

Active Speaker Localization with Circular Likelihoods and Bootstrap Filtering*

Ivan Marković¹, Alban Portello², Patrick Danès³, Ivan Petrović⁴, Sylvain Argentieri⁵

Abstract—This paper deals with speaker localization in two dimensions from a mobile binaural head. A bootstrap particle filtering scheme is used to perform active localization, i.e. to infer source location by fusing the binaural perception with the sensor motor commands. It relies on an original pseudo-likelihood of the source azimuth which captures both the inter-aural level and phase differences. Since the pseudo-likelihood is discrete, it is fitted with a mixture of circular distributions in order to enhance its resolution. For the fitting task two mixtures are compared and evaluated, namely the mixture of von Mises and wrapped Cauchy distributions. Furthermore, a solution is presented for calculating the von Mises curvefitting with low uncertainty, since the direct implementation can quickly surpass double precision floating number representation. The performance of the filter is compared using both the raw and fitted pseudo-likelihoods on experiments recorded in an acoustically prepared room with ground-truth obtained from a motion capture system. The results show that the proposed algorithm successfully localizes the speaker with an advantage in the direction of the fitted von Mises mixture likelihood.

I. INTRODUCTION

In the field of robotics, the subject of sound source localization has been approached and studied from aspects of many different fields, namely speech processing, estimation theory, and sensor fusion to name but a few. From the aspect of sensors, researchers have been using binaural setups, microphone arrays featuring several or more than a hundred of microphones, placing them on wheeled mobile robots, humanoid walking robots, and even autonomous aerial vehicles. Furthermore, when moving sensors are used, the seamless fusion of their motor commands with the binaural perception—active localization—has been acknowledged to overcome ambiguities inherent to the use of static sensors.

Tracking (localization) with bearing-only values is a challenging problem due to the non-linearity of the measurement equation and the unobservability of the target-observer distance which, indeed, can only be estimated by motion (activ-

ity) from the observer's side. In [1] it was shown that tracking in modified polar coordinates with an extended Kalman filter (EKF) provided better and more stable results than when tracking in Cartesian coordinates. This brought higher complexity in the motion model, but made the observation model linear and separated observable and unobservable entries in the state vector. This model was further developed in [2] where the tracking was performed with a bank of range-parameterized EKFs in modified polar coordinates. Although this problem has been studied for few decades, it still receives attention due to emerging new filtering methods. In [3] three different filters were compared for the task, while in [4] various methods for tracking and decentralized sensor fusion were studied, including bearing-only scenarios. In [5], relative localization is performed from a pair of moving microphones, based on a multiple-hypothesis square-root unscented Kalman filter. The filtering scheme uses time delays estimated from the sensed audio signals, together with information on the sensor's velocities to perform a consistent source localization. Results show that the strategy, together with a suitable sensor motion, allows to break front-back ambiguity and get accurate range information.

In the context of speaker localization, the bootstrap particle filter has been utilized in [6] for multiple speaker bearing and elevation estimation with an 8-channel microphone array mounted on a mobile robot. In [7], the authors address the problem of localizing multiple sound sources in an outdoor environment from a microphone array mounted on an aerial vehicle. An extension of the MUSIC algorithm is used, that uses adaptive estimation of the—dynamically changing—environment noise correlation matrix. The proposed method is tested with a Parrot AR.Drone and a Kinect device. In [8] the authors used a 4-channel array to localize narrow-band emergency signals from a micro air vehicle, where the sensor model was based on the cross-correlation and doppler shift in frequency due to the motion of the vehicle. In [9] the particle filter was used to estimate the bearing of a speaker from a von Mises (vM) mixture with a 4-channel array mounted on a mobile robot. In a non-robotic related context, in [10] particle filtering was utilized to estimate a position of a speaker in a room environment with 4 microphone pairs placed on the room walls, where the generalized cross-correlation and beamformer output power were used as pseudo-likelihood functions. In [11] the authors analyzed strategies for sensor motion in the context of speaker localization with PF in both range and bearing and performed evaluations in a simulated acoustic environment with single sources under both anechoic and reverberant conditions. In [12] the authors

*This work has been supported by European Community's Seventh Framework Programme under grant agreement no. 285939 (ACROSS) and the BINAHR (BINAural Active Audition for Humanoid Robots) project funded by ANR (France) and JST (Japan) under Contract nANR-09-BLAN-0370-02.

^{1,4}I. Marković and I. Petrović are with University of Zagreb, Faculty of Electrical Engineering and Computing, Department of Control and Computer Engineering, HR-10000 Zagreb, Croatia.

^{2,3}A. Portello and P. Danès are with CNRS, LAAS, 7 avenue du colonel Roche, F-31400 Toulouse, France, Univ de Toulouse, UPS, LAAS; F-31400 Toulouse, France.

⁵S. Argentieri is with UPMC Univ. Paris 06, 4 place Jussieu, F-75005, Paris, France and ISIR - CNRS UMR 7222, F-75005, Paris, France. (ivan.markovic, ivan.petrovic) at fer.hr, (alban.portello, patrick.danes) at laas.fr, sylvain.argentieri at upmc.fr

used a combination of direction-of-arrival estimates with speaker's fundamental frequencies (pitch) and gammatone prefiltering to form a pseudo-likelihood function for a 24-channel circular array in order to estimate the bearings of multiple speakers.

