



HAL
open science

Online Localization and Tracking of Multiple Moving Speakers in Reverberant Environments

Xiaofei Li, Yutong Ban, Laurent Girin, Xavier Alameda-Pineda, Radu Horaud

► **To cite this version:**

Xiaofei Li, Yutong Ban, Laurent Girin, Xavier Alameda-Pineda, Radu Horaud. Online Localization and Tracking of Multiple Moving Speakers in Reverberant Environments. *IEEE Journal of Selected Topics in Signal Processing*, 2019. hal-01851985v1

HAL Id: hal-01851985

<https://inria.hal.science/hal-01851985v1>

Submitted on 31 Jul 2018 (v1), last revised 1 Mar 2019 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Online Localization and Tracking of Multiple Moving Speakers in Reverberant Environment

Xiaofei Li*, Yutong Ban* Laurent Girin, Xavier Alameda-Pineda and Radu Horaud

Abstract—This paper addresses the problem of online multiple-speaker localization and tracking in reverberant environment. We propose to use the direct-path relative transfer function (DP-RTF) – a feature that encodes the inter-channel direct-path information robust against reverberation, hence well suited for reliable localization. A complex Gaussian mixture model (CGMM) is then used, such that each component weight represents the probability that an active speaker is present at a corresponding candidate source direction. Exponentiated gradient descent is used to update these weights online by minimizing a combination of negative log-likelihood and entropy. The latter imposes sparsity over the number of audio sources, since in practice only a few speakers are simultaneously active. The outputs of this online localization process are then used as observations within a Bayesian filtering process whose computation is made tractable via an instance of variational expectation-maximization. Birth and sleeping processes are used to account for the intermittent nature of speech. The method is thoroughly evaluated using several datasets.

Index Terms—Sound-source localization, multiple moving speakers, online tracking, reverberant environments.

I. INTRODUCTION

Online multiple-speaker localization and tracking are challenging tasks in real world environments and natural conversation scenarios, with reverberation and ambient noise, short sentences, speech pauses and frequent speech turns among speakers. Time difference of arrival (TDOA) is widely used for single source localization [1]. Many TDOA estimators exist, such as the generalized cross-correlation method [2]. For multiple-speaker localization, beamforming methods, e.g. steered-response power (SRP) [3], and subspace methods, e.g. multiple signal classification (MUSIC) [4], are widely used. The W-disjoint orthogonality (WDO) [5] assumes that the audio signal is dominated by only one speaker in each small region of the time-frequency (TF) domain, because of the natural sparsity of speech signals in this domain. Applying the short-time Fourier transform (STFT), or any TF decomposition, interchannel localization features, e.g. interaural phase difference (IPD) [5], can be extracted. In [5], multiple-speaker localization is based on the histogram of interchannel features, which however is suitable only in the case where

there is no wrapping of phase measures. In [6], a mixture of Gaussian mixture models (GMMs) is used as a generative model of the interchannel features of multiple speakers, with each GMM representing one speaker, and each GMM component representing one candidate interchannel time delay. An expectation-maximization (EM) algorithm iteratively estimates the component weights and assigns the features to their corresponding candidate time delays. This method overcomes the phase ambiguity problem by jointly considering all frequencies in the likelihood maximization procedure. After maximizing the likelihood, the azimuth of each speaker is given by the component that has the highest weight in the corresponding GMM. The complex-valued version of IPD, i.e. the pairwise relative phase ratio (PRP), is used in [7]. In addition, instead of setting one GMM for each speaker, a single complex GMM (CGMM) is used for all speakers with each component representing one candidate speaker location. After maximizing the likelihood of the PRP features, with an EM algorithm, the weight of each component represents the probability that there is an active speaker at the corresponding candidate location. Therefore, for an unknown number of speakers, counting and localization of active speakers can be jointly carried out by selecting the components with large weights.

The interchannel features and localization methods mentioned above assume direct-path propagation model, hence are poorly suited for reverberant environments. To overcome this, several TDOA estimators based on system identification were proposed in [8]–[11]. In [12], it is proposed to use the direct-path relative transfer function (DP-RTF) as a TF-domain localization feature robust against reverberation. The DP-RTF feature estimation is based on the identification of the STFT-domain representation of the room impulse response (RIR), i.e. the convolutive transfer function (CTF) [13], [14]. Overall, this method combines the merits of the robust TDOA estimators [8]–[11] and of the TF-domain WDO assumption, and thus allows multiple-speaker localization in reverberant environments.

To localize moving speakers, one-stage methods such as SRP and MUSIC can be directly used by applying frame-wise spatial spectrum estimation. In contrast, for the interchannel feature based methods, one needs to assign the frame-wise features to the speakers in an adaptive/recursive way, such as the smoothed histogram used in [15]. A CGMM model similar to [7] (but with one CGMM per predefined speaker) was used in [16] and plugged into a recursive EM algorithm to update online the CGMM component weights.

Taking the instantaneous outputs of a localization method

* Xiaofei Li and Yutong Ban contributed equally to this work.

Xiaofei Li, Yutong Ban, Xavier Alameda-Pineda and Radu Horaud are with INRIA Grenoble Rhône-Alpes, Montbonnot Saint-Martin, France. E-mail: first.last@inria.fr

Laurent Girin is with INRIA Grenoble Rhône-Alpes and with Univ. Grenoble Alpes, Grenoble-INP, GIPSA-lab, France. E-mail: laurent.girin@gipsa-lab.grenoble-inp.fr

This work was supported by the ERC Advanced Grant VHIA #340113.

as observations and using a speaker dynamic model, Bayesian tracking techniques estimate the posterior distribution of source directions, e.g. [17]–[19]. An improved Kalman filter and particle filtering were used for single sound source tracking in [20] and [21], respectively. To tackle the tracking problem with an unknown and time-varying number of speakers, additional model features such as observation-to-speaker assignments, speaker track birth/death processes and a model of speech activity can be included [22], [23]. Sampling-based tracking methods are also widely used for this aim, e.g. extended particle filtering frameworks [24]–[26] or sequential Monte Carlo implementation of the Probability Hypothesis Density (PHD) filter [27], [28]. However, in a general manner, the computation cost of sampling-based tracking methods can be very high. Under some assumptions, the Gaussian mixture PHD filter for multi-target tracking [29] has an analytical solution, and is computationally efficient. This filter was adopted for multiple-speaker tracking in [18]. In this paper we propose a new method for online localization and tracking of multiple moving speakers. The paper has the following contributions:

- First, CTF estimates are used for DP-RTF extraction. Unlike the batch formulation [12] which assumes static speakers, we adopt the online CTF estimation framework of [30] based on recursive least-squares (RLS). The latter has a better convergence rate, a crucial feature when dealing with moving speakers, than the least mean-squares (LMS) algorithms presented in [8], [9].
- Second, for frame-wise localization we adopt the CGMM model of [7] to assign DP-RTF features to speakers. Instead of using the recursive EM algorithm [16] to update the CGMM weights [30], we propose to use exponentiated gradient (EG) [31], which efficiently integrates likelihood maximization and entropy regularization.
- Finally, we propose a Bayesian multiple-speaker tracking method and an associated computational tractable solver, namely an instance of the variational expectation-maximization (VEM) algorithm. One novelty of the proposed algorithm is the embedding of birth and sleeping processes which enable to consider a time-varying number of active speakers.

The paper is organized as follows. Section II describes the online recursive DP-RTF estimation method. Section III presents the CGMM and EG based localization method. Section IV presents the Bayesian tracker based on VEM. Section V presents the experiments we carried out to evaluate the proposed method, using two “real world” datasets. Section VI concludes the paper.

II. RECURSIVE MULTICHANNEL DP-RTF ESTIMATION

A. Recursive Least Squares

For the sake of clarity, we first consider the noise-free single-speaker case. In the time domain $x^i(n) = a^i(n) \star s(n)$ is the i -th microphone signal, $i = 1, \dots, I$, where n is the time

index, $s(n)$ is the source signal, $a^i(n)$ is the RIR from the source to the i -th microphone, and \star denotes the convolution. Applying the STFT and using the CTF approximation, for each frequency index $f = 0, \dots, F - 1$ we have:

$$x_{t,f}^i = a_{t,f}^i \star s_{t,f} = \sum_{q=0}^{Q-1} a_{q,f}^i s_{t-q,f}, \quad (1)$$

where $x_{t,f}^i$ and $s_{t,f}$ are the STFT coefficients of the corresponding signals, and the CTF $a_{t,f}^i$ is a subband representation of $a^i(n)$. Here, the convolution is executed with respect to the frame index t . The number of CTF coefficients Q is related to the reverberation time of the RIR. The first CTF coefficient $a_{0,f}^i$ mainly consists of the direct-path information, thence the DP-RTF is defined as the ratio between the first CTF coefficients of two channels: $a_{0,f}^i/a_{0,f}^r$, where channel r is the reference channel.

