



# DoA Reliability for Distributed Acoustic Tracking

Christine Evers , *Senior Member, IEEE*, Emanuël A. P. Habets , *Senior Member, IEEE*, Sharon Gannot , *Senior Member, IEEE*, and Patrick A. Naylor , *Senior Member, IEEE*

**Abstract**—Distributed acoustic tracking estimates the trajectories of source positions using an acoustic sensor network. As it is often difficult to estimate the source-sensor range from individual nodes, the source positions have to be inferred from the direction-of-arrival (DoA) estimates. Due to reverberation and noise, the sound field becomes increasingly diffuse with increasing source-sensor distance, leading to a decreased Direction of Arrival (DoA)-estimation accuracy. To distinguish between accurate and uncertain DoA estimates, this letter proposes to incorporate the coherent-to-diffuse ratio as a measure of DoA reliability for single-source tracking. It is shown that the source positions, therefore, can be probabilistically triangulated by exploiting the spatial diversity of all nodes.

**Index Terms**—Acoustic sensors, Bayes' methods, smart homes.

## I. INTRODUCTION

**A**UTONOMOUS systems and smart devices rely on accurate knowledge of the positions of surrounding objects for human-machine interaction [1]. Acoustic scene mapping [2]–[4] provides three-dimensional (3-D) representations of the sound sources and microphone arrays in the surrounding environment. However, in realistic acoustic conditions, it is often difficult to accurately infer the positions of sound sources distant to a single microphone array [3]. In large enclosures, e.g., in smart homes, the networks of spatially distributed acoustic sensors can be exploited constructively for robust scene mapping [5]–[7]. Examples of network nodes include mobile phones [8], digital personal assistants [9], or robots [10].

For each node in the network, sound-source localization algorithms [11] estimate the instantaneous positional information of sources for a sequence of time frames. Devices exploited for ad hoc networks are typically equipped with compact microphone arrays, such that the range between the node's sensors and sources in the acoustic far field is difficult to determine. Sound-source localization, therefore, may only resort to estimates of the source direction-of-arrivals (DoAs) [12]–[14]. Furthermore, reverberation and noise often lead to missing and

false DoA estimates as well as estimation errors [15], [16]. Distributed acoustic tracking [17], [18] can be used to obtain smoothed source trajectories from the instantaneous DoA estimates by incorporating spatio-temporal models of the source motion. Within the Bayesian framework [19], the uncertainty in the source-motion model is traded off against the DoA reliability. However, most DoA estimators only provide point estimates of the source directions, but cannot quantify the reliability of the estimates [16], [20]. Therefore, source-tracking algorithms typically express the DoA reliability as a constant covariance matrix that is obtained by prior empirical experimentation [21], [22], supervised learning [23], or using the confidence intervals of DoA histograms [24], [25].

However, in practice, DoA-estimation accuracy is closely coupled to specific characteristics of each scenario, including the time-varying speech activity and source-sensor geometry. Reverberation [26] and ambient noise lead to diffuse noise fields [27], i.e., the acoustic energy approaches equal probability in every direction [28]. Furthermore, the energy in the direction of a source decreases with increasing source-sensor range. Therefore, the acoustic sound field at a microphone is decreasingly directional with increasing source-sensor range. The coupling between the DoA reliability and the sound field diffuseness is crucial for distributed sensor networks, where knowledge inferred from reliable nodes nearby a source must be distinguished from information at distant nodes.

The sound-field diffuseness can be quantified by the coherent-to-diffuse ratio (CDR) [29], [30]. The CDR was previously used to evaluate the speech-presence probability for blind speech separation in [31] and source extraction in [32]. In [33], the DoA estimates from a single microphone array were modeled by a von Mises (vM) distribution [34] whose concentration parameter is a function of the CDR.

This letter proposes a novel approach to distributed acoustic tracking that incorporates the CDR as a DoA reliability. Building on the model in [33] and using directional statistics [34], [35] we derive a Bayesian filter that distinguishes reliable DoA estimates at nearby nodes from DoA estimates at distant nodes. In contrast to [31], performing triangulation using instantaneous DoA estimates, the proposed tracking algorithm probabilistically triangulates the Cartesian source positions using the DoA estimates from all nodes within a centralized communication scheme. Node-specific information about the source directions is inferred from each node's DoA estimates. Assuming a synchronized network, estimates of the source position are obtained by fusing statistics of the node-specific information from the spatially diverse nodes.

Section II formulates the problem. Section III derives the proposed methodology. Section IV presents the experimental evaluation using realistic room simulations of a human talker. Conclusions are drawn in Section V.

Manuscript received March 12, 2018; revised May 11, 2018; accepted June 6, 2018. Date of publication June 21, 2018; date of current version July 30, 2018. This work was supported by the U.K. EPSRC Fellowship under Grant EP/P001017/1. The associate editor coordinating the review of this manuscript and approving it for publication was Prof. Jun Liu. (*Corresponding author: Christine Evers.*)

C. Evers and P. A. Naylor are with the Department of Electrical and Electronic Engineering, Imperial College London, London SW7 2AZ, U.K. (e-mail: c.evers@imperial.ac.uk; p.naylor@imperial.ac.uk).

