

# Acoustic Models and Kalman Filtering Strategies for Active Binaural Sound Localization

Alban Portello, Patrick Danès and Sylvain Argentieri

**Abstract**—This paper deals with binaural sound localization. An active strategy is proposed, relying on a precise model of the dynamic changes induced by motion on the auditive perception. The proposed framework allows motions of both the sound source and the sensor. The resulting stochastic discrete-time model is then exploited together with Unscented Kalman filtering to provide range and azimuth estimation. Simulations and experiments show the effectiveness of the method.

## I. INTRODUCTION

Within the large field of robot perception, Robot Audition is recent in comparison to vision or even haptics. The related topics have been widely envisioned for the last ten years inside a new paradigm, including original and specific constraints raised by the robotics context. This has led to various works dealing with speech recognition [1], speaker recognition [2], and/or sound source localization [3][4]. Two main approaches can be cited. On the one hand, array processing methods can be adapted so as to obtain effective and robust auditive systems [5][6]. On the other hand, recent works have been more concerned with binaural audition, based on only two microphones. But using only two microphones to perform the aforementioned auditive tasks is very difficult, and lots of work are now focusing on performances improvement.

Most of the above approaches have only focused on the auditory scene analysis from a static view of the world. Such an idealized situation greatly eases the problem, while it is clear that speech and hearing takes place in a world where none of the static assumptions hold [7]. Yet, active perception, exploiting the possible movement of the auditive system, may help in the environment analysis process and in the auditive tasks robustness. This paper is more concerned with this last topic. The movement of the binaural sensor is fused with the auditive perception to perform sound source localization, in the case when the source and/or the sensor are in motion. Though the static case has been widely studied in the literature, to our opinion the moving case is still open. Nevertheless, recent contributions have already proposed solutions to active auditive perception. In [8], a sound source tracking system integrates audition, vision and motor movements together with an adaptive ego-noise cancelling algorithm. This work has then been extended to

active speech recognition in [9]. An active sound source localization scheme is also proposed in [10], where successive perceptions for various positions are exploited to train a Self-Organizing Map. The same approach is used in [11], relying on successive intersections of the cone of confusion. These two last contributions, though active, are fundamentally different from the one proposed in this paper, where the motion and the perception are fused to infer the sound location. This idea has been recently advanced in [12], where a particle filter is used for binaural tracking on the basis of interaural delay and motion parallax. In [13], an Extended Kalman Filter performs source tracking using interaural delay and instantaneous frequencies of perceived signals. Nevertheless, in both papers, the sensor position and velocity expressed in the world frame are assumed known, which implies to endow the sensor with an absolute localization system.

This paper aims at a generic framework to active binaural localization, *i.e.* fusing motion and perception. As modeling issues are often overlooked in the literature, a state space model is first proposed in Section II. The Kalman strategy underlying the estimation process is outlined in Section III. Auditive cues are then discussed in Section IV. The focus is mainly put on Interaural Time Difference (ITD), yet the results can be straightly extended to cope with Interaural Level Difference (ILD). Section V assesses the effectiveness of the proposed approach on realistic simulations and experiments. A conclusion and prospects end the paper.

## II. PROBLEM STATEMENT AND OBJECTIVES

One aim of this paper is to provide a correct mathematical framework to the binaural active localization of a moving sound source in a free-field environment. This section delves further into modeling issues. First, the problem is described in more accurate terms. Then, state space models are discussed, upon which Kalman filtering strategies can apply.

### A. Problem statement

A pointwise sound emitter  $E$  and a binaural sensor move independently on a common plane parallel to the ground. The two transducers equipping the sensor are denoted by  $R_l$  and  $R_r$ . A frame  $\mathcal{F}_R : (R, \mathbf{x}_R, \mathbf{y}_R, \mathbf{z}_R)$  is rigidly linked to the sensor, with  $R$  the midpoint of the line segment  $[R_l; R_r]$ ,  $\mathbf{y}_R$  the vector  $\frac{R R_l}{\|R R_l\|}$  and  $\mathbf{x}_R$  the downward vertical vector. The frame  $\mathcal{F}_E : (E, \mathbf{x}_O, \mathbf{y}_O, \mathbf{z}_O)$  attached to the source is parallel to the world reference frame  $\mathcal{F}_O : (O, \mathbf{x}_O, \mathbf{y}_O, \mathbf{z}_O)$ , with  $\mathbf{x}_O = \mathbf{x}_R$  (Figure 1).  $\|R_l R_r\| = 2a$  terms the transducers interspace. The source undergoes a translational motion

A. Portello and P. Danès are with CNRS; LAAS; 7 avenue du colonel Roche, F-31077 Toulouse Cedex 4, France and Université de Toulouse; UPS, INSA, INP, ISAE ; UT1, UTM, LAAS; F-31077 Toulouse Cedex 4, France [alban.portello@laas.fr](mailto:alban.portello@laas.fr), [patrick.danes@laas.fr](mailto:patrick.danes@laas.fr)

S. Argentieri is with UPMC Univ. Paris 06, 4 place Jussieu, F-75005, Paris, France and ISIR - CNRS UMR 7222, F-75005, Paris, France [sylvain.argentieri@upmc.fr](mailto:sylvain.argentieri@upmc.fr)

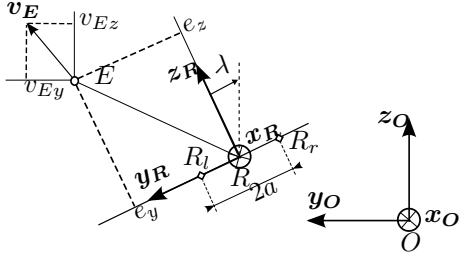


Fig. 1. The considered localization problem.