Based on the cross-relation method [32], using the CTF model of one microphone pair (i, j) , we have: $x_{t,f}^i \star a_{t,f}^j = x_{t,f}^j \star a_{t,f}^i$. This can be written in vector form as:

$$\mathbf{x}_{t,f}^{i\top} \mathbf{a}_f^j = \mathbf{x}_{t,f}^{j\top} \mathbf{a}_f^i, \quad (2)$$

with $\mathbf{a}_f^i = [a_{0,f}^i, \dots, a_{Q-1,f}^i]^\top$, where \top denotes matrix/vector transpose, and $\mathbf{x}_{t,f}^i = [x_{t,f}^i, \dots, x_{t-Q+1,f}^i]^\top$. The CTF vector involving all channels is defined as $\mathbf{a}_f = [\mathbf{a}_f^{1\top}, \dots, \mathbf{a}_f^{I\top}]^\top$. There is a total of $I(I-1)/2$ distinct microphone pairs, indexed by (i, j) with $i = 1, \dots, I-1$ and $j = i+1, \dots, I$. For each pair, we construct a cross-relation equation in terms of \mathbf{a}_f . For this aim, we define:

$$\mathbf{x}_{t,f}^{ij} = \underbrace{[0, \dots, 0]}_{(i-1)Q}, \underbrace{\mathbf{x}_{t,f}^{j\top}}_{(j-i-1)Q}, \underbrace{[0, \dots, 0]}_{(I-j)Q}. \quad (3)$$

Then, for each pair (i, j) , we have:

$$\mathbf{x}_{t,f}^{ij\top} \mathbf{a}_f = 0. \quad (4)$$

Let's assume, for simplicity, that the reference channel is $r = 1$. To avoid a trivial solution to (4), i.e. $\mathbf{a}_f = \mathbf{0}$, we constrain the first CTF coefficient of the reference channel to be equal to 1. This is done by dividing both sides of (4) by $a_{0,f}^1$ and by moving the first entry of $\mathbf{x}_{t,f}^{ij}$, denoted by $-y_{t,f}^{ij}$, to the right side of (4), which rewrites as:

$$\tilde{\mathbf{x}}_{t,f}^{ij\top} \tilde{\mathbf{a}}_f = y_{t,f}^{ij}, \quad (5)$$

where $\tilde{\mathbf{x}}_{t,f}^{ij}$ is $\mathbf{x}_{t,f}^{ij}$ with the first entry removed, and $\tilde{\mathbf{a}}_f$ is the relative CTF vector:

$$\tilde{\mathbf{a}}_f = \left[\frac{\tilde{\mathbf{a}}_f^{1\top}}{a_{0,f}^1}, \frac{\mathbf{a}_f^{2\top}}{a_{0,f}^1}, \dots, \frac{\mathbf{a}_f^{I\top}}{a_{0,f}^1} \right]^\top, \quad (6)$$

where $\tilde{\mathbf{a}}_f^1 = [a_{1,f}^1, \dots, a_{Q-1,f}^1]^\top$ denotes \mathbf{a}_f^1 with the first entry removed. For $i = 2, \dots, I$, the DP-RTFs appear in (6) as the first entries of $\frac{\mathbf{a}_f^{i\top}}{a_{0,f}^1}$. Therefore, the DP-RTF estimation amounts to solving (5).

Equation (5) is defined for one microphone pair and for one frame. In batch mode, the terms $\tilde{\mathbf{x}}_{t,f}^{ij\top}$ and $y_{t,f}^{ij}$ of this equation can be concatenated across microphone pairs and frames to

construct a least square formulation. For online estimation, we would like to update the $\tilde{\mathbf{a}}_f$ using the current frame t . For notational convenience, let $m = 1, \dots, M$ denote the index of a microphone pair, where $M = I(I - 1)/2$. Then let the superscript ij be replaced with m . The fitting error of (5) is

$$e_{t,f}^m = y_{t,f}^m - \tilde{\mathbf{x}}_{t,f}^m \top \tilde{\mathbf{a}}_f. \quad (7)$$

At the current frame t , for the microphone pair m , RLS aims to minimize the error

$$J_{t,f}^m = \sum_{t'=1}^{t-1} \sum_{m'=1}^M \lambda^{t-t'} |e_{t',f}^{m'}|^2 + \sum_{m'=1}^m |e_{t,f}^{m'}|^2, \quad (8)$$

which sums up the fitting error of all the microphone pairs for the past frames and the microphone pairs up to m for the current frame. The forgetting factor $\lambda \in (0, 1]$ gives exponentially lower weight to older frames, whereas at one given frame, all microphone pairs have the same weight. To minimize $J_{t,f}^m$, we set its complex derivative with respect to $\tilde{\mathbf{a}}_f^*$ to zero, where $*$ denotes complex conjugate. We obtain an estimate of $\tilde{\mathbf{a}}_f$ at frame t for microphone pair m as:

$$\tilde{\mathbf{a}}_{t,f}^m = \mathbf{R}_{t,f}^{m-1} r_{t,f}^m, \quad (9)$$

with

$$\begin{aligned} \mathbf{R}_{t,f}^m &= \sum_{t'=1}^{t-1} \sum_{m'=1}^M \lambda^{t-t'} \tilde{\mathbf{x}}_{t',f}^{m'} * \tilde{\mathbf{x}}_{t',f}^{m'} \top + \sum_{m'=1}^m \tilde{\mathbf{x}}_{t,f}^{m'} * \tilde{\mathbf{x}}_{t,f}^{m'} \top, \\ r_{t,f}^m &= \sum_{t'=1}^{t-1} \sum_{m'=1}^M \lambda^{t-t'} \tilde{\mathbf{x}}_{t',f}^{m'} * y_{t',f}^{m'} + \sum_{m'=1}^m \tilde{\mathbf{x}}_{t,f}^{m'} * y_{t,f}^{m'}. \end{aligned}$$

It can be seen that the covariance matrix $\mathbf{R}_{t,f}^m$ is computed based on the rank-one modification, thence its inverse, denoted by $\mathbf{P}_{t,f}^m$, can be computed using the Sherman-Morrison formula, without the need of matrix inverse. The recursion procedure is summarized in Algorithm 1, where \mathbf{g} is the *gain vector*. The current frame t is initialized with the previous frame $t - 1$. At the first frame, we initialize $\tilde{\mathbf{a}}_{1,f}^0$ as zero, and $\mathbf{P}_{1,f}^0$ as the identity. At each frame, all microphone pairs are related to the same CTF vector that corresponds to the current speaker direction, hence all microphone pairs should be simultaneously used to estimate the CTF vector of the current frame. In batch mode, this can be easily implemented by concatenating the microphone pairs. However, in RLS, to satisfy the rank-one modification of the covariance matrix, we need to process the microphone pairs one by one as shown in (8) and Algorithm 1. At the end of the iterations over all microphone pairs, $\tilde{\mathbf{a}}_{t,f}^M$ is the ‘‘final’’ CTF estimation for the current frame, and is used for speaker localization. The DP-RTF estimates, denoted as $\tilde{c}_{t,f}^i$, $i = 2, \dots, I$, are obtained from $\tilde{\mathbf{a}}_{t,f}^M$. Note that implicitly we have $\tilde{c}_{t,f}^1 = 1$.

