITD (top) and IID (bottom) for different sound source positions, integrated over different frequency channels. Note that for 0 degrees ITD is zero while IID is not.

B. Neighborhood adaptation

Learning can be improved considerably when updating not only those cue vectors for the current but also for neighboring positions. This is especially important when starting from a random initialization or after a stronger decalibration. We therefore update cue vectors for all positions $\phi_x \in [-90, 90]$ but modulate the step size by the distance from the measured position:



Fig. 5. Linear approximation model: average ITD as a function of relative sound source position and linear model approximation based on two ITD measurements at positions α_{M1} and α_{M2} . Based on these measurements the true sound source position ϕ_0 and the two relative positions ϕ_1 and ϕ_2 can be computed.

$$\Delta C^{k+1}(f,\phi_x) = \beta \cdot \exp\left(-\frac{(\phi_x - \phi_0)^2}{2\sigma^2}\right) \cdot \left(m_i(f) - C^k(f,\phi)\right).$$
(7)

The parameter σ defines the range of the adaptation neighborhood. The larger it is the quicker the map will adapt initially, but the lower the final quality of the map. It is therefore important to modify σ depending on the current performance of the audio motor map. Overall performance can be measured in many ways, we employ the following procedure: whenever the system tries to target an audio source, we measure the mean ITD after the move. In an ideal case it should be zero (assuming that the sound source hasn't moved). Any deviation from zero indicates an error in the sound localization, i.e. the audio motor map. Integrating these errors over many localization attempts, the system can estimate its own performance and update the neighborhood range σ . The ITD-based localization error E_{ITD} is calculated as the mean of measured ITD values in a time window of 6 seconds after a targeting move. Based on this error we can define an update rule for σ :

$$\sigma = \sigma_0 \cdot \frac{1}{1 + \exp\left(-s\left(E_{ITD} - t\right)\right)} \tag{8}$$

Here σ_0 is the maximum adaptation range, *s* the slope of the sigmoidal and *t* the ITD threshold. Example values are s = 20, $\sigma_0 = 80^\circ$, and t = 0.2. The advantage of this error function is that it can be computed online and does not require the knowledge of the true audio-motor map.

IV. RESULTS

The online adaptation was tested first in Matlab using simulated data and then in a real-world scenario. We investigated the performance when starting from a random audio-motor map in both cases. The online implementation runs in our own middleware system RTBOS [13] on a single CPU. The system graph is depicted in Fig. 1. To compare the results of online learning with the calibrated maps we introduce two new error measures:

Knowing the (correct) reference map C^{ref} the difference to the online adapted map after each updating step can be calculated using the following difference error:

$$E_{\text{diff}}^{k} = \frac{\sum_{f} \sum_{\phi} \left| C^{\text{ref}}(f,\phi) - C^{k}(f,\phi) \right|}{\sum_{f} \sum_{\phi} \left| C^{\text{ref}}(f,\phi) \right|},\tag{9}$$

with a sum over all positions ϕ , and frequency channels f. The difference error is normalized by the sum over all rectified entries in the offline calibrated reference map $C^{\text{ref.}}$

The localization error E_L indicates the mean deviation of the estimated position $\hat{\phi}$ from the true sound source position ϕ_0 . This estimation of the target sound source position ϕ is calculated using the main localization system (not the model based position estimation) and the current, learned audio-motor map. The localization error E_L is calculated as an average of the estimation error $\phi - \hat{\phi}$ for target sounds coming from all possible positions. In our simulations we considered only the front area and therefore $\phi_0 \in [-90^\circ; 90^\circ]$. This procedure can either be done using a simulation with an offline-calibrated map as the basis for cue generation or by testing the system online:

$$E_L = \frac{1}{N_{Positions}} \sum_{\phi_0 = -90}^{90} |\phi_0 - \hat{\phi}|.$$
 (10)

As a reference, offline calibrated maps in conjunction with our sound localization system [1] produce a mean localization error of around 2° under comparable conditions.

A. Offline tests

To study the system in a controlled environment we used Matlab to provide simulated inputs and head movements. Binaural cues were generated from an offline calibrated audio motor map with white noise added (with the same variance as those measured during calibration). Head movements were at random and binaural cues were produced for a fixed sound source position.

In the robot system the learning has to be precise and also fast. Figure 6 shows the learning progress using the linear model with parameters optimized for quick adaptation.



Fig. 6. Learning with parameters tuned for adaptation speed. Cues were simulated using an offline calibrated map.