(velocities  $v_{Ey}, v_{Ez}$  of  $\mathcal{F}_E$  w.r.t.  $\mathcal{F}_O$  expressed along axes  $\mathbf{y}_O, \mathbf{z}_O$ ), while the sensor is endowed with two translational and one rotational degrees-of-freedom (velocities  $v_{Ry}, v_{Rz}$  of  $\mathcal{F}_R$  w.r.t.  $\mathcal{F}_O$  expressed along axes  $\mathbf{y}_R, \mathbf{z}_R$ ; rotation velocity  $\omega$  of  $\mathcal{F}_R$  w.r.t.  $\mathcal{F}_O$  around  $\mathbf{x}_O = \mathbf{x}_R$ ).  $v_{Ry}, v_{Rz}, \omega$ , and, to simplify,  $v_{Ey}, v_{Ez}$  are assumed known. The aim is to localize the emitter ( $\mathcal{F}_E$ ) w.r.t. the binaural sensor ( $\mathcal{F}_R$ ) on the basis of the sensed data at  $R_l, R_r$ . Free-field condition is assumed. Importantly, the audio sensor is not localized w.r.t.  $\mathcal{F}_O$ .

### B. Mathematical modeling

To tackle binaural active localization through a well-posed filtering problem, a state space model must be defined where the state vector is minimal. The state space equation, describing the way the source and sensor velocities affect the location variables, comes from rigid body kinematics. The output equation relates these spatial variables to acoustic cues extracted from the sensed data. This last topic will be shortly discussed here, then examined in more depth in Section IV.

The relative position and attitude of  $\mathcal{F}_E$  w.r.t.  $\mathcal{F}_R$  will henceforth be fully parameterized by means of a minimal set of three parameters  $e_y \triangleq \langle \mathbf{RE}, \mathbf{y}_R \rangle$ ,  $e_z \triangleq \langle \mathbf{RE}, \mathbf{z}_R \rangle$  (where  $\langle \cdot, \cdot \rangle$  stands for the scalar product), and  $\lambda \triangleq (\mathbf{z}_R, \mathbf{z}_O)_{\mathbf{x}_O}$ . When the sensor and source velocities are zero-order held at the sampling period  $T_s$ , the exact discrete-time state space equation is given by

$$\begin{aligned} \mathbf{x}_{[k+1]} &= \mathbf{F}\mathbf{x}_{[k]} + \mathbf{G}_1\mathbf{u}_1[k] + \mathbf{G}_2(\mathbf{x}_{[k]})\mathbf{u}_2[k] \quad (1) \\ \text{with } \mathbf{x}_{[k]} &\triangleq \begin{pmatrix} e_y[k] \\ e_z[k] \\ \lambda[k] \end{pmatrix}, \quad \mathbf{u}_1[k] \triangleq \begin{pmatrix} v_{Ry}[k] \\ v_{Rz}[k] \\ \omega[k] \end{pmatrix}, \quad \mathbf{u}_2[k] \triangleq \begin{pmatrix} v_{Ey}[k] \\ v_{Ez}[k] \end{pmatrix}, \\ \mathbf{F} &= \begin{pmatrix} c[k] & s[k] & 0 \\ -s[k] & c[k] & 0 \\ 0 & 0 & 1 \end{pmatrix}, \quad \mathbf{G}_1 = \begin{pmatrix} -\frac{s[k]}{T_s} & \frac{c[k]-1}{\omega[k]} & 0 \\ -\frac{c[k]-1}{\omega[k]} & -\frac{s[k]}{T_s} & 0 \\ 0 & 0 & -T_s \end{pmatrix}, \quad c[k] = \cos(\omega[k]T_s), \\ & \quad s[k] = \sin(\omega[k]T_s), \\ \mathbf{G}_2(\mathbf{x}_{[k]}) &= T_s \begin{pmatrix} \cos(\lambda[k] - \omega[k]T_s) & -\sin(\lambda[k] - \omega[k]T_s) \\ \sin(\lambda[k] - \omega[k]T_s) & \cos(\lambda[k] - \omega[k]T_s) \\ 0 & 0 \end{pmatrix}. \end{aligned}$$

To account for differences w.r.t. this deterministic model—drifts, slips, etc.—the stochastic state equation

$$\mathbf{X}_{[k+1]} = \mathbf{F}\mathbf{X}_{[k]} + \mathbf{G}_1\mathbf{u}_1[k] + \mathbf{G}_2(\mathbf{X}_{[k]})\mathbf{u}_2[k] + \mathbf{W}_{[k]}, \quad (2)$$

is defined by inserting a Gaussian random dynamic noise  $\mathbf{W}_{[k]}$  with known statistics. Recall that  $\mathbf{u}_1, \mathbf{u}_2$  are given. The measurement vector  $\mathbf{z}$  is made up with binaural cues, such as Interaural Time Difference (ITD) and Interaural Level Difference (ILD) [14][8]. It is viewed as a sample of the

measurement process  $\mathbf{Z}$ , linked to  $\mathbf{X}$  and the measurement noise process  $\mathbf{V}$  by an output equation of the form

$$\mathbf{Z}_{[k]} = \mathbf{h}(\mathbf{X}_{[k]}, \mathbf{u}_1[k], \mathbf{u}_2[k]) + \mathbf{V}_{[k]}. \quad (3)$$

A fundamental problem is often overlooked in the literature: due to the finite speed of sound, *the acoustic signal sensed by a microphone at time  $t$  depends on the signal emitted by the source at time  $t - \tau$ , with  $\tau$  the time delay due to propagation;  $\tau$  itself depends on the distance between the source at time  $t - \tau$  and the microphone at  $t$* . This fact precludes the writing of any output equation such as (2) when the source and/or the sensor are moving. Consequently, to base the localization strategy on a conventional state space model, as is required by any filtering scheme, simplifications are necessary. This will be the topic of Section IV. Before going into details, extensions to Kalman filtering are first briefly reviewed.