B. Multiple Moving Speakers

So far, the proposed online DP-RTF estimation method has been presented in the noise-free single-speaker case. The noisy multiple-speaker case was considered in [12], but only for static speakers, i.e. batch mode, and in the two-channel

Algorithm 1 RLS at frame t

Input: $\tilde{\mathbf{x}}_{t,f}^m, y_{t,f}^m, m = 1, \dots, M$
Initialization: $\tilde{\mathbf{a}}_{t,f}^0 \leftarrow \tilde{\mathbf{a}}_{t-1,f}^M, \mathbf{P}_{t,f}^0 \leftarrow \lambda^{-1} \mathbf{P}_{t-1,f}^M$
for $m = 1$ to M **do**
 $e_{t,f}^m = y_{t,f}^m - \tilde{\mathbf{x}}_{t,f}^m \top \tilde{\mathbf{a}}_{t,f}^{m-1}$
 $\mathbf{g} = \mathbf{P}_{t,f}^{m-1} \tilde{\mathbf{x}}_{t,f}^{m*} / (1 + \tilde{\mathbf{x}}_{t,f}^m \top \mathbf{P}_{t,f}^{m-1} \tilde{\mathbf{x}}_{t,f}^{m*})$
 $\mathbf{P}_{t,f}^m = \mathbf{P}_{t,f}^{m-1} - \mathbf{g} \tilde{\mathbf{x}}_{t,f}^m \top \mathbf{P}_{t,f}^{m-1}$
 $\tilde{\mathbf{a}}_{t,f}^m = \tilde{\mathbf{a}}_{t,f}^{m-1} + e_{t,f}^m \mathbf{g}$
end for
Output: $\tilde{\mathbf{a}}_{t,f}^M, \mathbf{P}_{t,f}^M$

case. We summarize the principles of this method and then explain in details the present online/multi-channel extension. We assume that the speakers are static over a short time and that a single source is active within a small region of the TF domain, due to the sparsity of speech in this domain. Therefore, the CTF is assumed to be locally time-invariant and can be estimated using the current frame as well as a small number of recent frames. This can be done by adjusting a forgetting factor λ . To approximately have a memory of P frames, we set $\lambda = \frac{P-1}{P+1}$. To efficiently estimate the CTF vector $\tilde{\mathbf{a}}_{t,f}^M$ of length $IQ - 1$, we need $\rho(IQ - 1)$ equations. The factor ρ should be empirically set to achieve a good tradeoff between the validity of the above assumptions and a robust estimate of $\tilde{\mathbf{a}}_{t,f}^M$. To have $\rho(IQ - 1)$ equations, we need $P = \frac{\rho(IQ-1)}{I(I-1)/2} \approx \rho \frac{2Q}{I-1}$ frames. One may observe that the number of frames that are needed to estimate $\tilde{\mathbf{a}}_{t,f}^M$ decreases as the number of microphones is increased.

When noise is present, especially if the noise sources are directional, the CTF estimate can be contaminated. In addition, even in a low-noise case, many TF bins are dominated by noise due to the sparsity of speech spectra. To remove the noise, we use the inter-frame spectral subtraction algorithm proposed in [30], [33]. Briefly stated, at each frequency bin f , the frames are first classified into speech frames and noise frames, then inter-frame spectral subtraction is applied between them. In the RLS process, only the speech frames (after spectral subtraction) are used, and the noise frames are skipped. A speech frame with a preceding noise frame is initialized with the latest speech frame.

In practice, a DP-RTF estimate can sometimes be unreliable. Possible reasons are that in a small frame region, (i) the CTF is time-varying due to a fast movement of the speakers, (ii) multiple speakers are present, (iii) only noise is present due to a wrong noise-speech classification, or (iv) only reverberation is present at the end of a speech occurrence. In [12], a consistency test was proposed to tackle this problem: If a small frame region corresponds to an actual active speaker, the DP-RTFs estimated using different reference channels are consistent, otherwise the DP-RTFs are biased, with inconsistent bias values. In the present work, we use the first and second channels as references, we obtain the DP-RTF estimates $\tilde{c}_{t,f}^i$ (with $\tilde{c}_{t,f}^1 = 1$) and $\bar{c}_{t,f}^i$ (with $\bar{c}_{t,f}^2 = 1$), respectively. Then $\tilde{c}_{t,f}^i$ and $\bar{c}_{t,f}^i / \bar{c}_{t,f}^1$ are two estimates of the same DP-RTF $a_{0,f}^i / a_{0,f}^1$. If the similarity of these two estimates is large, they are said

to be consistent. They are then averaged and normalized as done in [12], resulting in a final complex-valued feature $\hat{c}_{t,f}^i$ with module in $[0, 1]$. The estimates that do not pass this consistency test are simply ignored.

Finally, at frame t , we obtain a set of features $\mathcal{C}_t = \{\{\hat{c}_{t,f}^i\}_{i \in \mathcal{I}_f}\}_{f=0}^{F-1}$, where $\mathcal{I}_f \subseteq \{2, \dots, I\}$ denotes the set of microphone indices that pass the consistency test. Note that \mathcal{I}_f is empty if frame t is a noise frame at frequency f , or if no channel passes the consistency test. Each of the features is assumed to be associated with a single speaker.

III. LOCALIZATION OF MULTIPLE MOVING SPEAKERS

In this section we describe the proposed frame-wise online multiple-speaker localizer. We start by briefly presenting the underlying complex Gaussian mixture model, followed by the recursive estimation of its parameters.

A. Generative Model for Multiple-Speaker Localization

In order to associate DP-RTF features from \mathcal{C}_t with speakers and to localize each active speaker, we adopt the generative model proposed in [7]. Let $\mathcal{D} = \{\tilde{\theta}_1, \dots, \tilde{\theta}_d, \dots, \tilde{\theta}_D\}$ be a set of D candidate source *directions*, e.g. azimuth angles. An observed feature $\hat{c}_{t,f}^i$ (cf. Section II), when emitted by a sound source located along the direction $\tilde{\theta}_d$, is assumed to be drawn from a complex-Gaussian distribution with mean $c_f^{i,d}$ and variance σ^2 , i.e. $\hat{c}_{t,f}^i | d \sim \mathcal{N}_c(\hat{c}_{t,f}^i, \sigma^2)$. The mean $c_f^{i,d}$ is the predicted feature at frequency f for channel i , and is precomputed based on direct-path propagation along azimuth $\tilde{\theta}_d$ to the microphones. The variance σ^2 is empirically set as a constant value. The marginal density of an observed feature $\hat{c}_{t,f}^i$ (taking into account all candidate directions) is a CGMM with each component corresponding to a candidate direction:

$$p(\hat{c}_{t,f}^i | \mathcal{D}) = \sum_{d=1}^D w_d \mathcal{N}_c(\hat{c}_{t,f}^i; c_f^{i,d}, \sigma^2), \quad (10)$$

where $w_d \geq 0$ is the prior probability (component weight) of the d -th component, with $\sum_{d=1}^D w_d = 1$. Let us denote the vector of weights with $\mathbf{w} = (w_1, \dots, w_D)^\top$. Note that this vector is the only free parameter of the model.

Assuming that the observations in \mathcal{C}_t are independent, the corresponding (normalized) negative log-likelihood function, as a function of w_d , is given by:

$$\mathcal{L}_t = -\frac{1}{|\mathcal{C}_t|} \sum_{\hat{c}_{t,f}^i \in \mathcal{C}_t} \log \left(\sum_{d=1}^D w_d \mathcal{N}_c(\hat{c}_{t,f}^i; c_f^{i,d}, \sigma^2) \right), \quad (11)$$

where $|\mathcal{C}_t|$ denotes the cardinality of \mathcal{C}_t . Once \mathcal{L}_t is minimized, each weight w_d represents the probability that a speaker is active in the direction $\tilde{\theta}_d$. Therefore, sound source localization amounts to the minimization of \mathcal{L}_t . In addition, taking into account the fact that the number of actual active speakers is much lower than the number of candidate directions, an

entropy term was proposed in [12] as a regularizer to impose a sparse solution for w_d . The entropy is defined as

$$H = -\sum_{d=1}^D w_d \log(w_d). \quad (12)$$

The concave-convex procedure [34] was adopted in [12], to minimize the objective function $\mathcal{L} + \gamma H$ w.r.t. \mathbf{w} , where \mathcal{L} is the normalized negative log-likelihood of the DP-RTF features of all frames, i.e. batch mode optimization, and the positive scalar γ was used to control the tradeoff between likelihood minimization and imposing sparsity over the weights. In the batch mode, the weight vector \mathbf{w} is shared across all frames. Hence this method is not suitable for moving speakers.