In the present paper, active speaker localization is performed with two microphones mounted on a spherical head by bootstrap particle filtering [13]. The underlying state space equation describing the evolution of the source position in the head frame is defined in both cartesian and polar coordinates. We propose a pseudo-likelihood function of the source bearing (azimuth) as the measurement model, which captures both the interaural phase difference (IPD) and interaural level difference (ILD) between the binaural signals. Since the pseudo-likelihood has no analytic expression and is only given for a discrete set of candidate bearings, the fitting of circular distributions to the discrete pseudo-likelihood is discussed in order to enhance its resolution for the purpose of estimation. Incidentally, this can give further ground for possible analytical filtering schemes. Two distributions are presented and compared for the task: namely the vM distribution, for which we also present a method for evaluation with a large concentration parameter, and the wrapped Cauchy (WC) distribution. Furthermore, we compare two bootstrap particle filtering schemes on experimental data—one using the raw discrete pseudo-likelihood, and the other based on the fitted circular distribution. As aforementioned, both fuse the known head velocities with binaural data in order to infer the speaker location.

The paper is organized as follows. First, the problem is stated in §II, while §III presents and compares the proposed fitting with the vM and WC distributions. In §IV the proposed speaker localization with the bootstrap algorithm is presented, §IV-B presents the experimental evaluation, and in the end §V concludes the paper.

II. PROBLEM STATEMENT

A. Kinematics and state space equation

A pointwise sound emitter E and a binaural sensor lie on a common plane parallel to the ground. The two receivers equipping the sensor are denoted by R_1 and R_2 . 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 $[R_1 R_2]$, \mathbf{y}_R the vector $\frac{\mathbf{R}R_1}{|\mathbf{R}R_1|}$ 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 $\mathbf{x}_O = \mathbf{x}_R$ (see Fig. 1). The source is assumed motionless w.r.t. the world frame, 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 ω of \mathcal{F}_R w.r.t. \mathcal{F}_O around $\mathbf{x}_O = \mathbf{x}_R$). Assuming v_{Ry}, v_{Rz}, ω are 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_1, R_2 . The audio sensor location w.r.t. \mathcal{F}_O is not required and in this paper the localization of the mobile base is not performed. The relative attitude of

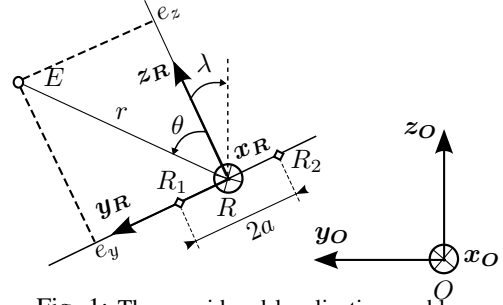


Fig. 1: The considered localization problem.

\mathcal{F}_R w.r.t. \mathcal{F}_E can be described, when v_{Ry}, v_{Rz}, ω are zero-order held at the sampling period T_s , by the discrete-time deterministic state space equation

$$\mathbf{x}_{t+1} = F\mathbf{x}_t + G_1\mathbf{u}_{1t}, \text{ with}$$

$$F = \begin{bmatrix} \cos(\omega T_s) & \sin(\omega T_s) & 0 \\ -\sin(\omega T_s) & \cos(\omega T_s) & 0 \\ 0 & 0 & 1 \end{bmatrix}, \quad G_1 = - \begin{bmatrix} \frac{\sin(\omega T_s)}{\omega} & \frac{1 - \cos(\omega T_s)}{\omega} & 0 \\ \frac{\cos(\omega T_s) - 1}{\omega} & \frac{\sin(\omega T_s)}{\omega} & 0 \\ 0 & 0 & T_s \end{bmatrix},$$

Therein, the state vector $\mathbf{x} \triangleq [e_y, e_z, \lambda]^T$ gathers the entries e_y and e_z (the \mathbf{y}_R and \mathbf{z}_R component of E in \mathcal{F}_R) and the orientation angle λ . The sensor velocities constituting $\mathbf{u}_1 \triangleq [v_{Ry}, v_{Rz}, \omega]^T$ are supposed known. When parameterizing the problem in terms of polar coordinates rather than cartesian, i.e. when using the variables $\theta \triangleq \text{atan2}(e_y, e_z)$, $r \triangleq \sqrt{e_y^2 + e_z^2}$, the state space equation comes as

$$\mathbf{r}_{t+1} = \sqrt{r_t^2 + \mathbf{u}_t^T G^T G \mathbf{u}_t + 2r_t [\sin\theta_t, \cos\theta_t] G^T \mathbf{u}_t} \quad (1)$$

$$\theta_{t+1} = \text{atan2}(r_t \sin(\theta_t + \omega T_s) + \mathbf{g}_1 \mathbf{u}_t, r_t \cos(\theta_t + \omega T_s) + \mathbf{g}_2 \mathbf{u}_t)$$

$$\lambda_{t+1} = \lambda_t - \omega T_s, \quad (2)$$

with $\mathbf{u} \triangleq [v_{Ry}, v_{Rz}]^T$, G the square matrix made up with the first two rows and columns of G_1 , \mathbf{g}_1 (resp. \mathbf{g}_2) the first (resp. second) row of G . To model uncertainty in the relative motion, a random white Gaussian noise of known statistics is added to (1).