## III. ESTIMATION STRATEGIES

An original estimation scheme is presented, which can ensure a convenient localization with no prior knowledge.

### A. Extended and Unscented Kalman filtering

Consider a stochastic state space model such as (2)–(3). In the linear Gaussian case, the Kalman filter enables the recursive closed-form computation of the first two moments  $\hat{\mathbf{x}}_{[k|k]}, P_{[k|k]}$  and  $\hat{\mathbf{x}}_{[k|k-1]}, P_{[k|k-1]}$  of the probability density functions (pdfs)  $p(\mathbf{X}_{[k]} | \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]})$  and  $p(\mathbf{X}_{[k]} | \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]})$  for a given sequence of measurements  $\mathbf{z}_{[1:k]} = z_{[1]}, \dots, z_{[k]}$ . In the nonlinear case, the extended Kalman filter (EKF) propagates over time approximations relying on first-order Taylor expansions. Though widely used, it is often overconfident, *i.e.* it outputs too small covariances. Contrarily, the unscented Kalman filter (UKF) provides approximations of  $P_{[k|k]}, P_{[k|k-1]}$  up to an order of accuracy which can be analyzed theoretically, with no additional complexity [15]. For numerical efficiency and stability, the square-root UKF (SRUKF) has been implemented [16].

### B. A multiple hypothesis filter

Two difficult problems remain. The first one concerns the characterization of the initial state prior mean and covariance, to be used in the SRUKF initialization. Indeed, an initialization inconsistent with the values of the genuine hidden state can cause filter divergence. A tightly connected second problem may occur even if the SRUKF is well-initialized: the unscented transform of widely spread distributions may lead to overconfident conclusions. This partly comes from the difficulty of tuning the so-called “scaling parameter” of the UT.

To avoid these pitfalls,  $J$  independent—noninteracting—SRUKFs  $\{F_j\}_{j=1, \dots, J}$  are started in parallel, each one propagating the posterior moments  $\{\hat{\mathbf{x}}_{[k|k]}^j, P_{[k|k]}^j\}_{j=1, \dots, J}$ . Their initial conditions  $\{\hat{\mathbf{x}}_{[0|0]}^j, P_{[0|0]}^j\}_{j=1, \dots, J}$  are defined from a partition of the admissible state space—*i.e.* of the admissible relative sensor-source locations—into  $J$  overlapping cells  $\{C_j\}_{j=1, \dots, J}$ , *e.g.* so that each 99%

---

**Algorithm 1** The MH-SRUKF.

---

$[\hat{\mathbf{x}}_{[k|k]}, P_{[k|k]}, \{P(F_j|z_{1:k})\}_{j=1\dots J}, \{\hat{\mathbf{x}}_{[k|k]}^j, P_{[k|k]}^j\}_{j=1\dots J}] = \text{MH-SRUKF}(z_k, \{P(F_j|z_{1:k-1})\}_{j=1\dots J}, \{\hat{\mathbf{x}}_{[k-1|k-1]}^j, P_{[k-1|k-1]}^j\}_{j=1\dots J}, \mathbf{u}_{1[k]}, \mathbf{u}_{2[k]})$

- 1: **IF**  $k = 0$  **THEN**
- 2:   Define the initial conditions  $\{\hat{\mathbf{x}}_{[0|0]}^j, P_{[0|0]}^j\}_{j=1,\dots,J}$  and the weights  $\{W_0^j \triangleq P(F_j)\}_{j=1,\dots,J}$  of the filters  $\{F_j\}_{j=1,\dots,J}$ , with  $\sum_{j=1}^J P(F_j) = 1$ .
- 3: **END IF**
- 4: **IF**  $k \geq 1$  **THEN**
- 5:   **FOR**  $j = 1, \dots, J$  **DO**
- 6:     Predict the moments  $\hat{\mathbf{x}}_{[k|k-1]}^j, P_{[k|k-1]}^j$  from  $\hat{\mathbf{x}}_{[k-1|k-1]}^j, P_{[k-1|k-1]}^j$  through  $F_j$  according to (2) (SRUKF Time Update).
- 7:     Inside  $F_j$ , on the basis of (3) fuse  $\mathbf{z}_{[k]}$  with the moments  $\hat{\mathbf{x}}_{[k|k-1]}^j, P_{[k|k-1]}^j$  so as to get  $\hat{\mathbf{x}}_{[k|k]}^j, P_{[k|k]}^j$  (SRUKF Measurement Update).
- 8:     (Filter Likelihood) Apply the UT on  $\mathcal{N}(\mathbf{X}_{[k]}; \hat{\mathbf{x}}_{[k|k-1]}^j, P_{[k|k-1]}^j)$  through (3), and get the predicted mean and covariance  $\hat{\mathbf{z}}_{[k|k-1]}^j, S_{[k|k-1]}^j$  of the output. Then, set  $p(\mathbf{Z}_{[k]} = \mathbf{z}_{[k]} | F_j, \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]}) = \frac{1}{\sqrt{(2\pi)^{n_{\mathbf{x}}} \det(S_{[k|k-1]}^j)}} \exp(-\frac{1}{2}(\mathbf{z}_{[k]} - \hat{\mathbf{z}}_{[k|k-1]}^j)^T (S_{[k|k-1]}^j)^{-1} (\mathbf{z}_{[k]} - \hat{\mathbf{z}}_{[k|k-1]}^j))$ , with  $n_{\mathbf{x}} = \dim(\mathbf{x}$  or  $\mathbf{X}) = 3$  after (1)–(2).
- 9:   **END FOR**
- 10:   **FOR**  $j = 1, \dots, J$  **DO**
- 11:     (Filters posterior probabilities) Update  $W_k^j \triangleq P(F_j | \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]})$  of  $F_j$  according to (5).
- 12:   **END FOR**
- 13:   (Filters collapsing) If some  $W_k^j$  are lesser than a given threshold  $\gamma$ , then suppress the corresponding filters  $F_j$ . Decrease  $J$  and renormalize all the filters posterior probabilities accordingly.
- 14:   (Output processing) From the set of active filters, compute the overall posterior mean and covariance according to (7).
- 15: **END IF**