B. Recursive Parameter Estimation

We now describe a recursive method for updating the weight vector from \mathbf{w}_{t-1} to \mathbf{w}_t , i.e. from frame $t-1$ to frame t , using the DP-RTF features at t . This can be formulated as the following online optimization problem [31]:

$$\begin{aligned} \mathbf{w}_t = \underset{\mathbf{w}}{\operatorname{argmin}} \quad & \chi(\mathbf{w}, \mathbf{w}_{t-1}) + \eta(\mathcal{L}_t + \gamma H), \\ \text{s.t.} \quad & w_d > 0, \quad \forall d \in \{1 \dots D\} \quad \text{and} \quad \sum_{d=1}^D w_d = 1, \end{aligned} \quad (13)$$

where $\chi(\mathbf{a}, \mathbf{b})$ is a distance between \mathbf{a} and \mathbf{b} . The positive scalar factor η controls the parameter update rate. To minimize (13), the derivative of the objective function w.r.t \mathbf{w} is set to zero, yielding a set of equations with no closed-form solution. To speed up the computation, it is assumed that \mathbf{w}_t is close to \mathbf{w}_{t-1} , hence the derivative of $\mathcal{L}_t + \gamma H$ at \mathbf{w} can be approximated with the derivative of $\mathcal{L}_t + \gamma H$ at \mathbf{w}_{t-1} . This assumption is reasonable when parameter evolution is not too fast. As a result, when the distance $\chi(\mathbf{w}, \mathbf{w}_{t-1})$ is Euclidean, the objective function leads to gradient descent with a step length equal to η . Nevertheless, the constraints (14) lead to an inefficient gradient descent procedure. To obtain an efficient solver, we exploit the fact that the weights w_d are probability masses, hence we replace the Euclidean distance with the more suitable Kullback-Leibler divergence, i.e. $\chi(\mathbf{w}, \mathbf{w}_{t-1}) = \sum_{d=1}^D w_d \log \frac{w_d}{w_{t-1,d}}$, which results in the exponentiated gradient algorithm [31].

The partial derivatives of \mathcal{L}_t and H w.r.t w_d at the point $w_{t-1,d}$ are computed with, respectively:

$$\begin{aligned} \frac{\partial(\mathcal{L}_t)}{\partial w_d} \Big|_{w_{t-1,d}} &= -\frac{1}{|\mathcal{C}_t|} \sum_{\hat{c}_{t,f}^i \in \mathcal{C}_t} \frac{\mathcal{N}_c(\hat{c}_{t,f}^i; c_f^{i,d}, \sigma^2)}{\sum_{d'=1}^D w_{t-1,d'} \mathcal{N}_c(\hat{c}_{t,f}^i; c_f^{i,d'}, \sigma^2)}, \\ \frac{\partial H}{\partial w_d} \Big|_{w_{t-1,d}} &= -(1 + \log(w_{t-1,d})), \quad \forall d \in \mathcal{D}. \end{aligned} \quad (15)$$

Then, the exponentiated gradient,

$$r_{t-1,d} = e^{-\eta \left(\frac{\partial(-\mathcal{L}_t)}{\partial w_d} \Big|_{w_{t-1,d}} + \gamma \frac{\partial H}{\partial w_d} \Big|_{w_{t-1,d}} \right)}, \quad \forall d \in \mathcal{D}, \quad (16)$$

is used to update the weights with:

$$w_{td} = \frac{r_{t-1,d} w_{t-1,d}}{\sum_{d'=1}^D r_{t-1,d'} w_{t-1,d'}}, \quad \forall d \in \mathcal{D}. \quad (17)$$

It is clear from (17) that the parameter constraints (14) are necessarily satisfied. The exponentiated gradient algorithm sequentially evaluates (15), (16) and (17) at each frame. At the first frame, the weights are initialized with the uniform distribution, namely $w_{1,d} = \frac{1}{D}$. When \mathcal{C}_t is empty, such as during a silent period, the parameters are recursively updated with $w_{td} = (1 - \eta')w_{t-1,d} + \eta' \frac{1}{D}$.

A plot of w_d as a function of d exhibits a curve with a few peaks that should correspond to active speakers. The use of the entropy regularization term has been shown to both suppress small spurious peaks, observed without regularization, and sharpen the peaks corresponding to actual active speakers, thus allowing to better discriminate between the true peaks and spurious ones. In the case of moving speakers, a peak should jump from a direction $\tilde{\theta}_d$ to one of its neighbors across frames. Spatial smoothing over the weight curve raises the weight values around a peak, which results in smoother peak jumps. For example, a one-dimensional smoothing filter $[0.02 \ 1 \ 0.02]/1.04$ is used in the experiments. Note that the effect of spatial smoothing is somehow contrary to the effect of entropy regularization, which imposes sparsity, but overall, the combination of these two processes has been found to be beneficial in practice.

C. Peak selection and frame-wise speaker localization

Frame-wise counting and localization of active speakers can be jointly carried out by selecting the peaks in the profile of w_t with a value larger than a predefined threshold [12], [30]. However, peak selection does not exploit the long-term movement information of speakers. Moreover, peak selection can be a risky process since a too high or too low setting of the threshold can lead to unsatisfying missed detection rate or false alarm rate, respectively. In the present work, rather than applying peak selection, the whole set of candidate directions together with their associated weights are fed to the Bayesian tracker described in the next section, which is in charge of transforming the frame-wise localization information into consistent sequences of speaker position, i.e. speaker tracks. The tracker refines the speaker localization trajectory and mitigates the risk of threshold setting. Peak selection is used in the experiments only to evaluate the proposed multiple speaker localizer, i.e. without tracking.

IV. MULTIPLE-SPEAKER TRACKING

Let N be the maximum number of speakers that can be simultaneously active at any time t , and let n be the speaker index. Moreover, let $n = 0$ denote *no speaker*. We now introduce the main variables and their notations. Upper case letters denote random variables while lower case letters denote their realizations.

Let \mathbf{S}_{tn} be a latent (or state) variable associated with speaker n at frame t , and let $\mathbf{S}_t = (\mathbf{S}_{t1}, \dots, \mathbf{S}_{tn}, \dots, \mathbf{S}_{tN})$. \mathbf{S}_{tn} is composed of two parts: the speaker direction and

the speaker velocity. In this work, speaker direction is defined by an azimuth θ_{tn} . To avoid phase (circular) ambiguity we describe the direction with the unit vector $\mathbf{U}_{tn} = (\cos(\theta_{tn}), \sin(\theta_{tn}))^\top$. Moreover, let $V_{tn} \in \mathbb{R}$ be the angular velocity. Altogether we define a realization of the state variable as $\mathbf{s}_{tn} = [\mathbf{u}_{tn}; v_{tn}]$ where the notation $[\cdot; \cdot]$ stands for vertical vector concatenation.

Let $\mathbf{O}_t = (\mathbf{O}_{t1}, \dots, \mathbf{O}_{td}, \dots, \mathbf{O}_{tD})$ be the observed variables at frame t . Each realization \mathbf{o}_{td} of \mathbf{O}_{td} is composed of a candidate location, or azimuth $\tilde{\theta}_{td} \in \mathcal{D}$, and a weight w_{td} . The weight w_{td} is the probability that there is an active speaker in the direction $\tilde{\theta}_{td}$, namely (13). As above, let the azimuth be described by a unit vector $\mathbf{b}_{td} = (\cos(\tilde{\theta}_{td}), \sin(\tilde{\theta}_{td}))^\top$. In summary we have $\mathbf{o}_{td} = [\mathbf{b}_{td}; w_{td}]$. Moreover, let Z_{td} be a (latent) assignment variable associated with each observed variable \mathbf{O}_{td} , such that $Z_{td} = n$ means that the observation indexed by d at frame t is assigned to active speaker $n \in \{0, \dots, N\}$. Note that $Z_{td} = 0$ is a “fake” assignment – the corresponding observation is assigned to an audio source that is either background noise or any other source that has not yet been identified as an active speaker.

The problem at hand can now be cast into the estimation of the filtering distribution $p(\mathbf{s}_t, \mathbf{z}_t | \mathbf{o}_{1:t})$, and further inference of \mathbf{s}_t and \mathbf{z}_t . In this work we make two hypotheses, namely (i) that the observations at frame t only depend on the assignment and state variables at t , and (ii) that the prior distribution of the assignment variables is independent of all the other variables. By applying the Bayes rule together with these hypotheses, and ignoring terms that do not depend on \mathbf{s}_t and \mathbf{z}_t , the filtering distribution is proportional to:

$$p(\mathbf{s}_t, \mathbf{z}_t | \mathbf{o}_{1:t}) \propto p(\mathbf{o}_t | \mathbf{s}_t, \mathbf{z}_t) p(\mathbf{z}_t) p(\mathbf{s}_t | \mathbf{o}_{1:t-1}), \quad (18)$$

which contains three terms: the observation model, the prior distribution of the assignment variable and the predictive distribution over the sources state. We now characterize each one of these three terms.