With tuned parameters the adaptation speed can be increased and the map is learned within 100 steps (with a localization error $E_L = 4^\circ$). For this simulation the following

learning parameters were used: $\sigma_0 = 80^\circ$, s = 20 and t = 0.2. The learning step size is $\beta = 0.3$.

Good performance (less than 5° mean localization error) is reached after just 100 adaptation steps and an acceptable value already after 20 steps. Considering that one adaptation can be done in an online system roughly every 15 seconds (see below), the system could adapt from random initialization within about 300 seconds (5 minutes) under ideal conditions.

B. Online scenario

We tested online adaptation in a real-world scenario by setting our system in a normal (3x5m) room (echo decay constant $T_{60} = 625$ ms) and putting a conventional radio in front of the system (at 0 degrees). The audio-motor map was initially set to random for all cues. Then radio and online adaptation were turned on and the head was moved horizontally in a random fashion every 15 s. This long period guarantees that sound signals can be measured for extended periods of time before and after the head moves. In this scenario the radio was the main, but not the only, sound source - there was some background noise from the computing hardware and from neighboring office rooms. Also the echoes were considerable. From our experience we know that most of the background noise and echoes are handled by our preprocessing architecture (see [1], [2]). We let the system do online adaptation for several hours and then evaluated the results. We first analyzed the performance of the position estimation, see Fig. 7. As can be seen, in most cases the position of the sound source is well estimated (see also Table I). In almost 50% of trials the position estimation was correct within 10°. The few, really large errors could be selected out easily. The mean position estimation is not exactly on target (estimated mean position is -1.3° and the histogram peaks at -5°). About this we have to note that it was difficult to position the radio exactly at zero degrees relative to the robot head. Therefore the small localization error is probably due to the limited precision of the set-up. Furthermore, minor asymmetries in the head shape might also contribute. In any way, as we will see, this had no negative effect on the quality of the learned maps.

Next we monitored the status of the audio-motor map over time, see Figs. 8 and 9. Since the correct result is not known, we can only compare the online adapted map with one gained through an offline calibration in the same settings (which also took more than 1 hour of time). Fig. 8 shows the difference error E_{diff} for ITD and IID. As can be seen, the adaptation is finished after about 400 update steps (less than 2 hours of real time). Fig. 9 shows the development of the online calibrated map (IID) over learning steps (there is, on average, one learning step roughly every 15 s) and in the bottom right corner the offline calibrated map. Columns 1 and 3 show the IID-motor map, while columns 2 and 4 the value for a single frequency channel (1000 Hz). Dashed lines indicate the online adapted map, the solid lines the offline calibrated maps. The adaptation width σ is updated according to eqn. 8 with s = 20, t = 0.3,

 $\sigma_0 = 0.5$. The ITD error E_{ITD} is initialized with zero and rises slowly. Therefore, in the beginning adaptation is very slow but speeds up after a few hundred learning steps. Starting with a higher σ would speed up the learning accordingly. Here, we used an adaptation rate of $\beta = 0.1$.

Estimation	Result
estimation error within 1°	5%
estimation error within 5°	26%
estimation error within 10°	47%
estimation error within 20°	70%
estimation error over 30°	20%
mean estimation error	13.99°
mean estimated radio position	-1.3°

TABLE I

ESTIMATION PERFORMANCE FOR THE RADIO SCENARIO.



Fig. 7. Histogram of position estimation errors.



Fig. 8. Difference between offline-calibrated and online-adapted maps.

In a real-world environment we can't guarantee that the main constraints of the position estimation model are satisfied (source at same position before and after the move, only one sound source). In order to minimize adaptation errors we employ a number of checks to sort out potentially wrong position estimations: Firstly we measure localization



Fig. 9. Adaptation of (IID) online map over time. In columns 1 and 3 the IID map is depicted for different updating steps. Columns 2 and 4 depict the values for a single frequency channel (1 kHz), comparing calibrated data (solid line) with the online adapted values (dashed line).

cues over a period of 6 seconds before and after the moveignoring all cues measured before and after these periods and also during the head movement. This increases the chance that all measured cues are from the same sound source. Furthermore we only use position estimations for head movements beyond a certain threshold value (more than $\theta_m = 30^\circ$ head rotation between two measurements). In our tests this led to a substantial improvement in position estimation precision. Finally, by comparing the estimated slope *a* from eqn. 3 for the current position estimation with the mean estimated slope from previous measurements, we can sort out measurements with a considerable deviation from the mean. Obviously, we also sort out position estimations outside the possible range ($\pm 90^\circ$).