E. A. P. Habets is with International Audio Laboratories Erlangen (a joint institution between the University of Erlangen-Nuremberg and Fraunhofer IIS), Erlangen 91058, Germany (e-mail: emanuel.habets@audiolabs-erlangen.de).

S. Gannot is with the Faculty of Engineering, Bar-Ilan University, Ramat Gan 5290002, Israel (e-mail: Sharon.Gannot@biu.ac.il).

Color versions of one or more of the figures in this letter are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/LSP.2018.2849579

## II. PROBLEM FORMULATION

### A. Likelihood of the DoA Estimates

Denote the 2-D Cartesian position of a sound source relative to node  $n = 1, \dots, N$  at position  $\mathbf{q}_n$  and time frame  $t$  as  $\mathbf{x}_{t,n} \triangleq [x_{t,n}, y_{t,n}]^T$ . The source position relative to node  $n$  can be obtained from  $\mathbf{x}_{t,\ell}$  relative to any node,  $\ell \neq n$ , as

$$\mathbf{x}_{t,n} = f_n(\mathbf{x}_{t,\ell}) \triangleq \mathbf{R}_{n|\ell} \mathbf{x}_{t,\ell} - \mathbf{q}_n + \mathbf{q}_\ell, \quad (1)$$

where  $\mathbf{R}_{n|\ell}$  is the rotation matrix. In the following we assume, without loss of generality, that the axes are orientated equally for all nodes, such that  $\mathbf{R}_{n|\ell}$  is a  $2 \times 2$  identity matrix.

The direction and range between the source and node,  $n$ , can be obtained from the source position using the Cartesian-to-spherical transformation, such that  $\mathbf{s}_{t,n} \triangleq [\phi_{t,n}, r_{t,n}]^T$  where  $\phi_{t,n} = \arctan(y_{t,n}/x_{t,n})$  is the source azimuth and  $r_{t,n} = \sqrt{x_{t,n}^2 + y_{t,n}^2}$  is the source-to-node range.

For distant sound sources, estimation of the range is difficult when using compact sensor arrays. Hence, at each array, estimates,  $\omega_{t,n} \in [-\pi, \pi)$ , of the true source DoA,  $\phi_{t,n} \in [-\pi, \pi)$ , are available, whereas the range,  $r_{t,n}$ , is unmeasured. The likelihood function of each DoA estimate,  $p(\omega_{t,n} | \mathbf{s}_{t,n}, \kappa_{t,n}^{(\text{DoA})})$ , can be modeled by a vM distribution [33], [34] with mean,  $\phi_{t,n}$ , and concentration parameter,  $\kappa_{t,n}^{(\text{DoA})} > 0$ :

$$\begin{aligned} p(\omega_{t,n} | \mathbf{s}_{t,n}, \kappa_{t,n}^{(\text{DoA})}) &= \mathcal{M}(\omega_{t,n} | \phi_{t,n}, \kappa_{t,n}^{(\text{DoA})}) \\ &\triangleq \frac{(2\pi)^{-1}}{I_0(\kappa_{t,n}^{(\text{DoA})})} \exp\{\kappa_{t,n}^{(\text{DoA})} \cos(\omega_{t,n} - \phi_{t,n})\}, \end{aligned} \quad (2)$$

where  $I_p(\kappa)$  is the modified Bessel function of the first kind and order  $p$ .

### B. Reliability Measure for DoA Estimates

Previous work showed that the accuracy of narrow-band DoA estimates is highly dependent on the CDR. The CDR is also known to be related to the source-sensor range and the reverberation time of the enclosure. In contrast to the range and reverberation time, the CDR can be estimated efficiently and accurately with a small microphone array using, e.g., [36].

In our model, the concentration of the vM distribution can be used to indicate the reliability of the DoA estimates. Therefore, as in [33], we use the CDR,  $\Gamma_{t,n}$ , to determine the concentration parameter,  $\kappa_{t,n}^{(\text{DoA})}$ , using

$$\kappa_{t,n}^{(\text{DoA})} = \ell_{\min} + (\ell_{\max} + \ell_{\min}) \frac{10^{c\varrho/10}}{10^{c\varrho/10} + [\Gamma_{t,n}]^\varrho}. \quad (3)$$

The terms  $\ell_{\min}$ ,  $\ell_{\max}$ ,  $c$ , and  $\varrho$ , respectively are the minimum and maximum concentration values, offset, and steepness of the transition region of the mapping from  $\Gamma_{t,n}$  to  $\kappa_{t,n}^{(\text{DoA})}$ .

The reliability of the DoA estimates, which is now reflected by the concentration of the vM distribution, can be used to distinguish between reliable DoA estimates from nearby nodes and less reliable DoA estimates from distant nodes.