1) *Audio observation model*: The audio observation model describes the distribution of the observations given speakers state and assignment. We assume the different observations are independent conditionally to speakers state and assignment, which can be written as:

$$p(\mathbf{o}_t | \mathbf{z}_t, \mathbf{s}_t) = \prod_{d=1}^D p(\mathbf{o}_{td} | \mathbf{z}_t, \mathbf{s}_t). \quad (19)$$

Since the weights describe the confidence associated with each observed azimuth, we adopt the weighted-data GMM model of [35]:

$$p(\mathbf{b}_{td} | Z_{td} = n, \mathbf{s}_{tn}; w_{td}) = \begin{cases} \mathcal{N}(\mathbf{b}_{td}; \mathbf{M}\mathbf{s}_{tn}, \frac{1}{w_{td}}\Sigma) & \text{if } n \in \{1, \dots, N\} \\ \mathcal{U}(\text{vol}(\mathcal{G})) & \text{if } n = 0 \end{cases}, \quad (20)$$

where the matrix $\mathbf{M} = [\mathbf{I}_{2 \times 2}, \mathbf{0}_{2 \times 1}]$ projects the state variable onto the space of source directions and Σ is a covariance matrix (set empirically to a fixed value in the present study). Note that the weight plays the role of a precision: The higher

the weight w_{td} , the more reliable the source direction \mathbf{b}_{td} . The case $Z_{td} = 0$ follows a uniform distribution over the volume of the observation space.

2) *Prior distribution*: The prior distribution of the assignment variable is independent over observations and is assumed to be uniformly distributed over all the speakers (including the case $n = 0$), hence:

$$p(\mathbf{z}_t) = \prod_{d=1}^D p(Z_{td} = n) \quad \text{with} \quad \pi_{dn} = p(Z_{td} = n) = \frac{1}{N+1}. \quad (21)$$

3) *Predictive distribution*: The predictive distribution describes the relationship between the state \mathbf{s}_t and the past observations up to frame t , $\mathbf{o}_{1:t-1}$. To calculate this distribution, we first marginalize $p(\mathbf{s}_t|\mathbf{o}_{1:t-1})$ over \mathbf{s}_{t-1} , writing:

$$p(\mathbf{s}_t|\mathbf{o}_{1:t-1}) = \int p(\mathbf{s}_t|\mathbf{s}_{t-1})p(\mathbf{s}_{t-1}|\mathbf{o}_{1:t-1})d\mathbf{s}_{t-1}, \quad (22)$$

where the two terms under the integral are the state dynamics and the marginal filtering distribution of the state variable at frame $t-1$, respectively. We model the state dynamics as a linear-Gaussian first-order Markov process, independent over the speakers, i.e. :

$$p(\mathbf{s}_t|\mathbf{s}_{t-1}) = \prod_{n=1}^N \mathcal{N}(\mathbf{s}_{tn}; \mathbf{D}_{t-1,n}\mathbf{s}_{t-1,n}, \mathbf{\Lambda}_{tn}), \quad (23)$$

where $\mathbf{\Lambda}_{tn}$ is the dynamics' covariance matrix and $\mathbf{D}_{t-1,n}$ is the state transition matrix. Given the estimated azimuth $\theta_{t-1,n}$ and angular velocity $v_{t-1,n}$ at frame $t-1$, we have the following relation:

$$\begin{pmatrix} \cos(\theta_{tn}) \\ \sin(\theta_{tn}) \end{pmatrix} = \begin{pmatrix} \cos(\theta_{t-1,n} + v_{t-1,n}\Delta t) \\ \sin(\theta_{t-1,n} + v_{t-1,n}\Delta t) \end{pmatrix}, \quad (24)$$

where Δt is the time increment between two consecutive frames. Expanding (24) and assuming that the angular displacement $v_{t-1,n}\Delta t$ is small, the state transition matrix can be written as:

$$\mathbf{D}_{t-1,n} = \begin{pmatrix} 1 & 0 & -\sin(\theta_{t-1,n}) \\ 0 & 1 & \cos(\theta_{t-1,n}) \\ 0 & 0 & \Delta t \end{pmatrix}. \quad (25)$$

In the following $\mathbf{D}_{t-1,n}$ is written as \mathbf{D} , only to lighten the equations.

A. Variational Expectation Maximization Algorithm

At this point, the standard solution to the calculation of the filtering distribution consists of using EM methodology. EM alternates between evaluating the expected complete-data log-likelihood and maximizing this expectation with respect to the model parameters. More precisely, the expectation writes:

$$J(\Theta, \Theta^o) = \mathbf{E}_{p(\mathbf{z}_t, \mathbf{s}_t|\mathbf{o}_{1:t}, \Theta^o)} [\log p(\mathbf{z}_t, \mathbf{s}_t, \mathbf{o}_{1:t}|\Theta)], \quad (26)$$

where Θ^o denotes the current parameter estimates and Θ denotes the new estimates, obtained via maximization of (26). However, given the hybrid combinatorial-and-continuous nature of the latent space, such solution is intractable in

practice, due to combinatorial explosion. We thus propose to use of a variational approximation to solve the problem efficiently. We inspire from [22] and propose the following factorization:

$$p(\mathbf{z}_t, \mathbf{s}_t|\mathbf{o}_{1:t}) \approx q(\mathbf{z}_t, \mathbf{s}_t) = q(\mathbf{z}_t) \prod_{n=0}^N q(\mathbf{s}_{tn}). \quad (27)$$

The optimal solution is then given by two E-steps, an E-S step for each individual state variable \mathbf{s}_{tn} and an E-Z step for the assignment variable \mathbf{z}_t :

$$\log q(\mathbf{s}_{tn}) = \mathbf{E}_{q(\mathbf{z}_t)} \prod_{m \neq n} q(\mathbf{s}_{tm}) [\log p(\mathbf{z}_t, \mathbf{s}_t|\mathbf{o}_{1:t})], \quad (28)$$

$$\log q(\mathbf{z}_t) = \mathbf{E}_{q(\mathbf{s}_t)} [\log p(\mathbf{z}_t, \mathbf{s}_t|\mathbf{o}_{1:t})]. \quad (29)$$

It is easy to see that in order to compute (28) and (29), two elements are needed: the predictive distribution (22) and the marginal filtering distribution at $t-1$, $p(\mathbf{s}_{t-1}|\mathbf{o}_{1:t-1})$. Remarkably, as a consequence of the factorization (27), we can replace $p(\mathbf{s}_{t-1}|\mathbf{o}_{1:t-1})$ with $q(\mathbf{s}_{t-1}) = \prod_{n=1}^N q(\mathbf{s}_{t-1,n})$ in (22) and compute the predictive distribution as follows:

$$p(\mathbf{s}_t|\mathbf{o}_{1:t-1}) \approx \int p(\mathbf{s}_t|\mathbf{s}_{t-1}) \prod_{n=1}^N q(\mathbf{s}_{t-1,n}) d\mathbf{s}_{t-1}. \quad (30)$$

This predictive distribution factorizes across speakers. Moreover, one prominent feature of the proposed variational approximation is that, if the posterior distribution at time $t-1$ $q(\mathbf{s}_{t-1,n})$ is assumed to be a Gaussian, say

$$q(\mathbf{s}_{t-1,n}) = \mathcal{N}(\mathbf{s}_{t-1,n}; \boldsymbol{\mu}_{t-1,n}, \boldsymbol{\Gamma}_{t-1,n}), \quad (31)$$

then (the approximation of) the predictive distribution (30) is a Gaussian. More specifically, the derivation of (30) leads to:

$$p(\mathbf{s}_{tn}|\mathbf{o}_{1:t-1}) = \mathcal{N}(\mathbf{s}_{tn}; \mathbf{D}\boldsymbol{\mu}_{t-1,n}, \mathbf{D}\boldsymbol{\Gamma}_{t-1,n}\mathbf{D}^\top + \mathbf{\Lambda}_{tn}). \quad (32)$$

Moreover, as we will see in the E-S-step below, the posterior distribution at time t , $q(\mathbf{s}_{tn})$, is also a Gaussian.