B. Acoustic model, measurement vector, pseudo-likelihood

Consider first a static world where the sensor is motionless. We assume that the source lies in the farfield (i.e. the source range $r = |\overline{RE}|$ is sufficiently high compared to the microphones interspace $2a$ so that the source wavefronts can be considered as planar in the vicinity of the microphone pair). We model the signals y_1, y_2 monitored at R_1, R_2 in the presence of additive noise as follows

$$\begin{cases} y_1(\tau) = s(\tau) + n_1(\tau) \\ y_2(\tau) = (s * h_\theta)(\tau) + n_2(\tau), \end{cases} \quad (3)$$

where the signal s (i.e. the contribution of the emitter at R_1) and the noises n_1, n_2 are real, band-limited, individually and jointly stationary random processes, and $*$ denotes convolution. The deterministic impulse response h_θ between R_1, R_2 , is parameterized by θ , and captures free-field propagation of the emitted signal as well as head scattering. H_θ , the Fourier transform of h_θ , is supposed known for every θ within a discrete set of values (say, it has been learned

from calibration, or is known theoretically). The process $\mathbf{y}(\tau) \triangleq [y_1(\tau), y_2(\tau)]^T$ is observed over N adjacent non-overlapping rectangular T/N -width time windows. Denote \mathbf{y}_n the observation of \mathbf{y} over the n^{th} window. A data vector \mathbf{Z} is made up by stacking the values of

$$\mathbf{Y}_n[k] = \sqrt{\frac{N}{T}} \int_{\mathbb{R}} \mathbf{y}_n(\tau) e^{-2i\pi k \frac{N}{T} \tau} d\tau, \quad n = 1, \dots, N \quad (4)$$

at $k = k_1, \dots, k_B$, the B frequency indexes within the bandwidth of s . \mathbf{Z} is hence defined as $\mathbf{Z} \triangleq [\mathbf{Y}[k_1]^T, \dots, \mathbf{Y}[k_B]^T]^T$, with $\mathbf{Y}[k] \triangleq [\mathbf{Y}_1[k]^T, \dots, \mathbf{Y}_N[k]^T]^T$. Assume now that s, n_1, n_2 are zero-mean jointly Gaussian and that n_1, n_2 are identically distributed, uncorrelated with each other and with s . Then, under general mild conditions on the power spectra of s, n_1, n_2 and on H_θ , the maximum likelihood estimate of θ can be obtained, given a sample \mathbf{z} of \mathbf{Z} , by maximizing the following criterion [14], hereafter referred to as the “pseudo log-likelihood function”

$$J(\mathbf{z}|\theta) = c_2 - N \sum_{k=k_1}^{k_B} \left(\ln |P_\theta[k] \hat{C}[k] P_\theta[k] + \hat{\sigma}_\theta^2[k] P_\theta^\perp[k]| \right), \quad (5)$$

with $c_2 \triangleq -2NB(\ln(\pi) + 1)$, $\hat{C}[k] \triangleq \frac{1}{N} \sum_n \mathbf{y}_n[k] \mathbf{y}_n[k]^\dagger$, $P_\theta[k] \triangleq \mathbf{V}_\theta[k] (\mathbf{V}_\theta[k]^\dagger \mathbf{V}_\theta[k])^{-1} \mathbf{V}_\theta[k]^\dagger$, $P_\theta^\perp[k] \triangleq \mathbb{I}_2 - P_\theta[k]$, $\mathbf{V}_\theta[k] \triangleq [1, H_\theta[k]]^T$, $\hat{\sigma}_\theta^2[k] \triangleq \text{tr}(P_\theta^\perp[k] \hat{C}[k])$.

Therein, $^\dagger, ^\perp, |\cdot|, \text{tr}(\cdot)$ respectively stand for Hermitian transpose, orthogonal complement, determinant and trace; $\mathbf{y}_n[k]$ denotes a sample of $\mathbf{Y}_n[k]$, and the sample covariance matrix \hat{C} is assumed full rank.

Consider now a real world where the sensor moves and where the source signal and environment noise are possibly nonstationary. All the fundamental hypotheses leading to (3)–(4)–(5) are consequently violated. Nevertheless, the problem can still be handled if, at each process time index t , the data vector \mathbf{z}_t is made up from audio data collected over a time window matched to t , sufficiently short so that, along this window, the sensor motion is negligible and the recorded signals can be regarded as finite-time samples of stationary processes. Hence, at each time index t , the pseudo likelihood of θ_t w.r.t. \mathbf{z}_t , denoted $p(\mathbf{z}_t|\theta_t)$, can be output and will henceforth be used in a Bayesian filtering scheme in §IV. Importantly, $p(\mathbf{z}_t|\theta_t)$ has in the general case no analytic expression. Its numerical values are just given for a discrete set of tested azimuths. This precludes the use of Bayesian filtering schemes requiring an analytic form of the pseudo likelihood, e.g. Gaussian mixture filters, unless an analytic function is fitted to the discrete values. Alternatively, with particle filters, the pseudo likelihood in its discrete form can be utilized as a sensor model. However, low azimuth resolution can affect the particle filter performance and consistency, and it may be useful to fit some distribution to the discrete pseudo likelihood. §III is thus dedicated to the fitting of Von Mises and wrapped Cauchy mixtures models to the discrete pseudo likelihood.

III. FITTING THE CIRCULAR DISTRIBUTIONS

A. Circular distributions

In this section we present two solutions to fitting the pseudo likelihood function, namely fitting with the vM distribution and with the wC distribution. The motivation behind using circular distributions lies in the fact that they intrinsically take noneuclidian properties of the angular data into account. For an example, this property proves useful in the optimization since a circular distribution close to π continues contributing to points larger than $-\pi$. Furthermore, in the present paper we do not require the component weights to sum up to one, since the pseudo likelihood function itself is not a valid probability distribution.