### C. Early Fusion of $\omega_{t,n}$ and $\kappa_{t,n}^{(\text{DoA})}$ Across Frequencies

Performing estimation in the short time Fourier transform (STFT) domain results in one DoA and one CDR estimate per

frequency bin,  $k \in 1, \dots, K$ , time frame, and node. We assume that source tracking is provided with estimates per time frame and per node. The frequency-dependent estimates are, therefore, spectrally fused by evaluating their median values across all frequency bins, resulting in  $\omega_{t,n}$  and  $\kappa_{t,n}^{(\text{DoA})}$ .

### D. Source-Position Inference From DoA Estimates

Bayesian acoustic source tracking estimates the posterior probability density function (pdf),  $p(\mathbf{x}_{t,n_0} | \Omega_{1:t})$ , of the current source position,  $\mathbf{x}_{t,n_0}$ , relative to a reference point,  $\mathbf{q}_{n_0}$ , with index,  $n_0$ , from  $\Omega_{1:t} \triangleq \{\omega_{1:t}, \kappa_{1:t}^{(\text{DoA})}\}$ , where  $\omega_{1:t} \triangleq [\omega_{1,1}, \dots, \omega_{1,N}, \dots, \omega_{t,N}]^T$ , and  $\kappa_{1:t}^{(\text{DoA})}$  is structured similarly to  $\omega_{1:t}$ . A point estimate of the source position at  $t$  is given by the maximum *a posteriori* (MAP) estimate:

$$\hat{\mathbf{x}}_{t,n_0}^{\text{MAP}} = \arg \max_{\mathbf{x}_{t,n_0}} p(\mathbf{x}_{t,n_0} | \Omega_{1:t}). \quad (4)$$

However, 1)  $\kappa_{t,n}^{(\text{DoA})}$  is unknown in practice; 2) each node,  $n$ , only has access to its own DoA estimates,  $\omega_{1,n}, \dots, \omega_{t,n}$ , such that (4) cannot be evaluated directly; and 3) the 2-D source position,  $\mathbf{x}_{t,n_0}$ , must be inferred from the 1-D DoA estimates.

## III. PROPOSED METHODOLOGY

This section derives a novel approach for distributed source tracking that exploits (3) to distinguish reliable DoA estimates at nearby nodes from less reliable estimates at distant nodes.

### A. Network Fusion for Distributed Acoustic Tracking

In practice, each node only has access to its own DoA estimates. Hence, the source position is tracked by separately propagating a node-specific pdf in time using that node's DoA estimates. The statistics of the node-specific Probability Density Functions (pdfs) are fused within the network. For approaches that do not account for the reliability of DoA estimates, network fusion can lead to track divergence as knowledge inferred from nearby nodes cannot be distinguished from uncertain information from distant nodes.

In contrast to approaches agnostic to DoA-estimation reliability, exploiting the CDR in (3) ensures that the node-specific pdfs of nearby nodes correspond to peaked distributions, whereas the pdfs of distant nodes do not contain significant modes. Therefore, the node-specific source pdfs are fused within the network, relative to a reference point,  $\mathbf{x}_{t,n_0}$ , as:

$$p(\mathbf{x}_{t,n_0} | \Omega_{1:t}) \approx \prod_{n=1}^N p(\mathbf{x}_{t,n_0} | \omega_{t,n}, \kappa_{t,n}^{(\text{DoA})}, \Omega_{1:t-1}), \quad (5)$$

where the nodes are assumed to be independent due to their spatial diversity. By probability transformation of (1), we have

$$\begin{aligned} p(\mathbf{x}_{t,n_0} | \omega_{t,n}, \kappa_{t,n}^{(\text{DoA})}, \Omega_{1:t-1}) \\ = p(f_{n_0}(\mathbf{x}_{t,n}) | \omega_{t,n}, \kappa_{t,n}^{(\text{DoA})}, \Omega_{1:t-1}). \end{aligned}$$

### B. Probabilistic Triangulation for Range Inference

In order to estimate the node-specific pdfs required in (5), the density of the Cartesian source position must be updated from new information inferred from the DoA estimates. Contrary to bearing-only tracking, e.g., [37], we propose to infer the missing

**Algorithm 1:** PROST incorporating DoA reliability.

---

```

1: for  $t = 1, \dots, \infty$  do
2:   for  $n = 1, \dots, N$  do
3:     for  $k = 1, \dots, K$  do
4:       Estimate narrowband DoA,  $\omega_{t,n,k}$  [38];
5:       Estimate narrowband CDR,  $\Gamma_{t,n,k}$ , [36]
         and  $\kappa_{t,n,k}^{(\text{DoA})}$ , (3);
6:     end for
7:     Spectral fusion of  $\omega_{t,n}$  and  $\kappa_{t,n}^{(\text{DoA})}$ ,
         Section II-B;
8:     for  $j = 1, \dots, J$  do
9:       Transform  $\mathbf{x}_{t-1,n_0}^{(j)}$  to  $\mathbf{s}_{t-1,n}^{(j)}$ , (1);
10:      Predict  $\tilde{\mu}_{t,n}^{(j)}$ ,  $\tilde{\kappa}_{t,n}^{(j)}$ , (9);
11:      Update  $w_{t,n}^{(j)}$ ,  $\mu_{t,n}^{(j)}$ ,  $\kappa_{t,n}^{(j)}$ , (13);
12:      Transform  $\mathbf{s}_{t,n}^{(j)}$  to  $\mathbf{x}_{t,n_0}^{(j)}$ , (1);
13:    end for
14:  end for
15:  Network fusion, (5);
16:  Resample node-specific pdfs, Section III-B3;
17: end for