1) *E-S step*: The computation of the variational posterior distribution $q(\mathbf{s}_{tn})$, for all currently tracked speakers, is carried out by developing (28) as follows. We first exploit (18), (19), (21) and (32) to rewrite $\log p(\mathbf{z}_t, \mathbf{s}_t|\mathbf{o}_{1:t})$ in (28) as a sum of individual log-probabilities. Then we eliminate all terms not depending on \mathbf{s}_{tn} and we evaluate the expectation of the remaining terms. Because the terms not depending on \mathbf{s}_{tn} were disregarded, the expectation is computed only with respect to $q(\mathbf{z}_t)$. This nicely makes the computation of $q(\mathbf{s}_{tn})$ independent of the structure of $q(\mathbf{s}_{tm})$ for $m \neq n$. Eventually, this yields a Gaussian distribution:

$$q(\mathbf{s}_{tn}) = \mathcal{N}(\mathbf{s}_{tn}; \boldsymbol{\mu}_{tn}, \boldsymbol{\Gamma}_{tn}), \quad (33)$$

with the following parameters:

$$\begin{aligned} \boldsymbol{\Gamma}_{tn} = & \left(\left(\sum_{d=1}^D \alpha_{tdn} w_{td} \right) \mathbf{M}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M} \right. \\ & \left. + \left(\mathbf{\Lambda}_{tn} + \mathbf{D}\boldsymbol{\Gamma}_{t-1,n}\mathbf{D}^\top \right)^{-1} \right)^{-1}, \end{aligned} \quad (34)$$

$$\begin{aligned} \boldsymbol{\mu}_{tn} = & \boldsymbol{\Gamma}_{tn} \left(\mathbf{M}^\top \boldsymbol{\Sigma}^{-1} \left(\sum_{d=1}^D \alpha_{tdn} w_{td} \mathbf{b}_{td} \right) \right. \\ & \left. + \left(\mathbf{\Lambda}_{tn} + \mathbf{D}\boldsymbol{\Gamma}_{t-1,n}\mathbf{D}^\top \right)^{-1} \mathbf{D}\boldsymbol{\mu}_{t-1,n} \right), \end{aligned} \quad (35)$$

where $\alpha_{tdn} = q(Z_{td} = n)$ is the variational posterior distribution of the assignment variable, which will be detailed in Section IV-A2. Importantly, the first two entries of $\boldsymbol{\mu}_{tn}$ in (35) represent the estimated azimuth of speaker n . The ‘‘final’’ azimuth estimate at frame t is thus given by this subvector at the end of the VEM iterations. Since we use a unit-vector representation, we normalize this vector at each iteration of the algorithm. Finally, note that since we have shown that $q(\mathbf{s}_{t-1,n})$ being Gaussian leads to $q(\mathbf{s}_{tn})$ being Gaussian as well, it is sufficient to assume that $q(\mathbf{s}_{1n})$ is Gaussian, namely at $t = 1$: $q(\mathbf{s}_{1n}) = \mathcal{N}(\mathbf{s}_{1n}; \boldsymbol{\mu}_{1n}, \boldsymbol{\Gamma}_{1n})$.

2) *E-Z step*: Developing (29) with the same principles as above, one can easily find that the variational posterior distribution of the assignment variable factorizes as:

$$q(\mathbf{z}_t) = \prod_{d=1}^D q(z_{td}). \quad (36)$$

In addition, we obtain a closed-form expression for $q(z_{td})$:

$$\alpha_{tdn} = q(Z_{td} = n) = \frac{\rho_{tdn} \pi_{dn}}{\sum_{i=0}^N \rho_{tdi} \pi_{di}}, \quad (37)$$

where π_{dn} was defined in (21), and ρ_{tdn} is given by:

$$\rho_{tdn} = \begin{cases} \mathcal{N}(\mathbf{b}_{td}; \mathbf{M}\boldsymbol{\mu}_{tn}, \frac{1}{w_{td}}\boldsymbol{\Sigma}) \\ \times e^{-\frac{1}{2}\text{tr}(w_{td}\mathbf{M}^\top \boldsymbol{\Sigma}^{-1} \mathbf{M}\boldsymbol{\Gamma}_{tn})} & \text{if } 1 \leq n \leq N \\ \mathcal{U}(\text{vol}(\mathcal{G})) & \text{if } n = 0. \end{cases} \quad (38)$$

3) *M-step*: Once the two expectation steps are executed, we maximize J in (26) with respect to the model parameters, i.e. the covariance matrix of the state dynamics $\boldsymbol{\Lambda}_{tn}$. By exploiting again the proposed variational approximation, the dependency of J on $\boldsymbol{\Lambda}_{tn}$ can be written as:

$$J(\boldsymbol{\Lambda}_{tn}) = \mathbf{E}_{q(\mathbf{s}_{tn})} [\log \mathcal{N}(\mathbf{s}_{tn}; \mathbf{D}\boldsymbol{\mu}_{t-1,n}, \mathbf{D}\boldsymbol{\Gamma}_{t-1,n}\mathbf{D}^\top + \boldsymbol{\Lambda}_{tn})],$$

which can be further developed as:

$$J(\boldsymbol{\Lambda}_{tn}) = \log |\boldsymbol{\Lambda}_{tn}| + \text{Tr} [\boldsymbol{\Lambda}_{tn}^{-1} (\boldsymbol{\Gamma}_{tn} - \mathbf{D}\boldsymbol{\Gamma}_{t-1,n}\mathbf{D}^\top + (\boldsymbol{\mu}_{tn} - \mathbf{D}\boldsymbol{\mu}_{t-1,n})(\boldsymbol{\mu}_{tn} - \mathbf{D}\boldsymbol{\mu}_{t-1,n})^\top)]. \quad (39)$$

By equating to zero the gradient of (39) w.r.t. $\boldsymbol{\Lambda}_{tn}$, we obtain:

$$\boldsymbol{\Lambda}_{tn} = \boldsymbol{\Gamma}_{tn} - \mathbf{D}\boldsymbol{\Gamma}_{t-1,n}\mathbf{D}^\top + (\boldsymbol{\mu}_{tn} - \mathbf{D}\boldsymbol{\mu}_{t-1,n})(\boldsymbol{\mu}_{tn} - \mathbf{D}\boldsymbol{\mu}_{t-1,n})^\top. \quad (40)$$

B. Speaker-Birth Process

A birth process is used to initialize new tracks, i.e. speakers that become active. We take inspiration from the birth process for visual tracking proposed in [22] and adapt it to audio tracking. The general principle is the following. In a short period of time, say from $t - L$ to t , with L being small (typically 3), we assume that at most one new (yet untracked) speaker becomes active. For each frame from $t - L$ to t , among the observations assigned to $n = 0$ we select the one with the highest weight, and thus obtain an observation sequence $\tilde{\mathbf{o}}_{t-L:t}$. We then compute the marginal likelihood of this sequence according to our model, $\tau_0 = p(\tilde{\mathbf{o}}_{t-L:t})$. If these

observations have been generated by a speaker that has not been detected yet (hypotheses H_1), then they are assumed to be consistent with the model, i.e. exhibit smooth trajectories, and τ_0 will be high; otherwise, i.e. if they have been generated by background noise (hypotheses H_0), they will be more randomly spread over the range of possible observations, and τ_0 will be low. Giving birth to a new speaker track amounts to setting a threshold τ_1 and deciding between the two hypotheses:

$$\tau_0 \underset{H_0}{\overset{H_1}{>}} \tau_1. \quad (41)$$

This process is applied continuously over time to detect new speakers. This includes speaker track initialization at $t = 1$. Note that initially all the assignment variables are set to $n = 0$ (background noise), namely $Z_{1d} = 0, \forall d$.

As for the computation of $p(\tilde{\mathbf{o}}_{t-L:t})$, we first rewrite it as the marginalization of the joint probability of the selected observations and the state trajectory $\hat{\mathbf{s}}_{t-L:t}$ of a potential speaker:

$$\tau_0 = \int p(\tilde{\mathbf{o}}_{t-L:t}, \hat{\mathbf{s}}_{t-L:t}) d\hat{\mathbf{s}}_{t-L:t}, \quad (42)$$

which, under the proposed model, is given by:

$$\tau_0 = \int \left(\prod_{i=t-L+1}^t p(\tilde{\mathbf{o}}_i | \hat{\mathbf{s}}_i) p(\hat{\mathbf{s}}_i | \hat{\mathbf{s}}_{i-1}) \right) p(\tilde{\mathbf{o}}_{t-L} | \hat{\mathbf{s}}_{t-L}) p(\hat{\mathbf{s}}_{t-L}) d\hat{\mathbf{s}}_{t-L:t}. \quad (43)$$

All the terms in the above equation have been defined during the derivation of our model except the marginal prior distribution of the state $p(\hat{\mathbf{s}}_{t-L})$, and all these terms are Gaussian. For the track-birth process, we just want to test if the trajectory of observations from $t-L$ to t is coherent, and we can define here $p(\hat{\mathbf{s}}_{t-L})$ as a non-informative distribution, such as a uniform distribution. In practice we choose a Gaussian distribution with a very large covariance, to ensure a closed-form solution to (43). Due to room limitation, we do not present more details. Let us just mention that in practice we set $L = 3$, which enables efficient speaker birth detection.