A probability distribution on the unit circle with density function given by [15]

$$p(\theta; \mu, \kappa) = \frac{1}{2\pi I_0(\kappa)} \exp \{ \kappa \cos(\theta - \mu) \}, \quad (6)$$

is called the von Mises distribution, where $0 \leq x \leq 2\pi$, μ is the mean direction, $\kappa \geq 0$ is the concentration parameter, and $I_0(\kappa)$ is the modified Bessel function of the first kind and of order zero. The distribution is unimodal and symmetric around μ and is often referred to as the circular analogue of the Gaussian distribution. When $\kappa \rightarrow 0$ the vM becomes the uniform distribution, while if $\kappa \rightarrow \infty$ it becomes concentrated at $\theta = \mu$.

We used the vM distribution in the context of robot audition in [9] where the sensor model was represented as a mixture of vM distribution in particle filtering, while in [16] we extended this approach to model the entire Bayesian tracking procedure in the analytical domain of the distribution. However, both of the aforementioned works were only concerned with tracking the bearing value of the speaker and not its position in two dimensions which is one of the goals of the present paper.

The second distribution that we analyze for the task is a distribution wrapped on a circle. Given a distribution on the line we can wrap it around the circumference of a circle with unit radius. If a random variable θ is defined on a line, then the corresponding random variable of the wrapped distribution is $\theta_w = \theta \pmod{2\pi}$. Furthermore, if θ has a probability density function (pdf) p , then the corresponding wrapped pdf p_w is defined as $p_w(\theta) = \sum_{k=-\infty}^{k=\infty} p(\theta + 2k\pi)$ [15], from which we can note practical issues when dealing with the infinite number of terms in the summation. However, it can be shown that the Cauchy distribution on the line has an interesting property that its wrapped counterpart, due to certain geometric series expansion property, reduces to [15]

$$p(\theta; \mu, \rho) = \frac{1}{2\pi} \frac{1 - \rho^2}{1 + \rho^2 - 2\rho \cos(\theta - \mu)}, \quad (7)$$

where μ is the mean direction and ρ is called the mean resultant length. When $\rho \rightarrow 0$ the wC tends to uniform distribution, while if $\rho \rightarrow \infty$ the distributions becomes concentrated at point μ .

Naturally, the pseudo likelihood function will suffer from front-back ambiguity since in the present paper we utilize

a binaural setup. Hence, our sensor model will contain at least two distinct modes on the interval 0 to 2π and for this reason we chose to model the likelihood as a mixture of distributions. If we denote with \mathcal{X} a set of distributions parameters, then an N component mixture can be defined as $p(\theta; \mathcal{X}) = \sum_{i=1}^N \omega_i p(\theta; \mathcal{X}_i)$, where the set \mathcal{X} consists of $\cup_i \{\mu_i, \kappa_i\}$ for the vM distribution and $\cup_i \{\mu_i, \rho_i\}$ for the wC distribution.

B. Calculation of von Mises distributions with large κ

The direct form of the vM distribution suffers from numerical issues when working with large concentration parameter κ , i.e. with sharp distributions which may be necessary to fit the pseudo likelihood in the vicinity of its local modes. The main problem is that for large κ both the exponent and the modified Bessel function of the first kind quickly reach the maximum value that can be stored in double precision floating point representation.

To solve this problem, we move the normalizer of the vM distribution in the exponent as follows

$$p(\theta; \mu, \kappa) = \exp \{ \kappa \cos(\theta - \mu) - \log(2\pi I_0(\kappa)) \}, \quad (8)$$

and approximate $\log(I_0(\kappa))$ as [17]

$$\log(I_0(x)) = m(x) + \log \sum_{k=0}^{\infty} \exp \{ t_k(x) - m(x) \}, \quad (9)$$

where $t_k(x) = 2k \log \frac{x}{2} - 2 \sum_{r=0}^k \log r$ and $m(x) = \max\{t_k(x)\}$. The number of the terms in (9) required to have an accurate approximation depends on the κ . For the present application we have found that the maximal practical value of the concentration parameter for fitting the pseudo likelihood is $\kappa = 2000$, for which an accurate approximation was empirically determined to be for $k \leq 1100$. But for smaller parameters, e.g. $\kappa = 1000$, the number of terms $k \leq 600$ was sufficient. We did not notice any increase in the computational time when compared to Matlab's implementation based on [18].

C. Evaluation of the fitting performance

The fitting of a mixture of distributions to the pseudo likelihood function $\hat{p}(\theta)$ comes down to solving the following optimization problem

$$\min_{\omega, \mathcal{X}} \sum_{i=1}^N (\omega_i p(\theta; \mathcal{X}_i) - \hat{p}(\theta))^2$$

with the constraints $0 \leq \omega_i \leq 1$ and $0 \leq \mathcal{X}_i \leq \mathcal{B}$ for $i = 1, \dots, N$, and where the upper bound \mathcal{B} depends on the parameter and the distribution. For both distributions the upper bound of the mean directions μ is $\mathcal{B} = 2\pi$, while for the vM distribution the upper bound was $\mathcal{B} = 2000$ for the concentration parameter, and for the wC distribution $\mathcal{B} = 1$ for the mean resultant length.