---

probability ellipsoid defined from the Gaussian prior  $p(\mathbf{X}_{[0]} | F_j) = \mathcal{N}(\mathbf{X}_0; \hat{\mathbf{x}}_{[0|0]}^j, P_{[0|0]}^j)$  covers  $\mathcal{C}_j$ . Each filter  $F_j$  is assigned a given initial probability  $W_0^j \triangleq P(F_j)$ . So, the initial state prior pdf is described by the Gaussian mixture

$$p(\mathbf{X}_{[0]}) = \sum_{j=1}^J P(F_j) \mathcal{N}(\mathbf{X}_{[0]}; \hat{\mathbf{x}}_{[0|0]}^j, P_{[0|0]}^j). \quad (4)$$

At each time  $k$ , after the measurement update, the posterior probability —or “weight”—  $W_{[k]}^j \triangleq P(F_j | \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]})$  of each filter  $F_j$  is computed from its likelihood  $p(\mathbf{Z}_{[k]} = \mathbf{z}_{[k]} | F_j, \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]})$  w.r.t. the measurement  $\mathbf{z}_{[k]}$  and from the weights at time  $k-1$   $\{P(F_j | \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]})\}_{j=1,\dots,J}$  through (see [17])

$$\begin{aligned}
 W_k^j &\triangleq P(F_j | \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]}) & (5) \\
 &= \frac{p(\mathbf{Z}_{[k]} = \mathbf{z}_{[k]} | F_j, \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]}) P(F_j | \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]})}{\sum_{l=1}^J p(\mathbf{Z}_{[k]} = \mathbf{z}_{[k]} | F_l, \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]}) P(F_l | \mathbf{Z}_{[1:k-1]} = \mathbf{z}_{[1:k-1]})} \\
 &= \frac{p(\mathbf{Z}_k = \mathbf{z}_k | F_j, \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]}) W_{[k-1]}^j}{\sum_{l=1}^J p(\mathbf{Z}_k = \mathbf{z}_k | F_l, \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]}) W_{[k-1]}^l}.
 \end{aligned}$$

Then, those filters whose weights fall below a given threshold  $\gamma$  (e.g.  $\gamma = 0.01$ ) collapse, and the weights of the remaining ones are renormalized. The posterior pdf then writes as

$$\begin{aligned}
 p(\mathbf{X}_{[k]} | \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]}) &= \sum_{j=1}^J P(F_j | k) p(\mathbf{X}_{[k]} | F_j, k) & (6) \\
 &= \sum_{j=1}^J P(F_j | k) \mathcal{N}(\mathbf{X}_{[k]}; \hat{\mathbf{x}}_{[k|k]}^j, P_{[k|k]}^j),
 \end{aligned}$$

where  $P/p(\cdot | k)$  are shortcuts for  $P/p(\cdot | \mathbf{Z}_{[1:k]} = \mathbf{z}_{[1:k]})$ . The overall posterior mean and covariance follow:

$$\begin{aligned}
 \hat{\mathbf{x}}_{[k|k]} &= \sum_{j=1}^J W_{[k]}^j \hat{\mathbf{x}}_{[k|k]}^j & (7) \\
 P_{[k|k]} &= \sum_{j=1}^J W_{[k]}^j (P_{[k|k]}^{(j)} + (\hat{\mathbf{x}}_{[k|k]}^j - \hat{\mathbf{x}}_{[k|k]})(\hat{\mathbf{x}}_{[k|k]}^j - \hat{\mathbf{x}}_{[k|k]})^T).
 \end{aligned}$$

The consequent Multiple Hypothesis SRUKF (MH-SRUKF) is summarized in Algorithm 1.

## IV. ACOUSTIC CUES

### A. Assumptions and General Equations

The air medium is assumed linear, and free of reflectors or scatterers. The signal  $s_T(t)$  received at time  $t$  by a transducer  $T \in \{R_l, R_r\}$  is a propagated and attenuated transform of the signal  $s_E(t)$  emitted by  $E$  according to

$$s_T(t) = \frac{1}{c\tau(t)} s_E(t - \tau(t)), \quad (8)$$

with  $c = 340\text{m}\cdot\text{s}^{-1}$  the speed of sound and

$$\tau(t) = \frac{1}{c} \|\mathbf{O}\mathbf{T}(t) - \mathbf{O}\mathbf{E}(t - \tau(t))\| \quad (9)$$

the propagation delay from  $E$  to  $T$ . Define  $\mathbf{v}_E$  and  $\mathbf{v}_T$  as the velocity vectors of  $E$  and  $T$  w.r.t. frame  $\mathcal{F}_O$ . Deriving  $(c\tau(t))^2$  w.r.t. time leads to, after some manipulations,

$$\dot{\tau}(t) = \frac{\langle (\mathbf{v}_E(t - \tau(t)) - \mathbf{v}_T(t)), \mathbf{n}_{\mathbf{T}(t) \rightarrow \mathbf{E}(t - \tau)} \rangle}{c + \langle \mathbf{v}_E(t - \tau(t)), \mathbf{n}_{\mathbf{T}(t) \rightarrow \mathbf{E}(t - \tau)} \rangle} \quad (10)$$

with

$$\mathbf{n}_{\mathbf{T}(t) \rightarrow \mathbf{E}(t - \tau)} \triangleq \frac{\mathbf{O}\mathbf{E}(t - \tau(t)) - \mathbf{O}\mathbf{T}(t)}{\|\mathbf{O}\mathbf{E}(t - \tau(t)) - \mathbf{O}\mathbf{T}(t)\|}. \quad (11)$$