C. Speaker Activity Detection

A very interesting feature of the proposed model is that, once speaker tracks have been estimated, the posterior distribution of the assignment variables \mathbf{Z}_t can be used for speech activity detection, i.e. who are the active speakers at each frame, a task also referred to as *speaker diarization* in the multi-speaker context. This can be formalized as testing for each frame t and each speaker n between the two following hypotheses: H_1 : Speaker n is active at frame t , and H_0 : Speaker n is silent at frame t . In the present work, this is done by computing the following *weighted sum of weights*, averaged over a small number of frames L' to take into account speaker activity inertia, and comparing with a threshold δ , a test formally written as:

$$\sum_{i=t-L'+1}^t \sum_{d=1}^D \alpha_{idn} w_i^d \underset{H_0}{\overset{H_1}{>}} \delta. \quad (44)$$

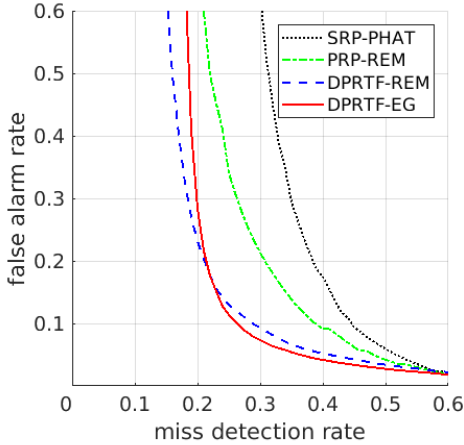


Fig. 2: ROC curve for the LOCATA dataset.

V. EXPERIMENTS

In this section, experiments with multiple moving speakers are conducted to evaluate the performance of the proposed online localization and tracking methods.

A. Experimental setup

1) *Datasets*: We used the LOCATA and the Kinovis datasets. The LOCATA (the IEEE-AASP Challenge on sound source localization and tracking) [36] data were recorded in the Computing Laboratory of the department of Computer Science of the Humboldt University Berlin. The room size is 7.1 m \times 9.8 m \times 3 m, with a reverberation time $T_{60} \approx 0.55$ s. We report the results of the development corpus for tasks #3 and #5 with a single moving speaker, and tasks #4 and #6 with two moving speakers, each task comprising three recorded sequences. In this work, we use four microphones with indices {5, 8, 11, 12} out of the twelve microphones of a spherical array built in the head of a humanoid robot, i.e. NAO, to perform azimuth localization and tracking. These four microphones are mounted on the top of the robot head, and they approximately lie in a horizontal plane parallel. An optical motion capture system is used to provide ground-truth positions of the robot and the moving speakers. The participants are supposed to continuously speak during the entire recordings. However, the phonetic pauses are inevitable, which sometimes last several seconds. We apply the voice activity detector proposed in [37] to the source signal (recorded by a reference microphone close to the speaker) of each participant to obtain ground-truth voice activity information.

The Kinovis data are recorded in the Kinovis motion capture laboratory at INRIA Grenoble [38]. The room size is 10.19 m \times 9.87 m \times 5.6 m, with $T_{60} \approx 0.53$ s. A version 5 NAO robot with four microphones [39] is used. The topology and layout of the microphones is similar to the ones used in LOCATA. The speakers were moving around the robot with a speaker-to-robot distance ranging between 1.5 m and 3.5 m. As with LOCATA, a motion capture system was used to obtain

ground-truth trajectories of the moving participants and the location of the robot. Ten sequences were recorded with up to three participants, for a total length of about 357 s. Here the participants behave more naturally than in the LOCATA data, i.e. they take speech turns in a natural multi-party dialog. When one participant is silent, he/she manually hides the infrared marker located on his head to make it invisible to the motion capture system. This provides ground-truth speech activity information for each participant. This dataset and the associated annotations allows us to test the proposed tracking algorithm when the number of active speakers varies over time.

2) *Parameter setting*: For both datasets, we perform 360°-wide azimuth estimation and tracking: $D = 72$ azimuth directions at every 5° in $[-175^\circ, 180^\circ]$ are used as candidate directions. The CGMM mean $c_f^{i,d}$ is the HRTF ratio between two microphones, which are precomputed based on the direct-path propagation model for each candidate direction. In the Kinovis dataset, the HRTFs have been measured to compute the CGMM means. For LOCATA, the TDOAs are computed based on the coordinate of microphones, which are then used to compute the phase of the CGMM means, while the magnitude of the CGMM means are set to a constant, e.g. 0.5, for all the frequencies. All the recorded signals are resampled to 16 kHz. The STFT uses the Hamming window with length of 16 ms and shift of 8 ms. The CTF length is $Q = 8$ frames. The RLS forgetting factor λ is computed using $\rho = 1$. The exponentiated gradient update factor is $\eta = 0.07$. The smoothing factor η' is set to 0.065. The entropy regularization factor is $\gamma = 0.1$. For the tracker, the covariance matrix is set to be isotropic $\Sigma = 0.03\mathbf{I}_2$. The threshold giving birth to a new identity is $\tau_1 = 0.75$ and $L = 3$. To decide whether a person is speaking or is silent, $L' = 3$ frames are used, with a threshold $\delta = 0.15$. At each time instance, the VEM algorithm has 5 iterations.

3) *Comparison with Baseline Methods*: The proposed method is evaluated both in “frame-wise localization” mode and in “tracker” mode. In the first mode, the frame-wise online localization module of Section III is applied without the tracker of Section IV. Instead, it is followed by the peak selection process described in [12]. This method is referred to as DPRTF-EG which stands for direct-path relative transfer function (DPRTF) using exponentiated gradient (EG). In tracker mode, DPRTF-EG is directly followed by the proposed VEM tracker, without peak selection. It is then simply referred to as VEM-tracker. In that case, the directions of active speakers are given by the state variable, and the continuity of the speaker tracks is given by the assignment variable. We compare DPRTF-EG with several baseline methods:

- The standard beamforming-based localization method called Steered Response Power using the PHase Transform (SRP-PHAT) [3]. The same STFT configuration and candidate directions are used for SRP-PHAT and for the proposed method. The steering vector for each candidate direction is derived from the HRTFs and TDOAs for the Kinovis and LOCATA datasets, respectively. The frame-wise SRP is recursively smoothed with a smoothing factor set to 0.065.