Concerning the number of the mixture components all the results were obtained with $N = 4$. Initial conditions for the mean directions were determined by searching recursively for N most dominant peaks in the vein of [6] where the

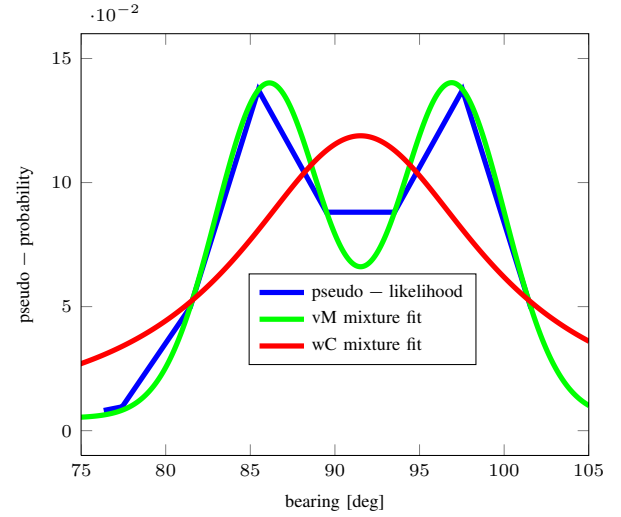


Fig. 2: Fitting the pseudo likelihood for a single frame with vM and wC mixture

authors searched for the number of active speakers. Once the dominant peak is found, an area around it is removed and the search continues until the predetermined number of modes is found. Since in the pseudo likelihood function we expect two peaks to be dominant we set the initial conditions for the first two dominant peaks to be $\kappa = 1500$ or $\rho = 0.9$, while for the rest we set $\kappa = 10$ or $\rho = 0.1$. The weights are initially set to $\omega_i = 0.5, \forall i$.

In Fig. 2 we can see the result of fitting for a single relatively difficult frame when the speaker was close to the end-fire position of the array and the two dominant modes started overlapping. Empirically we noticed that this is the more difficult case for the wC distribution and that often the two distinct nodes tend to be fitted with a single component in between them. Overall, the whole dataset consisted of four experiments with a talking speaker as the source. The average root-mean-square-error (RMSE) of fitting for the speaker scenario was $1.6 \cdot 10^{-3}$ for the vM mixture and $3.7 \cdot 10^{-3}$ for the wC mixture, respectively. Given that, for the rest of the paper we have chosen to work with the vM mixture since it provided better fitting in terms of the average RMSE.

IV. SPEAKER LOCALIZATION IN 2D

A. Particle filtering

Particle filtering is a versatile method to recursive Bayesian state estimation. It can handle nonlinear prior dynamics and measurements models, as well as nonGaussian noises. The posterior pdf of the state at any time t conditioned on the sequence of observed measurements up to t is estimated by means of a point-mass probability distribution with stochastic support, or “weighted particle set”. Let $\{\mathbf{x}^p, w^p\}_{p=1}^P$ depict the random measure that characterizes the posterior state pdf $p(\mathbf{x}_t | \mathbf{z}_{1:t})$, where each particle in the set $\{\mathbf{x}^p\}_{p=1}^P$ is associated to the respective weight in $\{w^p\}_{p=1}^P$. The weights satisfy $\sum_p w^p = 1$, so that $p(\mathbf{x}_t | \mathbf{z}_{1:t})$

can be approximated as [13], [19]

$$p(\mathbf{x}_t | \mathbf{z}_{1:t}) \approx \sum_{p=1}^P w_t^p \delta(\mathbf{x}_t - \mathbf{x}_t^p), \quad (10)$$

with $\delta(\cdot)$ the Dirac delta measure. In other words, sampling from $p(\mathbf{x}_t | \mathbf{z}_{1:t})$ returns to sampling a particle with a probability equal to its associated weight.

The particles are drawn according to a so-called importance function, then weighted so that the consequent random measure constitutes a sound approximation to the posterior pdf. As, for any recursive particle filter, the significant weights tend to concentrate on a limited set of particles after few iterations, a resampling step is inserted, which consists in turning $\{\mathbf{x}_t^p, w_t^p\}_{p=1}^P$ into the equivalent evenly weighted set $\{\mathbf{x}_t^{*p}, \frac{1}{P}\}_{p=1}^P$ by independently sampling (with replacement) \mathbf{x}_t^{*p} according to $P(\mathbf{x}_t^{*p} = \mathbf{x}_t^p) = w_t^p$.

In the sequential importance resampling (SIR) scheme [13], or bootstrap filter, the importance function matches the prior dynamics $p(\mathbf{x}_t | \mathbf{x}_{t-1})$, calculated via (1), i.e. each particle \mathbf{x}_t^p at time t is drawn from its predecessor \mathbf{x}_{t-1}^p at time $t-1$ according to the proposal density $\mathbf{x}_t^p \sim p(\mathbf{x}_t | \mathbf{x}_{t-1}^p)$. Then, its weight is updated by evaluating its likelihood $p(\mathbf{z}_t | \mathbf{x}_t^p)$ prior to setting

$$w_t^p \propto w_{t-1}^p p(\mathbf{z}_t | \mathbf{x}_t^p), \quad (11)$$

where $p(\mathbf{z}_t | \mathbf{x}_t)$ represents the sensor model, i.e. the fitted vM mixture:

$$p(\mathbf{z}_t | \mathbf{x}_t) = \sum_{i=1}^N \omega_i \frac{1}{2\pi I_0(\kappa_i)} \exp[\kappa_i \cos(\mathbf{x}_t - \mathbf{z}_{t,i})]. \quad (12)$$

Then, all the particle weights are normalized so that they sum up to unity.

Once the random measure approximating the posterior pdf of the state is computed, the posterior mean and posterior covariance can be estimated via

$$\begin{aligned} \hat{\mathbf{x}}_t &= E[\mathbf{x}_t | \mathbf{z}_{1:t}] \approx \sum_{p=1}^P w_t^p \mathbf{x}_t^p, \\ \hat{\mathbf{P}}_t &= E[(\mathbf{x}_t - E[\mathbf{x}_t | \mathbf{z}_{1:t}])(\mathbf{x}_t - E[\mathbf{x}_t | \mathbf{z}_{1:t}])^T | \mathbf{z}_{1:t}] \\ &\approx \sum_{p=1}^P w_t^p (\mathbf{x}_t^p - \hat{\mathbf{x}}_t)(\mathbf{x}_t^p - \hat{\mathbf{x}}_t)^T. \end{aligned} \quad (13)$$

To avoid a loss of diversity in the particle cloud, the resampling step is applied only when the number of effective weights $P_{\text{eff}} = 1 / \sum_p (w_t^p)^2$ is less than a given threshold, e.g. 33% of the total number of particles P .