```

---

range estimates by exploiting the spatial diversity of the network nodes. As information is inferred from the DoA estimates, the node-specific pdfs must be propagated in polar coordinates as  $p(\mathbf{s}_{t,n} | \Omega_{1:t})$ . The posterior pdf of the polar coordinates can be estimated by: 1) predicting the pdf using a model of the source dynamics; and 2) updating the pdf with information inferred from the DoA estimates.

1) *Prediction:* The predicted pdf relative to node  $n$  is

$$p(\mathbf{s}_{t,n} | \Omega_{1:t-1}) = \int_{\mathbb{R}^3} p(\mathbf{s}_{t-1,n} | \Omega_{1:t-1}) p(\mathbf{s}_{t,n} | \mathbf{s}_{t-1,n}) d\mathbf{s}_{t-1,n}, \quad (6)$$

where the prior pdf,  $p(\mathbf{s}_{t,n} | \mathbf{s}_{t-1,n})$ , models the source dynamics. Assuming the DoA and range are independent:

$$p(\phi_{t,n} | \phi_{t-1,n}) = \mathcal{M}(\phi_{t,n} | \phi_{t-1,n}, \kappa_{t|t-1})$$

$$p(r_{t,n} | r_{t-1,n}) = \mathcal{U}(r_{t-1,n} - \rho_{t|t-1}, r_{t-1,n} + \rho_{t|t-1}), \quad (7)$$

where  $\mathcal{U}(\cdot)$  denotes the uniform distribution, and the transition concentration parameter and range,  $\kappa_{t|t-1}$  and  $\rho_{t|t-1}$ , respectively, capture the source dynamics between time frames.

The range pdf is approximated by importance sampling  $J$  hypotheses,  $\{\rho_{t,n}^{(j)}\}_{j=1}^J$ , of  $r_{t,n}$  from (7). Hence, the predicted pdf is given by the probability mass function (PMF):

$$p(\mathbf{s}_{t,n} | \Omega_{1:t-1}) \approx \sum_{j=1}^J w_{t-1,n}^{(j)} \delta_{\rho_{t,n}^{(j)}}(r_{t,n}) \mathcal{M}(\phi_{t,n} | \tilde{\mu}_{t,n}^{(j)}, \tilde{\kappa}_{t,n}^{(j)}), \quad (8)$$

where  $\delta_a$  is the Dirac measure at a given state  $a$ , and the predicted mean  $\tilde{\mu}_{t,n}^{(j)}$  and covariance  $\tilde{\kappa}_{t,n}^{(j)}$  are given by (see Supplementary material, Section I-A)

$$\tilde{\mu}_{t,n}^{(j)} = \mu_{t-1,n}^{(j)}, \quad \tilde{\kappa}_{t,n}^{(j)} = A^{-1}(A(\kappa_{t-1,n})A(\kappa_{t|t-1})), \quad (9)$$

with  $A(\kappa) \triangleq I_1(\kappa)/I_0(\kappa)$  [34], and where, for any  $\bar{r}$  [35]

$$A^{-1}(\bar{r}) \approx \frac{2\bar{r} - \bar{r}^3}{1 - \bar{r}^2}. \quad (10)$$

2) *Update:* Using Bayes' theorem:

$$p(\mathbf{s}_{t,n} | \omega_{t,n}, \kappa_{t,n}^{(\text{DoA})}, \Omega_{1:t-1}) = \frac{p(\omega_{t,n} | \mathbf{s}_{t,n}, \kappa_{t,n}^{(\text{DoA})}) p(\mathbf{s}_{t,n} | \Omega_{1:t-1})}{\int p(\omega_{t,n} | \mathbf{s}_{t,n}, \kappa_{t,n}^{(\text{DoA})}) p(\mathbf{s}_{t,n} | \Omega_{1:t-1}) d\mathbf{s}_{t,n}}. \quad (11)$$