The celebrated formula of the Doppler shift straightly follows from (10). Indeed, defining the instantaneous frequency  $f$  of a signal as  $2\pi$  times the derivative of its phase, one gets

$$f_T(t) = f_E(t - \tau(t)) \left( \frac{c + \langle \mathbf{v}_T(t), \mathbf{n}_{\mathbf{T}(t) \rightarrow \mathbf{E}(t - \tau)} \rangle}{c + \langle \mathbf{v}_E(t - \tau(t)), \mathbf{n}_{\mathbf{T}(t) \rightarrow \mathbf{E}(t - \tau)} \rangle} \right). \quad (12)$$

Eqs. (8)–(9)–(10)–(12) are graphically represented on Figure 2. As they link the characteristics of the emitted and received signals to their spatial positions and velocities, they are important for sound source localization. Nevertheless, as outlined before, they cannot be used as they are in the output equation of the state space model for they do not express a memoryless mapping between the output and the state/control vectors. The forthcoming section discusses some approximations.

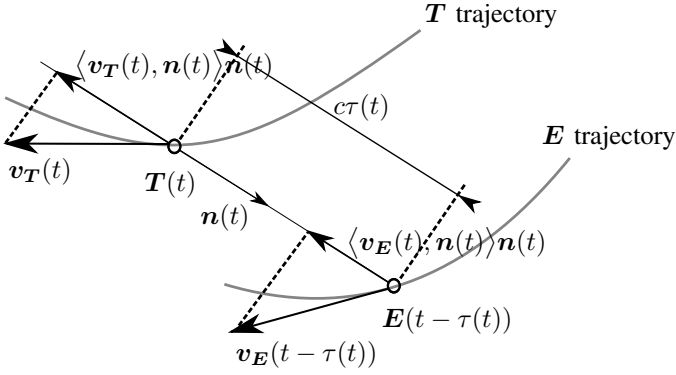


Fig. 2. Geometry of sound propagation

### B. Approximations

Developing the square of (9) and using the relationship

$$\mathbf{OE}(t - \tau(t)) = \mathbf{OE}(t) - \int_{t-\tau(t)}^t \mathbf{v}_E(u) du \quad (13)$$

leads to the equation

$$\begin{aligned} c^2\tau^2(t) &= \|\mathbf{OT}(t)\|^2 + \|\mathbf{OE}(t)\|^2 - 2\langle \mathbf{OT}(t), \mathbf{OE}(t) \rangle \\ &+ \left\| \int_{t-\tau(t)}^t \mathbf{v}_E(u) du \right\|^2 \\ &+ 2\langle (\mathbf{OT}(t) - \mathbf{OE}(t)), \int_{t-\tau(t)}^t \mathbf{v}_E(u) du \rangle, \end{aligned} \quad (14)$$

which can be approximated in various ways.

1) *Zero-th order approximation:* For a slowly moving source, considering that the vector  $\mathbf{OE}$  is constant during the time delay  $\tau$  in (14) leads to

$$c\tau(t) \approx \sqrt{\|\mathbf{OT}(t)\|^2 + \|\mathbf{OE}(t)\|^2 - 2\langle \mathbf{OT}(t), \mathbf{OE}(t) \rangle} \quad (15)$$

2) *First order approximation:* If the emitter velocity is not small but constant during the time delay  $\tau$ , then taking  $\mathbf{v}_E(t) \equiv \mathbf{v}_E$  out of integral terms in (14) results in

$$\alpha\tau^2 + \beta\tau(t) + \gamma \approx 0 \quad (16)$$

$$\begin{aligned} \text{with } \alpha &= \|\mathbf{v}_E(t)\|^2 - c^2 \\ \beta &= 2\langle (\mathbf{OT}(t) - \mathbf{OE}(t)), \mathbf{v}_E(t) \rangle \\ \gamma &= \|\mathbf{OT}(t)\|^2 + \|\mathbf{OE}(t)\|^2 - 2\langle \mathbf{OT}(t), \mathbf{OE}(t) \rangle. \end{aligned}$$

As only subsonic motions are considered, (16) always admits a single positive solution for  $\tau$ , which depends both on positions and velocities.

### C. Measurement equation

Recall that two transducers  $R_l, R_r$  are used, and that the source signal is assumed unknown. The theoretical binaural Interaural Time Difference (ITD) then comes as

$$\text{ITD}(t) = \tau_l(t) - \tau_r(t), \quad (17)$$

with  $\tau_l(t), \tau_r(t)$  the propagation delays from the source to  $R_l, R_r$ , respectively. Define  $T \in \{R_l, R_r\}$ ,  $\varepsilon_{R_l} = -1$ ,  $\varepsilon_{R_r} = +1$ . From (15)–(16) and the additional relations

$$\begin{pmatrix} \langle \mathbf{OT}, \mathbf{y}_O \rangle \\ \langle \mathbf{OT}, \mathbf{z}_O \rangle \end{pmatrix} = \begin{pmatrix} \langle \mathbf{OR}, \mathbf{y}_O \rangle \\ \langle \mathbf{OR}, \mathbf{z}_O \rangle \end{pmatrix} - \varepsilon_T a \begin{pmatrix} \cos\lambda \\ -\sin\lambda \end{pmatrix}, \quad (18)$$

$$\begin{pmatrix} e_y \\ e_z \end{pmatrix} = \begin{pmatrix} \cos\lambda & -\sin\lambda \\ \sin\lambda & \cos\lambda \end{pmatrix} \begin{pmatrix} \langle \mathbf{OE} - \mathbf{OT}, \mathbf{y}_O \rangle \\ \langle \mathbf{OE} - \mathbf{OT}, \mathbf{z}_O \rangle \end{pmatrix}, \quad (19)$$

approximate static equations linking the time delays with the location variables can be obtained. On this basis, an ordinary—static—output equation can be built for ITD cues. This equation will then be part of the state space model underlying the filtering based localization. The approximate output equations are as follows:

1) *Zero-th order approximation:*

$$c\tau_{l,r} = \sqrt{a^2 + e_y^2(t) + e_z^2 + \varepsilon_{R_{l,r}} 2ae_y}. \quad (20)$$

2) *First order approximation:*

$\tau_{l,r}$  is the positive root of  $\alpha\tau_{l,r}^2 + \beta\tau_{l,r} + \gamma$

$$\text{with } \alpha = v_{E_y}^2 + v_{E_z}^2 - c^2 \quad (21)$$

$$\beta = 2\left\langle \begin{pmatrix} -\cos\lambda & \sin\lambda \\ \sin\lambda & \cos\lambda \end{pmatrix} \begin{pmatrix} e_y \\ e_z \end{pmatrix} - \varepsilon_{R_{l,r}} a \begin{pmatrix} \cos\lambda \\ -\sin\lambda \end{pmatrix}, \begin{pmatrix} v_{E_y} \\ v_{E_z} \end{pmatrix} \right\rangle$$

$$\gamma = a^2 + e_y^2 + e_z^2 + \varepsilon_{R_{l,r}} 2ae_y.$$

### D. Extraction of the ITD from the raw signals

Despite an output equation can now be defined, the way the measurements—viz. the binaural cues—can be elaborated still need to be discussed. This process is based on the estimation of the temporal shift between the raw signals sensed by the binaural sensor while it is moving. One approach to measure a constant time delay between the signals  $s_{R_l}(t)$  and  $s_{R_r}(t)$ , provided they are individually and jointly stationary, consists in extracting the argmax of their cross-correlation  $C_{lr}(u) = \mathbb{E}[s_{R_l}(t)s_{R_r}(t-u)]$ . From the data gathered on a finite time window of length  $T$  during a single experiment, an estimate of  $C_{lr}$  can be computed as follows:

$$\hat{C}_{lr}(u) = \frac{1}{T} \int_{\tau}^T s_{R_l}(t)s_{R_r}(t-u)dt. \quad (22)$$

However, a number of potential issues underlying the position of the peak in  $\hat{C}_{lr}$  must be taken into account. First, the time delay to be estimated must vary slowly within the observation time  $T$ . This is a fundamental concern when the source and sensor move. Next, the finiteness of the time window may result in a wide and inaccurate peak of  $\hat{C}_{lr}$ , e.g. when the measurement noise is important or if the source is narrowband. As a solution, [18] suggests to use the so-called “Generalized Cross Correlation” (GCC), which consists in weighting the signals Cross Power Spectral Density. Among all the frequency-weighting functions  $\psi$  proposed in [18], the Phase Transform (PHAT) processor may be the most popular. An enhancement, called Reliability-Weighted Phase Transform (RW-PHAT) consists in reducing the weighting function at frequencies where the signal-to-noise ratio is low [19]. One can also cite the Maximum Likelihood (ML) processor [18], which exhibits an asymptotic Gaussianity, making possible to consider the measurement noise in (3) as Gaussian.

In all the following, the considered sound source will emit a Gaussian white noise so as to simplify the auditory cues estimation problem. Consequently, only the digital version of (22) will be considered. But it is known that the sampling frequency  $f_s$  affects the temporal shift estimator statistics, introducing quantization errors in the estimation,

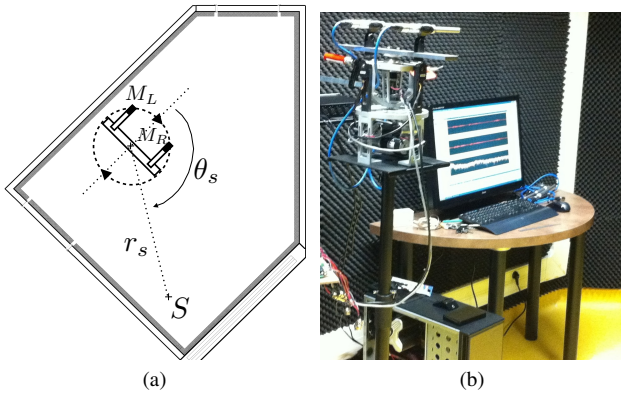


Fig. 3. (a) Experimental setup scheme. The two left ( $M_L$ ) and right ( $M_R$ ) omnidirectional microphones are circularly moved by a motoreductor. The loudspeaker ( $S$ ) emits a white noise during the movement. (b) Experimental setup picture, showing the two microphones and the acquisition computer.

as noninteger multiples of sampling period delays cannot be reached. This problem can be partially solved by the simple approximation of fitting a parabola in the neighborhood of the cross correlation peak [20].

## V. CASE STUDY

In this section, experimental and simulation results are shown to confirm the proposed approach. In simulations as in the experiment, the position of a static pointwise source is estimated w.r.t a moving binaural sensor. The results reveal that the proposed filtering strategy, by using information coming from relative motion, overcomes front-back ambiguity problems inherent to ITD, and can provide range estimation despite ITD does not carry information about distance.