Consequently, particle filtering can be implemented even if a closed-form measurement model is not available, in that the particle likelihoods just need to be evaluated. In our case, the sensor model comes as the pseudo likelihood digitized with a resolution of 4° . However, we assert that the fitting utilized in the present paper constitutes a form of interpolation which yields better resolution. So, we henceforth compare the performance of the bootstrap particle filter which directly utilizes the discrete pseudo likelihood against the particle

filter utilizing the fitted vM mixture. Importantly, fitting with a vM mixture would be a prerequisite if the tracking was performed in the vein of [16].

B. Experimental results

Experiments were conducted in an acoustically prepared room, equipped with 3D pyramidal pattern studio foams placed on the roof and on the walls. Two surface microphones were mounted at the antipodes of a 8.9 cm radius plastic rigid sphere, itself place on a tripod. The two microphones outputs were synchronously acquired at 44.1 kHz. The sphere tripod was moved manually with a wheeled cart while the source, a loudspeaker placed at the same height as the microphones, was emitting voice recordings from a French radio programme. The true source and sensor positions were acquired at 200 Hz with a motion capture system, providing a less than 1 mm position error. For that purpose, small infrared active markers were placed on the sphere and the loudspeaker, and their signals were beamed to three infrared camera units placed at the corners of the room.

For the considered case of a rigid sphere, H_θ is shown to have the following analytic expression [20]

$$H_\theta(f) = \frac{\psi_{\frac{\pi}{2} + \theta}(f)}{\psi_{-\frac{\pi}{2} - \theta}(f)}, \quad \text{with} \quad (14)$$

$$\psi_\alpha(f) \triangleq \frac{1}{\left(\frac{2\pi f a}{c}\right)^2} \sum_{m=1}^{\infty} \frac{(-i)^{m-1} (2m+1) P_m(\cos \alpha)}{h'_m\left(\frac{2\pi f a}{c}\right)}.$$

Therein, ψ_β is the normalized head related transfer function (HRTF) to the microphone at angle β —with respect to boresight—on the sphere, where α stands for the angle between the source bearing and the direction to the considered microphone, P_m is the Legendre polynomial of degree m , h_m is the m th-order spherical Hankel function and h'_m is its first derivative. This expression was thus used in the pseudo likelihood computation. In practice, the infinite sum in (14) is approximated by a finite sum, the minimum order required to make the approximation reasonable depending on the maximum frequency considered. To avoid cumbersome computation during localization, H_θ was precomputed and stored offline for a discrete set of bearings.

In order to assess the performance of the PFs, we ran 50 Monte-Carlo runs on the sensed binaural data using either the discrete pseudo likelihood or the vM fitted pseudo likelihood. The runs were performed on four scenarios with different trajectories of the sensor, out of which one scenario included an intermittent sound source. In Fig. 3 we can see the results of range estimation for the four cases, while Fig. 4 shows the estimation of the bearing. By analyzing the results we can see that on average the PF with the vM fitted likelihood gave smaller error in terms of the range estimation although the performance in the bearing was similar for both PFs. The explanation lies in the fact that estimating the range from bearing-only measurements benefited from having an analytical likelihood compared to the 4° resolution of the discrete pseudo likelihood.