Substituting (8) and (2) into (11):

$$p(\mathbf{s}_{t,n} | \omega_{t,n}, \kappa_{t,n}^{(\text{DoA})}, \Omega_{1:t}) = \sum_{j=1}^J w_{t,n}^{(j)} \delta_{\rho_{t,n}^{(j)}}(r_{t,n}) \mathcal{M}(\phi_{t,n} | \mu_{t,n}^{(j)}, \kappa_{t,n}^{(j)}), \quad (12)$$

where the updated mean,  $\mu_{t,n}^{(j)}$ , concentration,  $\kappa_{t,n}^{(j)}$ , and weights  $w_{t,n}^{(j)}$ , are (see supplementary material, Section I-B)<sup>1</sup>

$$\mu_{t,n}^{(j)} = -\tan^{-1} \left( \frac{\tilde{\kappa}_{t,n}^{(j)} \sin \tilde{\mu}_{t,n}^{(j)} + \kappa_{t,n}^{\text{DoA}} \sin \omega_{t,n}}{\tilde{\kappa}_{t,n}^{(j)} \cos \tilde{\mu}_{t,n}^{(j)} + \kappa_{t,n}^{\text{DoA}} \cos \omega_{t,n}} \right) \quad (13a)$$

$$\kappa_{t,n}^{(j)} = \sqrt{[\tilde{\kappa}_{t,n}^{(j)}]^2 + [\kappa_{t,n}^{\text{DoA}}]^2 + 2\tilde{\kappa}_{t,n}^{(j)} \kappa_{t,n}^{\text{DoA}} (\cos(\tilde{\mu}_{t,n}^{(j)} - \omega_{t,n}))} \quad (13b)$$

$$w_{t,n}^{(j)} = w_{t-1,n}^{(j)} \frac{\tilde{w}_{t,n}^{(j)}}{\sum_{\ell=1}^J \tilde{w}_{t,n}^{(\ell)}} \quad \text{and} \quad \tilde{w}_{t,n}^{(j)} = \frac{I_0(\kappa_{t,n}^{(j)})}{I_0(\tilde{\kappa}_{t,n}^{(j)})}. \quad (13c)$$

3) *Network Fusion:* Due to the product in (5), network fusion requires a continuous representation of the discrete PMF in (12). A continuous approximation,  $\hat{p}(\mathbf{x}_{t,n_0} | \Omega_{1:t})$ , of each node-specific pdf is obtained from the weighted kernel density estimate (KDE) [39]. The resulting Kernel Density Estimates (KDEs) are used to evaluate (5). For propagation to  $t+1$ , the node-specific PMFs resulting from the network fusion are extracted by randomly selecting, with equal probability,  $J$  components of the nonfused PMFs in (12) from within the intersection region in (5).

Although the unmeasured range is only predicted in (11), resampling from the network posterior ensures that only stochastically relevant range hypotheses are propagated. Hence, the source positions are probabilistically triangulated by exploiting the spatial diversity of the nodes. PRObabilistic Source Triangulation (PROST) is summarized in Algorithm 1.

## IV. PERFORMANCE EVALUATION

### A. Experimental Setup

Simulations of a  $10 \times 7 \times 2.5$  m<sup>3</sup> room at 16 kHz sampling frequency are used for evaluation. The positions of a source moving at 0.5 m/s are generated at 32 sample intervals along a straight-line trajectory. The distributed sensor network contains four circular microphone arrays, each with three microphones spaced by 0.025 m. The array positions are simulated for a midrange scenario (S1) and a far-range scenario (S2) as detailed in Table I. The room impulse responses (RIRs) for each source-sensor configuration are simulated using the image-source method [40] for a reverberation time [26] of 0.5 s. The

<sup>1</sup>A vM filter for bearing-only tracking was previously proposed in [37]. It is important to note that the additive term in [37, eq. (8)] may result in  $|\mu_{t,n}| > \pi$ . In contrast, by definition of the inverse tangent, (13a) ensures valid directions within the region of support  $\mu_{t,n} \in [-\pi, \pi]$ .



TABLE I  
AVERAGE TRACKING ACCURACY FOR BOTH SCENARIOS S1 AND S2, AND VOICE ACTIVITY PERIODS P1 AND P2.

	$q_1$ [m]	$q_2$ [m]	$q_3$ [m]	$q_4$ [m]	LS				PROST-const				PROST			
					Mean [m]		Std [m]		Mean [m]		Std [m]		Mean [m]		Std [m]	
					P1	P2	P1	P2	P1	P2	P1	P2	P1	P2	P1	P2
S1	6, 5.5, 1.8	8, 3, 1.8	4, 2, 1.8	2, 2.5, 1.8	0.63	0.68	0.44	0.5	0.48	0.26	<b>0.1</b>	0.13	<b>0.38</b>	<b>0.18</b>	0.12	<b>0.06</b>
S2	6, 5.5, 1.8	8, 2, 1.8	4, 1, 1.8	2, 2.5, 1.8	0.97	1.12	0.66	0.65	<b>0.46</b>	0.52	0.17	0.18	0.5	<b>0.32</b>	<b>0.1</b>	<b>0.12</b>