V. SUMMARY AND OUTLOOK

We have presented a method for online adaptation of audio-motor maps that relies solely on audio cues and motor feedback. We use a simple linear model of the mean ITD value. Although more complex models can be chosen, the position estimation performance with our approach was already satisfying. Our system is capable of self-adaptation in scenarios with comparatively few constraints on the environment and shows a robust and quick convergence of the audio-motor map. The approach is therefore well suited for applications in complex robotic systems, where frequent recalibration is not feasible. The adaptation can work while the robot is used for some other task, allowing a more efficient use of a generally very limited hardware resource (on-the-fly operation). Under ideal conditions adaptation can be done within 5-10 minutes and even under real-world conditions, the audio-motor maps can be learned in less than 2 hours. With some tuning even lower values will probably be possible. The number of required updating steps is relatively low, however it takes some time to robustly collect enough auditory cue measurements for an efficient training. The performance of the system could be extended in a number of ways, by e.g. using other auditory cues to verify that the sound source before the move is the same as the one after the move. This could be based for example on pitch tracking. It would also be possible to use an additional visual feedback signal for higher precision and increased robustness, by integrating the two position estimations and updating the map based on the combined results. Another possible option is to replace the linear model with a sinusoidal or sigmoidal function. This would provide a better fit to the observed ITD curve especially for more lateral positions, but would require to estimate more parameters.

REFERENCES

- T. Rodemann, M. Heckmann, B. Schölling, F. Joublin, and C. Goerick, "Real-time sound localization with a binaural head-system using a biologically-inspired cue-triple mapping," in *Proceedings of the International Conference on Intelligent Robots & Systems (IROS)*. IEEE, 2006.
- [2] M. Heckmann, T. Rodemann, B. Schölling, F. Joublin, and C. Goerick, "Binaural auditory inspired robust sound source localization in echoic and noisy environments," in *Proceedings of the International Conference on Intelligent Robots & Systems (IROS)*. IEEE, 2006.
- [3] J. Hörnstein, M. Lopes, J. Santos-Victor, and F. Lacerda, "Sound localization for humanoid robots - building audio-motor maps based on HRTF," in *Proceedings of the International Conference on Intelligent Robots & Systems (IROS)*. IEEE, 2006.
- [4] T. Rodemann, F. Joublin, and C. Goerick, "Continuous and robust saccade adaptation in a real-world environment," *KI-Künstliche Intelligenz*, vol. 03, pp. 23–26, 2006.
- [5] S. Yamamoto, K. Nakadai, M. Nakano, H. Tsujino, J.-M. Valin, K. Komatani, T. Ogata, and H. G. Okuno, "Real-time robot audition system that recognizes simultaneous speech in the real world," in *Proceedings of the International Conference on Intelligent Robots & Systems (IROS)*. IEEE, 2006.
- [6] K. Nakadai, H. Nakajima, M. Murase, H. Okuno, Y. Hasegawa, and H. Tsujino, "Real-time tracking of multiple sound sources by integration of in-room and robot-embedded microphone arrays," in *Proceedings of the International Conference on Intelligent Robots & Systems (IROS)*. IEEE, 2006.
- [7] H. Nakashima, N. Ohnishi, and T. Mukai, "Self-organization of a sound source localization robot by perceptual cycle," in 9th International Conference on Neural Information Processing (ICONIP'02), 2002.
- [8] H. Nakashima and T. Mukai, "3D sound source localization system based on learning of binaural hearing," in *IEEE International Conference on Systems, Man and Cybernetics*. IEEE SMC'05, October 2005.
- [9] E. I. Knudsen, P. F. Knudsen, and S. D. Esterly, "A critical period for the recovery of sound localization accuracy following monaural occlusion in the barn owl," *The Journal of Neuroscience*, vol. 4, pp. 1012–1020, 1984.
- [10] O. Kacelnik, F. Nodal, C. Parsons, and A. King, "Training-induced plasticity of auditory localization in adult mammals," *Public Library* of Science (PLoS) Biology, vol. 4, no. 4, pp. 627–638, 2006.
- [11] M. Slaney, "An efficient implementation of the patterson-holdsworth auditory filterbank,," Apple Computer Co., Technical Report 35, 1993.
- [12] Y.-I. Kim, S. J. An, R. M. Kil, and H.-M. Park, "Sound segregation based on binaural zero-crossings," in *Proc. Int. Conf. on Spoken Lang. Proc. (ICSLP)* 05, Lisboa, Portugal, 2005, pp. 2325–2328.
- [13] A. Ceravola, F. Joublin, M. Dunn, J. Eggert, M. Stein, and C. Goerick, "Integrated research and development environment for realtime distributed embodied intelligent systems," in *Proceedings of the International Conference on Intelligent Robots & Systems (IROS)*. IEEE, 2006.