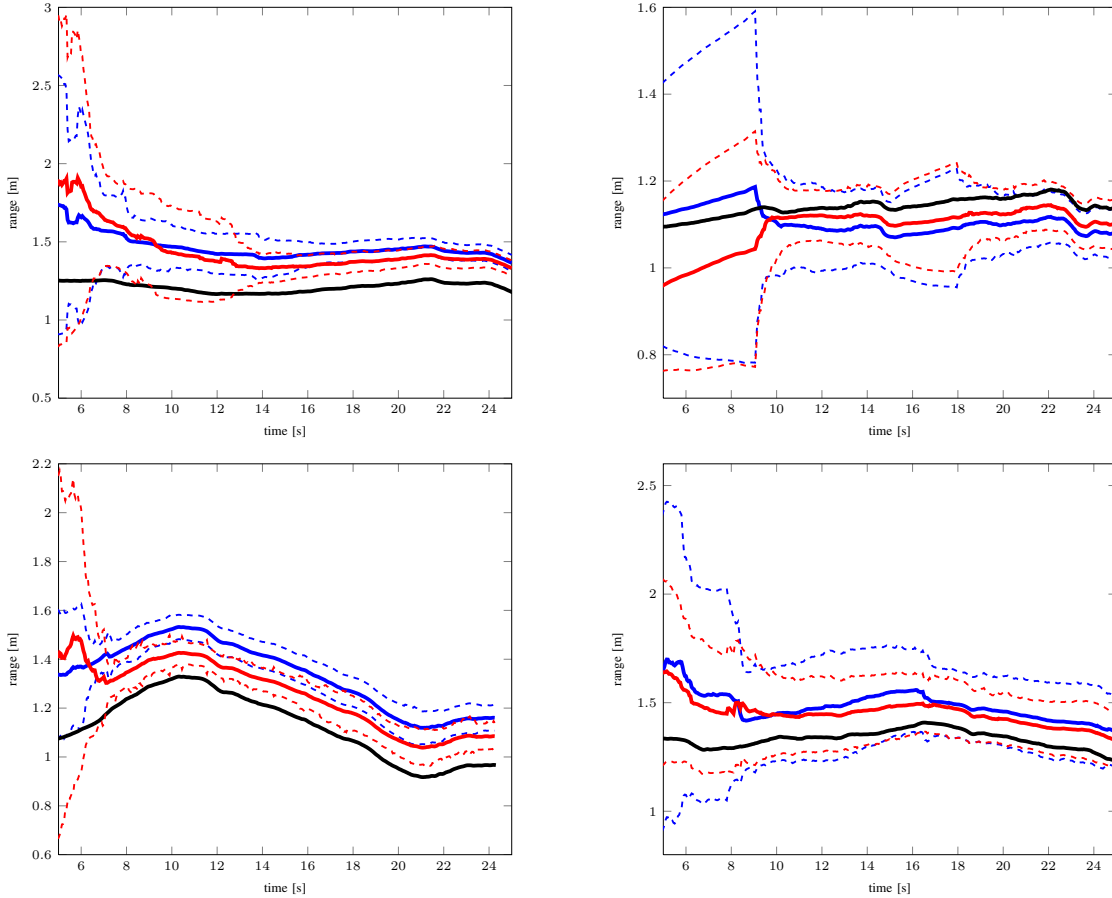


Fig. 3: Mean value of range estimates (solid) and pertaining three standard deviations (dashed) of 50 Monte-Carlo runs of the PF with pseudo likelihood (blue), VM fitted pseudo likelihood (red) and true range (black)

Then, for each entry of the posterior mean output by the filter, a minimum-width confidence interval for a moment matched Gaussian distribution of the estimation error was drawn (from the posterior particle set) which should approximately enclose the genuine hidden state vector with 99% probability. By analyzing the obtained plots concerning the range estimation error, we can see that the present implementation of the PF was not consistent over all the runs, since the true range is outside of the filters $\pm 3\sigma$ interval calculated from the estimated covariance matrix and that bias is present, which would indicate that the particle filter diverged at several instances of MC runs. This problem could be alleviated by utilizing a higher number of particles and/or more elaborate initialization and maneuvering strategies (cf. [21] for a deeper study of the problem).

V. CONCLUSION

In the present paper we have studied and proposed a solution for the problem of active speaker localization with a head mounted binaural microphone sensor. The solution was based on calculating a discrete pseudo likelihood function in speaker bearing based on the geometrical properties of the spherical head. The resulting likelihood was fitted with

a mixture of circular distributions, namely the VM and wrapped Cauchy distributions, whose comparison showed better results in the case of the VM distribution. A bootstrap algorithm was utilized with the direct and VM fitted pseudo likelihood in order to estimate the location of the speaker. We performed an experimental evaluation on four different data sets with accurate ground-truth, due to an active motion capture system, whose analysis showed that on average from 50 Monte Carlo runs both algorithms localized the speaker successfully, while the estimate with the VM fitted likelihood showed better accuracy in range and equal or better accuracy in bearing. However, a careful analysis revealed that at times the algorithms were inconsistent deviation-wise and that robust variants of the PF could be utilized, but which were eschewed in the present paper in order to guarantee the true posterior statistics. For future work, we will analyze different filtering schemes based on Gaussian filter mixtures and different members of the family of exponential distributions.

REFERENCES

- [1] V. Aidala and S. Hammel, "Utilization of modified polar coordinates for bearings-only tracking," *IEEE Trans. on Automatic Control*, vol. 28, no. 3, pp. 283–194, 1983.