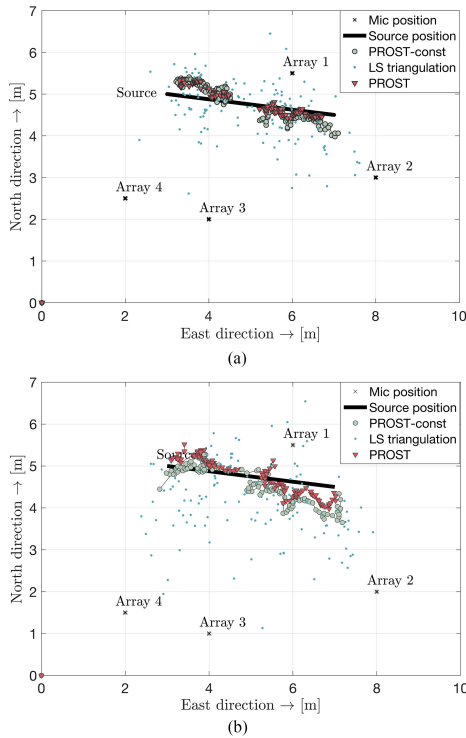


Fig. 1. Birdseye view of the acoustic scene map for: (a) Midrange. (b) Far-range scenarios. (a) Midrange scenario (S1). (b) Far-range scenario (S2).

Room Impulse Responses (RIRs) are convolved with anechoic speech of 8 s duration, consisting of two utterances by a female talker. The convolved signals are distorted by sensor and isotropic noise [41] with signal-to-noise ratios of 20 dB and 40 dB, respectively. The STFT is evaluated for 64 ms frame lengths, 1024 discrete Fourier transform points, and using a Hamming window with 25% overlap. The voice activity detector (VAD) in [42] is evaluated for a window length of 10 ms. Due to the compact array aperture, the propagation delay between microphones for each array is negligible in the context of Voice Activity Detector (VAD). The VAD is, hence, evaluated only for channel 3 at each node. The VAD results are combined within the network by evaluating the maximum start time and minimum duration of each voice activity period (VAP) across all nodes. During VAPs, one CDR and one DoA per time frame and frequency bin are estimated. The DoAs are estimated by minimizing the discrepancy between the observed and expected intermicrophone phase differences, as proposed in [38]. The CDR is estimated as a function of the spatial coherence between microphones 1 and 2 for each node using the approach in [36]. The corresponding estimates per frame are evaluated by the median across all frequency bins.  $\kappa_{t,n}^{(\text{DoA})}$  is obtained using (3) with  $l_{\min} = 0$ ,  $l_{\max} = 25$ ,  $c = 6$ , and  $\rho = -2$ . Algorithm 1 is initialized for each VAP by drawing  $J = 350$

source-position hypotheses for each node. The initial hypotheses are sampled along each  $\omega_{t,n}$  between  $q_n$  and the intersection of the DoA vector with the room boundary. For all subsequent frames, Algorithm 1 is evaluated using the transition parameters  $\rho_{t|t-1} = 1.5$  m and  $\kappa_{t|t-1} = 500$ , chosen to enforce confidence in the dynamical model for robustness against false DoA estimates.

The performance of PROST is compared against: 1) PROST-const, which uses a constant  $\kappa_{t,n}^{(\text{DoA})} = 5$ , corresponding to the mean of  $\kappa_{t,n}^{(\text{DoA})}$  estimates during the first utterance; and 2) least squares triangulation (LST) from the DoA estimates.

## B. Results

1) *Midrange Scenario*: The VAD detects two VAPs, P1 and P2. The DoA estimates correspond to average errors of  $10^\circ$  for P1 and  $14^\circ$  for P2. Fig. 1(a) compares the ground truth source positions against the estimates of PROST and the two benchmarks. Least Squares Triangulation (LST) produces a large-volume cloud of points. By tracking the source position across time, PROST-const results in tracks with large variation between positions. PROST results in smooth tracks near the source by incorporating the DoA reliability, achieving improvements in position accuracy of up to 0.5 m over LST and 0.1 m over PROST-const (see Table I).

2) *Far-Range Scenario*: Fig. 1(b) shows the estimated scene map for the far-range scenario. For the first VAP, all four arrays initially correspond to large source-sensor distances of over 4.5 m, such that PROST-const and PROST lead to comparable performance results. In the second VAP, Array 1 is within 1.04 m, leading to reliable DoA estimates and, hence, increased concentration parameters compared to the first VAP. PROST leads to an improvement of 0.2 and 0.8 m compared to PROST-const and LST, respectively. Hence, the benefits of incorporating the DoA reliability increase for far-range scenarios, where at least one node is nearby the source.

## V. CONCLUSION

We proposed to exploit the CDR as the DoA reliability for distributed acoustic tracking. The CDR is incorporated as the concentration parameter of the DoA-likelihood function, modeled by a vM distribution. The source DoAs are then tracked in time at each individual node using a vM filter. To infer the unmeasured source-sensor range, the vM filter is evaluated for a cloud of uninformative range hypotheses. By network fusion, spatial diversity of the nodes is exploited in order to probabilistically triangulate the stochastically relevant source positions, and hence, range hypotheses. Realistic simulation results demonstrate improvements of up to 39% compared with the classical approach of a constant concentration parameter, and up to 74% compared with least-squares source triangulation.