### A. Experimental results

1) *Experimental setup*: In order to assess the proposed approach with real binaural signals, several experiments have been performed in an acoustically prepared room, equipped with 3D pyramidal pattern studio foams placed on the roof and on the walls. Two omnidirectional microphones, spaced by 18cm, are then placed on the top of a simple mechanical system made of a single rotation moved by a motoreductor. The angular axis position is recorded during the movement with a  $0.5^\circ$ -precision, producing the proprioception of the platform. The two microphone outputs are simultaneously acquired by a National Instruments PCI acquisition card through 24 bits delta-sigma converters operating at the sampling frequency  $f_s = 44.1\text{kHz}$ . A loudspeaker, diffusing a white gaussian noise, is placed in front of the mechanical system at the position  $(r_s, \theta_s) = (1.48\text{m}, 120^\circ)$  (see Figure 3). The recording of the proprioception and of the auditory perception is simultaneously started at the beginning of the movement, which consists in a rotation performed at successive various constant angular velocities  $\omega$ .

2) *Sound source estimation results*: Figure 4b shows the source location estimate (azimuth and distance) as well as the 99% confidence interval as a function of time in the

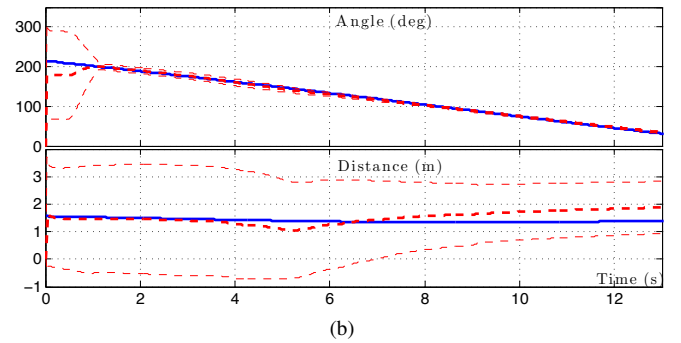
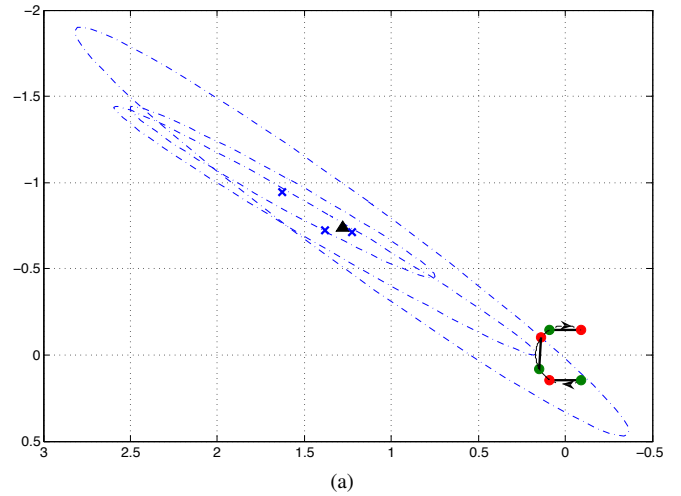


Fig. 4. Experimental results. (a) Sensor motion and emitter position estimate in the world frame (b) Range and azimuth estimates in the sensor frame as a function of time.

sensor frame. The azimuth estimate exhibits good performance, reaching its steady state in about two seconds, with a  $3^\circ$  standard deviation. The same does not apply for the distance estimate, which shows a 1m standard deviation. This large uncertainty in distance is depicted in Figure 4a, where very flat confidence ellipses (dotted lines) are reported, extending towards the source (triangle)-sensor (line segment) direction. In fact, some filters of the MH-SRUKF which were initialized with an azimuth close to the real source azimuth but with a different range might produce a likely predicted output w.r.t the measurement, because, again, ITD does not provide much information about range. If the type of sensor motion does not bring sufficient additional information, the likelihood weights of these filters will remain above the minimum threshold, they will not collapse, and the global estimate, which is a combination of all active filters estimates, will remain inaccurate in distance. Thus, the motion considered in this experiment is not suited to segregate all inconsistent estimates w.r.t range. Because of limitations of the mechanical system, which is only able to perform a simple rotation of the interaural axis, a realistic simulation has been performed with a different type of motion, in order to bring to the fore the sensor motion influence.

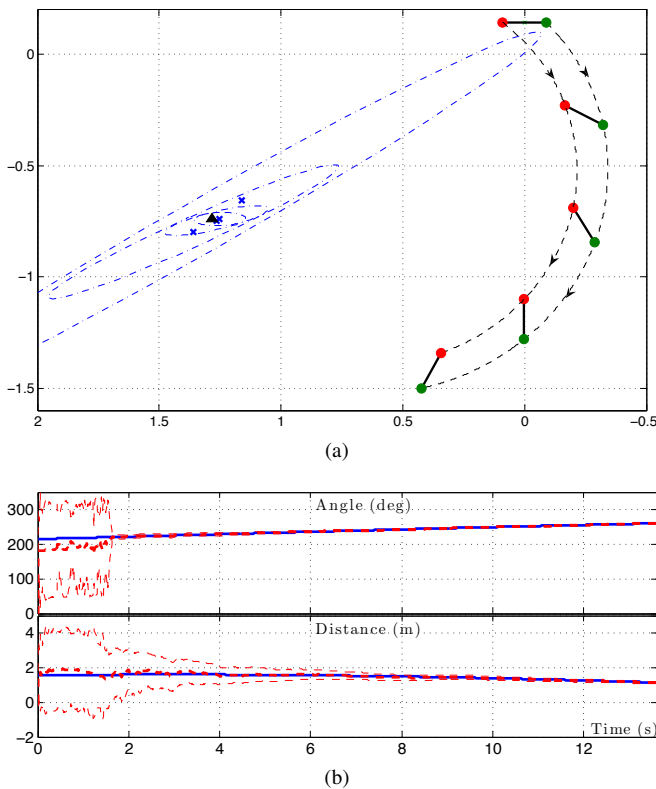


Fig. 5. Simulation : the sensor follows a circular trajectory. (a) Sensor motion and emitter position estimate in the world frame (b) Range and azimuth estimates in the sensor frame as a function of time.