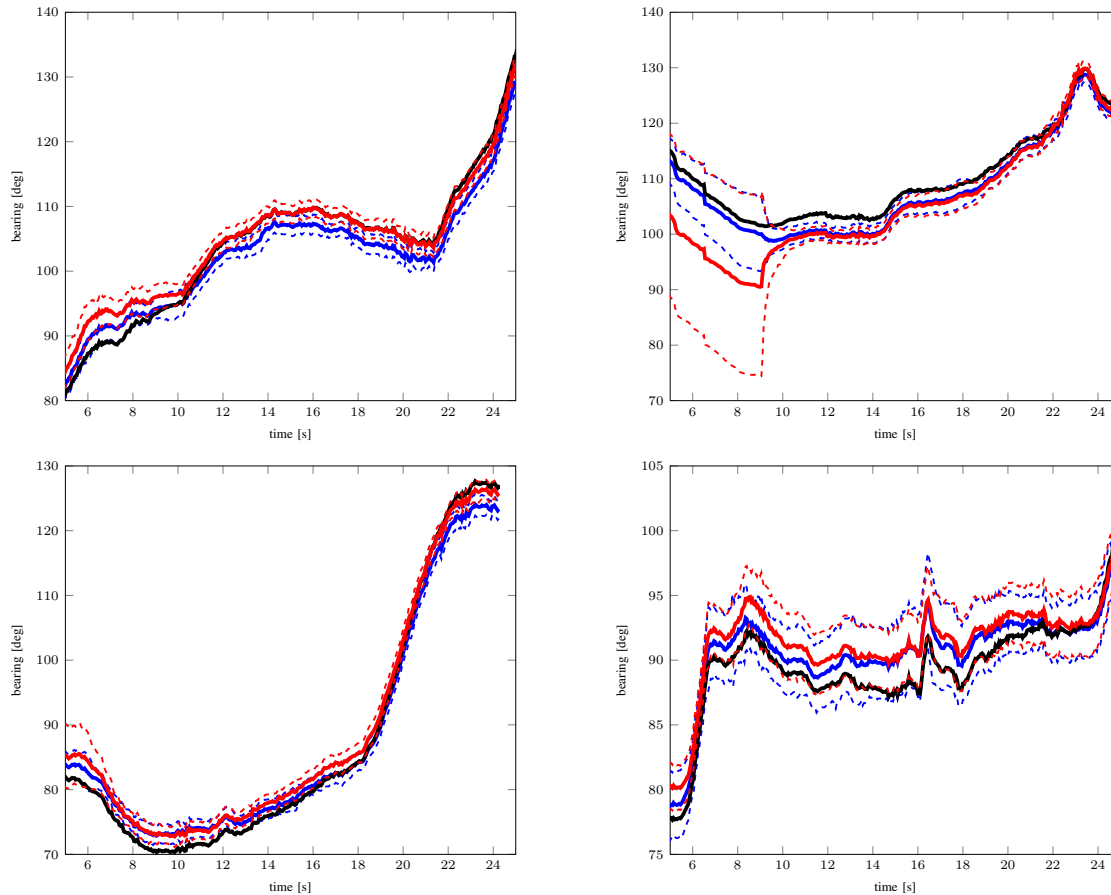


Fig. 4: Mean value of bearing estimates (solid) and pertaining three standard deviations (dashed) of 50 Monte-Carlo runs of the PF with pseudo likelihood (blue), VM fitted pseudo likelihood (red) and true bearing (black)

- [2] N. Peach, "Bearings-only tracking using a set of range-parameterised extended Kalman filters," *IEEE Proceedings – Control Theory and Applications*, vol. 142, no. 1, pp. 73–80, 1995.
- [3] B. Scala and M. Morelande, "An analysis of the single sensor bearings-only tracking problem," in *Int. Conf. on Information Fusion (FUSION'2008)*, 2008.
- [4] L.-L. Ong, "Non-Gaussian representations for decentralised bayesian estimation," Ph.D. dissertation, Univ. of Sydney, 2007.
- [5] A. Portello, P. Danès, and S. Argentieri, "Acoustic models and Kalman filtering strategies for active binaural sound localization," in *IEEE/RSJ Int. Conf. on Intell. Robots and Systems (IROS'2011)*.
- [6] J. Valin, F. Michaud, and J. Rouat, "Robust localization and tracking of simultaneous moving sound sources using beamforming and particle filtering," *Robotics and Autonomous Systems*, vol. 55, no. 3, pp. 216–228, 2007.
- [7] K. Okutani, T. Yoshida, K. Nakamura, and K. Nakadai, "Outdoor auditory scene analysis using a moving microphone array embedded in a quadcopter," in *IEEE/RSJ IROS'2012*.
- [8] M. Basiri, F. Schill, P. Lima, and D. Floreano, "Robust acoustic source localization of emergency signals from micro air vehicles," in *IEEE/RSJ Int. Conf. on Intell. Robots and Systems (IROS'2012)*.
- [9] I. Marković and I. Petrović, "Speaker localization and tracking with a microphone array on a mobile robot using von Mises distribution and particle filtering," *Robotics and Autonomous Systems*, vol. 58, no. 11, pp. 1185–1196, 2010.
- [10] D. Ward, E. Lehmann, and R. Williamson, "Particle filtering algorithms for tracking an acoustic source in a reverberant environment," *IEEE Trans. on Speech and Audio Processing*, vol. 11, no. 6, pp. 826–836, 2003.
- [11] Y.-C. Lu and M. Cooke, "Motion strategies for binaural localisation of speech sources in azimuth and distance by artificial listeners," *Speech Communication*, vol. 53, no. 5, pp. 622–642, 2010.
- [12] T. Habib and H. Romsdorfer, "Comparison of SRP-PHAT and multiband-PoPi algorithms for speaker localization using particle filters," in *Int. Conf. on Digital Audio Effects (DAFx-10)*.
- [13] S. Arulampalam, S. Maskell, N. Gordon, and T. Clapp, "A tutorial on particle filters for on-line non-linear/non-gaussian Bayesian tracking," *IEEE Transactions on Signal Processing*, vol. 50, pp. 174–188, 2001.
- [14] A. Portello, P. Danès, and S. Argentieri, "HRTF-based source azimuth estimation and activity detection from a binaural sensor," in *IEEE/RSJ Int. Conf. on Intelligent Robots and Systems (IROS'2013)*.
- [15] K. Mardia and P. Jupp, *Directional Statistics*. Wiley, 1999.
- [16] I. Marković and I. Petrović, "Bearing-only tracking with a mixture of von Mises distributions," in *IEEE/RSJ IROS'2012*.
- [17] A. Tanabe, K. Fukumizu, S. Oba, T. Takenouchi, and S. Ishii, "Parameter estimation for von Mises-Fisher distributions," *Computational Statistics*, vol. 22, no. 1, pp. 145–157, 2007.
- [18] D. Amos, "A portable package for Bessel functions of a complex argument and nonnegative order," *ACM Trans. on Math. Software*, 1986.
- [19] A. Doucet, N. de Freitas, and N. Gordon, "An introduction to sequential Monte Carlo methods," in *Sequential Monte Carlo Methods in Practice*. Springer-Verlag, 2001, ch. 1, pp. 3–14.
- [20] R. Duda and W. Martens, "Range dependence of the response of a spherical head model," *Jour. of the Acoustical Society of America*, vol. 104, pp. 3048–3058, 1998.
- [21] T. Brehard and J. Le Cadre, "Hierarchical particle filter for bearings-only tracking," *IEEE Trans. on Aerospace and Electronic Systems*, vol. 43, no. 4, pp. 1567–1585, 2007.