## REFERENCES

- [1] R. Parasuraman, T. B. Sheridan, and C. D. Wickens, "A model for types and levels of human interaction with automation," *IEEE Trans. Syst., Man, Cybern. A, Syst., Humans*, vol. 30, no. 3, pp. 286–297, May 2000.
- [2] C. Evers and P. A. Naylor, "Optimized self-localization for SLAM in dynamic scenes using probability hypothesis density filters," *IEEE Trans. Signal Process.*, vol. 66, no. 4, pp. 863–878, Feb. 2018.
- [3] C. Evers and P. A. Naylor, "Acoustic SLAM," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 26, no. 9, pp. 1484–1498, Sep. 2017.
- [4] C. Evers, A. H. Moore, and P. A. Naylor, "Acoustic simultaneous localization and mapping (a-SLAM) of a moving microphone array and its surrounding speakers," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, Shanghai, China, Mar. 2016, pp. 6–10.
- [5] C. Meesookho, U. Mitra, and S. Narayanan, "On energy-based acoustic source localization for sensor networks," *IEEE Trans. Signal Process.*, vol. 56, no. 1, pp. 365–377, Jan. 2008.
- [6] A. Canclini, F. Antonacci, A. Sarti, and S. Tubaro, "Acoustic source localization with distributed asynchronous microphone networks," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 21, no. 2, pp. 439–443, Feb. 2013.
- [7] Y. Dorfan, A. Plinge, G. Hazan, and S. Gannot, "Distributed expectation-maximization algorithm for speaker localization in reverberant environments," *IEEE/ACM Trans. Audio, Speech, Lang. Process.*, vol. 26, no. 3, pp. 682–695, Mar. 2018.
- [8] M. H. Hennecke and G. A. Fink, "Towards acoustic self-localization of ad hoc smartphone arrays," in *Proc. Joint Workshop Hands-Free Speech Commun. Microphone Arrays*, May 2011, pp. 127–132.
- [9] H. T. Sullivan and S. Sahasrabudhe, "Envisioning inclusive futures: Technology-based assistive sensory and action substitution," *Futures*, vol. 87, pp. 140–148, 2017.
- [10] C. Evers, A. H. Moore, P. A. Naylor, J. Sheaffer, and B. Rafaely, "Bearing-only acoustic tracking of moving speakers for robot audition," in *Proc. IEEE Int. Conf. Digital Signal Process.*, Jul. 2015, pp. 1206–1210.
- [11] S. Argentieri, P. Danès, and P. Souères, "A survey on sound source localization in robotics: From binaural to array processing methods," *Comput. Speech Lang.*, vol. 34, no. 1, pp. 87–112, 2015.
- [12] J. C. Chen, K. Yao, and R. E. Hudson, "Acoustic source localization and beamforming: Theory and practice," *EURASIP J. Adv. Signal Process.*, vol. 2003, no. 4, Mar. 2003, Art. no. 926837.
- [13] A. Griffin, A. Alexandridis, D. Pavlidi, Y. Mastorakis, and A. Mouchtaris, "Localizing multiple audio sources in a wireless acoustic sensor network," *Signal Process.*, vol. 107, pp. 54–67, 2015.
- [14] A. Canclini, F. Antonacci, A. Sarti, and S. Tubaro, "Distributed 3D source localization from 2D DOA measurements using multiple linear arrays," *Wireless Commun. Mobile Comput.*, vol. 2017, 2017, Art. no. 1049141.
- [15] D. Jarrett, E. A. P. Habets, and P. A. Naylor, *Theory and Applications of Spherical Microphone Array Processing*. New York, NY, USA: Springer-Verlag, 2016.
- [16] A. H. Moore, C. Evers, and P. A. Naylor, "Direction of arrival estimation in the spherical harmonic domain using subspace pseudointensity vectors," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 25, no. 1, pp. 178–192, Jan. 2017.
- [17] O. Hlinka, F. Hlawatsch, and P. M. Djurić, "Consensus-based distributed particle filtering with distributed proposal adaptation," *IEEE Trans. Signal Process.*, vol. 62, no. 12, pp. 3029–3041, Jun. 2014.
- [18] A. Mohammadi and A. Asif, "Distributed consensus + innovation particle filtering for bearing/range tracking with communication constraints," *IEEE Trans. Signal Process.*, vol. 63, no. 3, pp. 620–635, Feb. 2015.
- [19] B. Ristic, S. Arulampalam, and N. Gordon, *Beyond the Kalman filter: Particle Filters for Tracking Applications*. Norwood, MA, USA: Artech House, 2004.
- [20] O. Nadiri and B. Rafaely, "Localization of multiple speakers under high reverberation using a spherical microphone array and the direct-path dominance test," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 22, no. 10, pp. 1494–1505, Oct. 2014.