### B. Simulation results

In the simulation reported on Figure 5a, the sensor center follows a circular trajectory, with a constant interaural axis velocity. This motion enables a much more accurate range estimation, as shown in Figure 5b. Indeed, the confidence ellipsoid shrinks along all directions during the estimation. An auditive feedback control of the sensor motion could be designed in order to improve source localization accuracy.

## VI. CONCLUSION

A new active approach to binaural sound source localization has been proposed. It relies on an accurate modeling of the acoustic cues, providing ground credible equations for Kalman filtering applications. Theoretical developments have been supported by simulation and experiments.

Although an accurate model for ITD has been put forward in this paper, some limitations appear in the ITD extraction when the source-sensor relative velocity becomes important. In fact, high time delay variations induced by fast relative motions decrease the accuracy of classical cross correlation methods, which have been mainly designed for slowly varying ITDs [18]. An interesting future work could consist in using an enhanced cross correlation method accounting for source or sensor motion, as described in [21].

Another improvement will study the coupling of the proposed filtering scheme with detections of wrong ITDs/ILDs values computed during source silences (e.g. based on

residual monitoring or hypotheses testing). Various re-initialization strategies will then be evaluated.

Other prospective issues would consider a more sophisticated binaural system including a head, and account for room reverberation. Acoustic considerations then would have to be revisited accordingly. Finally, MUSIC-based super-resolution time delay estimation methods developed in [22] would be a possible alternative to Generalized Cross Correlation, exhibiting a higher accuracy in a reverberant or multi-source context.

### ACKNOWLEDGMENT

This work was conducted within the BINAHR (BINAural Active Audition for Humanoid Robots) project funded by ANR (France) and JST (Japan) under Contract n° ANR-09-BLAN-0370-02. Detailed analytical computations are not included for space reasons, but can be obtained on request.

### REFERENCES

- [1] T. Takahashi, K. Nakadai, K. Komatani, T. Ogata, and H. G. Okuno, "An improvement in automatic speech recognition using soft missing feature masks for robot audition," *IEEE/RSJ IROS'2010*, pp. 964–969.
- [2] K. Youssef, S. Argentieri, and J.-L. Zarader, "From monaural to binaural speaker recognition for humanoid robots," in *IEEE Humanoids'2010*, pp. 580–586.
- [3] T. Rodemann, "A study on distance estimation in binaural sound localization," in *IEEE/RSJ IROS'2010*, pp. 425–430.
- [4] F. Keyrouz and K. Diepold, "An enhanced binaural 3D sound localization algorithm," in *IEEE Symp. on Sig. Proc. and Info. Tech. 2006*.
- [5] J.M. Valin, F. Michaud, J. Rouat and D. Létourneau, *Robust sound source localization using a microphone array on a mobile robot*, in *IEEE/RSJ IROS'2003*, pp. 1228–1233.
- [6] H. L. Van Trees, *Optimum Array Processing*, ser. Detection, Estimation, and Modulation Theory. John Wiley & Sons, Inc., 2002, vol. IV.
- [7] M. Cooke, Y.-C. Lu, Y. Lu, and R. P. Horaud, "Active hearing, active speaking," in *Int. Symp. on Auditory and Audiological Res.*, 2007.
- [8] K. Nakadai, T. Lourens, H. Okuno, and H. Kitano, "Active audition for humanoid," in *17th Nat. Conf. on Artificial Intelligence*, 2000.
- [9] K. Nakadai, H. G. Okuno, and H. Kitano, "Robot recognizes three simultaneous speech by active audition," in *IEEE ICRA'2003*.
- [10] E. Berglund and J. Sitte, "Sound source localisation through active audition," in *IEEE/RSJ IROS'2005*, pp. 509–514.
- [11] L. Kneip and C. Baumann, "Binaural model for artificial spatial sound localization based on interaural time delays and movements of the interaural axis," *Jour. Acoust. Soc. America*, 124(5), pp. 3108–19, 2008.
- [12] Y.-C. Lu and M. Cooke, "Motion strategies for binaural localisation of speech sources in azimuth and distance by artificial listeners," *Speech Comm.*, 2010.
- [13] M. Kumon and S. Uozumi, "Binaural localization for a mobile sound source," *Jour. of Biomechanical Science and Engineering*, 6(1), 2011.
- [14] H. Viste and G. Evangelista, "Binaural source localization," in *Int. Conf. on Digital Audio Effects (DAFx'04)*.
- [15] S. Julier and J. Uhlmann, "Unscented filtering and nonlinear estimation," *Proceedings of IEEE*, 92(3), 2004.
- [16] R. Van der Merwe and E. Wan, "The square-root unscented kalman filter for state and parameter estimation," in *IEEE ICASSP'01*.
- [17] Y. Bar-Shalom and X. Li, *Estimation and Tracking : Principles, Techniques and Software*. YBS editions, Norwood, 1993.
- [18] C. Knapp and G. Carter, "The generalized correlation method for estimation of time delay," *IEEE Trans. ASSP*, 24(4), pp. 320–327.
- [19] K. Nakadai, H. Okuno, and H. Kitano, "Auditory fovea based speech separation and its application to dialog system," *IEEE/RSJ IROS'2002*.
- [20] R. E. Boucher and J. C. Hassab, "Analysis of discrete implementation of generalized cross correlator," *IEEE Trans. ASSP*, 29(3).
- [21] C. Knapp and G. Carter, "Time delay estimation in the presence of relative motion," *IEEE Trans. ASSP*, pp. 280–283, 1977.
- [22] F. Ge, D. Shen, Y. Peng, and V. Li, "Super-resolution time delay estimation in multipath environments," *IEEE Wireless Communications and Networking Conf.*, pp. 1121–1126, 2004.