- [21] X. Zhong and J. R. Hopgood, "A time-frequency masking based random finite set particle filtering method for multiple acoustic source detection and tracking," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 23, no. 12, pp. 2356–2370, Dec. 2015.
- [22] M. F. Fallon and S. J. Goddard, "Acoustic source localization and tracking of a time-varying number of speakers," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 20, no. 4, pp. 1409–1415, May 2012.
- [23] B. Laufer-Goldshtein, R. Talmon, and S. Gannot, "A hybrid approach for speaker tracking based on TDOA and data-driven models," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 26, no. 4, pp. 725–735, Apr. 2018.
- [24] J. Nikunen, A. Diment, and T. Virtanen, "Separation of moving sound sources using multichannel NMF and acoustic tracking," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 26, no. 2, pp. 281–295, Feb. 2018.
- [25] C. Evers, B. Rafaely, and P. A. Naylor, "Speaker tracking in reverberant environments using multiple directions of arrival," in *Proc. Joint Workshop Hands-Free Speech Commun. Microphone Arrays*, Mar. 2017, pp. 91–95.
- [26] P. A. Naylor and N. D. Gaubitch, Eds., *Speech Dereverberation*. New York, NY, USA: Springer-Verlag, 2010.
- [27] H. Kuttruff, *Room Acoustics, 5th ed.* Boca Raton, FL, USA: CRC Press, Jun. 2009.
- [28] R. K. Cook, R. V. Waterhouse, R. D. Berendt, E. Seymour, and M. C. Thompson, "Measurement of correlation coefficients in reverberant sound fields," *J. Acoust. Soc. Amer.*, vol. 27, no. 6, pp. 1072–1077, 1955.
- [29] M. Jeub, C. Nelke, C. Beaugeant, and P. Vary, "Blind estimation of the coherent-to-diffuse energy ratio from noisy speech signals," in *Proc. Eur. Signal Process. Conf.*, Aug 2011, pp. 1347–1351.
- [30] S. Braun *et al.*, "Evaluation and comparison of late reverberation power spectral density estimators," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 26, no. 6, pp. 1056–1071, Jun. 2018.
- [31] M. Taseska and E. A. P. Habets, "Blind source separation of moving sources using sparsity-based source detection and tracking," *IEEE/ACM Trans. Audio, Speech, Language Process.*, vol. 26, no. 3, pp. 657–670, Mar. 2018.
- [32] M. Taseska and E. A. P. Habets, "DOA-informed source extraction in the presence of competing talkers and background noise," *EURASIP J. Adv. Signal Process.*, vol. 2017, no. 1, p. 60, Aug. 2017.
- [33] M. Taseska and E. A. P. Habets, "Minimum Bayes risk signal detection for speech enhancement based on a narrowband DOA model," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, Apr. 2015, pp. 539–543.
- [34] K. V. Mardia and P. E. Jupp, *Directional Statistics*. Hoboken, NJ, USA: Wiley, 1999.
- [35] A. Banerjee, I. S. Dhillon, J. Ghosh, and S. Sra, "Clustering on the unit hypersphere using von Mises-Fisher distributions," *J. Mach. Learn. Res.*, vol. 6, pp. 1345–1382, Dec. 2005.
- [36] O. Thiergart, G. D. Galdo, and E. Habets, "On the spatial coherence in mixed sound fields and its application to signal-to-diffuse ratio estimation," *J. Acoust. Soc. Amer.*, vol. 132, no. 4, pp. 2337–2346, Dec. 2012.
- [37] I. Marković, J. Česić, and I. Petrović, "Von Mises mixture PHD filter," *IEEE Signal Process. Lett.*, vol. 22, no. 12, pp. 2229–2233, Dec. 2015.
- [38] O. Thiergart, W. Huang, and E. A. P. Habets, "A low complexity weighted least squares narrowband DOA estimator for arbitrary array geometries," in *Proc. IEEE Int. Conf. Acoust., Speech Signal Process.*, Mar. 2016, pp. 340–344.
- [39] F. J. G. Gisbert, "Weighted samples, kernel density estimators and convergence," *Empirical Econ.*, vol. 28, no. 2, pp. 335–351, Apr. 2003.
- [40] J. B. Allen and D. A. Berkley, "Image method for efficiently simulating small-room acoustics," *J. Acoust. Soc. Amer.*, vol. 65, no. 4, pp. 943–950, Apr. 1979.
- [41] E. A. P. Habets and S. Gannot, "Generating sensor signals in isotropic noise fields," *J. Acoust. Soc. Amer.*, vol. 112, no. 6, pp. 3464–3470, Dec. 2007.
- [42] J. Sohn, N. S. Kim, and W. Sung, "A statistical model-based voice activity detection," *IEEE Signal Process. Lett.*, vol. 6, no. 1, pp. 1–3, Jan. 1999.