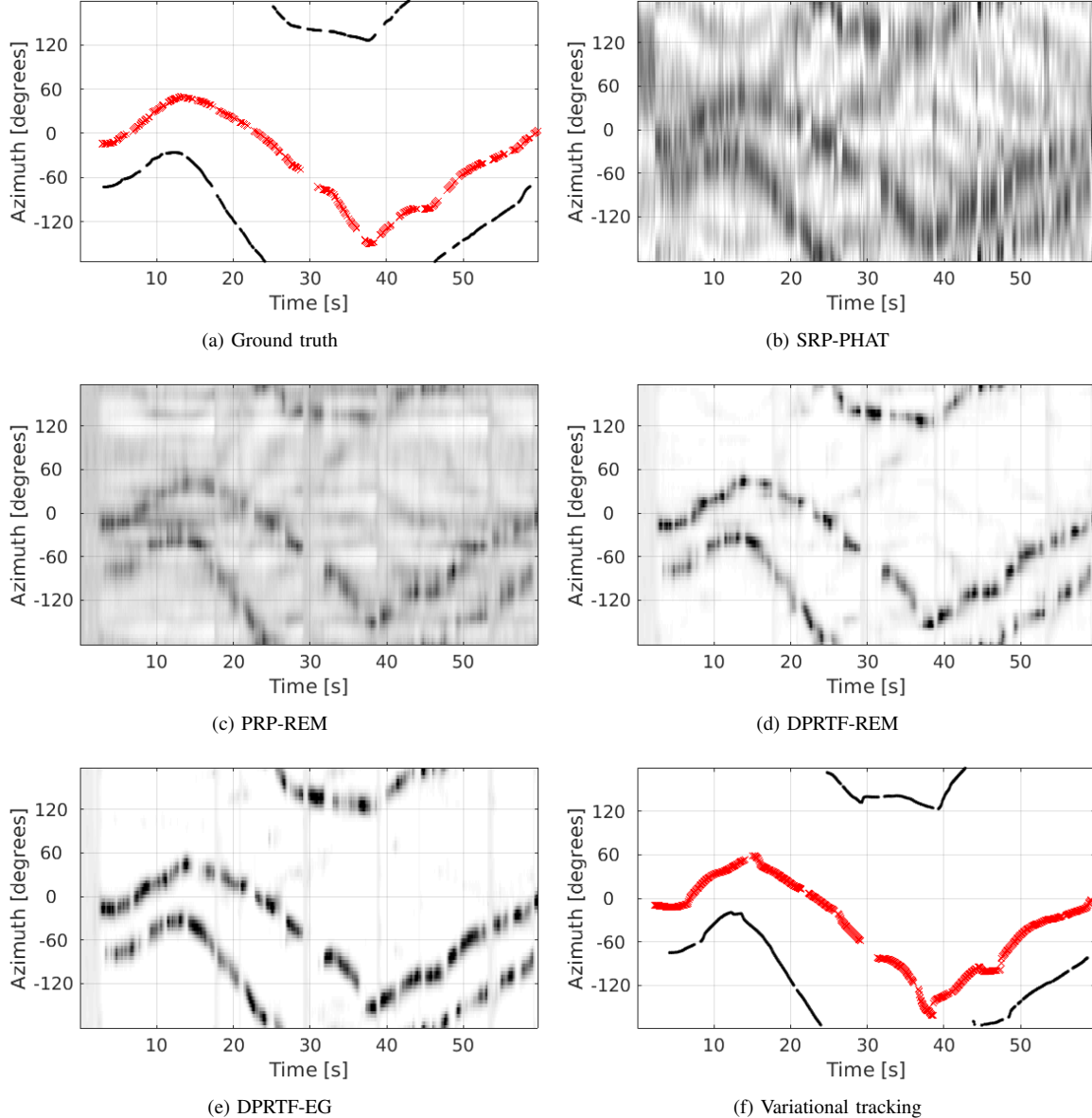


Fig. 1: Results of speaker localization and tracking for Recording 1 / Task 6 of LOCATA data. (a) Ground truth trajectory and voice activity (red for speaker 1, black for speaker 2). Intervals in the trajectories are speaking pauses. (b)-(e) One-dimensional heatmaps as a function of time for the four tested localization methods. (f) Results for the proposed VEM-based tracker. Black and red colors demonstrate a successful tracking, i.e. continuity of the tracks despite of speech pauses.

- A method combining PRP features, CGMM model and parameter update using recursive EM [16], referred to as PRP-REM. We also combine the DPRTF features and CGMM with the recursive EM (a combination referred to as DPRTF-REM). This is to evaluate the proposed DPRTF feature w.r.t. PRP, and the EG-based online parameters update method w.r.t. REM. For both baselines, the STFT and CGMM settings are the same as for the proposed method. The updating factor for recursive EM is set to 0.065.

4) *Evaluation Metrics*: The detected speakers should be assigned to the actual speakers for performance evaluation. This is done using a greedy matching algorithm. First the azimuth difference for all possible detected-actual speaker

pairs are computed, then the detected-actual speaker pair with the smallest difference is picked out as a matched pair. This procedure is iterated until the detected or actual speakers are all picked out. For each matched pair, the detected speaker is then considered to be successfully localized if the azimuth difference is not larger than 15° . The absolute error is calculated for the successfully localized sources. The mean absolute error (MAE) is computed by averaging the absolute error of all speakers and frames. For the unsuccessful localizations, we count the miss detections (MD) (speaker active but not detected) and false alarms (FA) (speaker detected but not active). Then the MD and FA rates are computed, using all the frames, as the percentage of the total MDs and FAs out of the total number of actual speakers, respectively. In addition to

TABLE I: Localization and tracking results for the LOCATA data.

	MD rate (%)	FA rate (%)	MAE ($^{\circ}$)	IDs
SRP-PHAT	37.1	23.6	5.2	-
PRP-REM	30.9	19.6	5.0	-
DPRTF-REM	23.1	15.5	4.6	-
DPRTF-EG	24.1	12.7	4.0	-
Variational Tracker	22.7	12.4	4.1	10

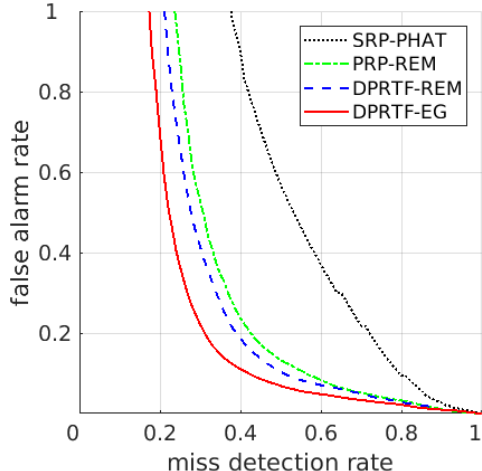


Fig. 4: ROC curve for the Kinovis dataset.

these localization metrics, we also count the identity switches (IDs) to evaluate the tracking continuity. IDs is an absolute number. It represents the number of the identity changes in the tracks for a whole test sequence.

B. Results for LOCATA Dataset

For convenience, the spatial spectrum of SRP-PHAT and the CGMM component weights profile will be referred to as heatmaps in the following. Fig. 1 shows an example of result for a LOCATA sequence. Two speakers are moving and continuously speaking with short pauses. The SRP-PHAT heatmap (Fig. 1 (b)) is cluttered due to the non-ideal beampattern of the microphone array and the influence of reverberation and noise. For most of the time, SRP-PHAT has prominent response power for the true speaker directions. Localization of the most dominant speaker can be made by selecting the direction with the largest response power. However, it is difficult to correctly count the number of active speakers and localize less dominant speakers, since there exist a number of spurious peaks. PRP-REM (Fig. 1 (c)) exhibits a clearer heatmap compared to SRP-PHAT, but there exist some spurious trajectories as well, since the PRP features are contaminated by reverberation. DPRTF-REM (Fig. 1 (d)) removes most of the spurious trajectories, which illustrates the robustness of the proposed DP-RTF feature against reverberation. From Fig. 1 (e), it can be seen that the proposed EG algorithm further removes the interferences by applying the entropy regularization. In addition, the peak evolution is smoother compared with Fig. 1 (d), which is mainly due to the use of the spatial smoothing. Fig. 1 (f) illustrates the result obtained with the proposed VEM tracker,

TABLE II: Localization and tracking results for the Kinovis dataset.

	MD rate (%)	FA rate (%)	MAE ($^{\circ}$)	IDs
SRP-PHAT	64.6	30.0	5.4	-
PRP-REM	45.2	17.2	5.1	-
DPRTF-REM	42.0	16.3	5.4	-
DPRTF-EG	35.4	14.4	5.3	-
Variational Tracker	31.1	11.7	4.9	11

with DPRTF-EG providing the observations. The proposed tracker gives smoother and cleaner results compared with the other methods. Even when the observations have a low weight, the tracker is still able to give the correct speaker trajectories. This is ensured by the second term in (35) which exploits the source dynamics model and continues to provide localization information even when w_{td} (and/or α_{tdn}) becomes small. As a result, the tracker is able to preserve the identity of speakers in spite of the (short) speech pauses. In the presented sequence example, the estimated speaker identities are quite consistent with the ground truth.

To quantitatively evaluate the quality of the heatmaps provided by the tested localization methods, we plot in Fig. 2 the receiver operating characteristic (ROC) curve (MD rate versus FA rate) obtained on the LOCATA dataset by varying the peak selection threshold, for each tested method. For each curve, a good balance between FA rate and MD rate is achieved at the left-bottom corner, i.e. the point of the curve having the maximum curvature. The average localization results corresponding to this optimal left-bottom point are summarized in Table I for each tested method. It can be seen that DPRTF-REM and DPRTF-EG achieve notably better performance than SRP-PHAT and PRP-REM according to all measures (MD, FA and, to a lower extent, MAE), since the heatmaps for SRP-PHAT and PRP-REM are much more corrupted by the reverberations. DPRTF-EG has a %1 higher MD rate over DPRTF-REM, but it has a nearly 3% lower FA rate, and 0.6° lower MAE, mainly because of its smoother peak evolution. The proposed tracker performs the best in terms of MD and FA rates. The tracker slightly reduces the FA rate compared to DPRTF-EG alone mainly by eliminating some spurious peaks that are present in the DPRTF-EG outputs. It also reduces the MD rate since some correct speaker trajectories can be recovered even when the observations have (very) weak weights, as explained above. In addition, the proposed tracker achieves quite consistent speaker ID estimation. For the whole LOCATA dataset, only 10 identity switches were observed. These identity switches are mainly due to the crossing of speaker trajectories, a hard case for the source dynamics model.

C. Results for Kinovis Dataset

Fig. 3 shows an example of result for a Kinovis sequence. Three participants are moving and intermittently speaking. It can be seen that, for many frames, the response power of SRP-PHAT and the CGMM component weights of PRP-REM corresponding to the true active speakers are not prominent, compared to the spurious trajectories. Again, DPRTF-REM and DPRTF-EG provide much better heatmaps, though they

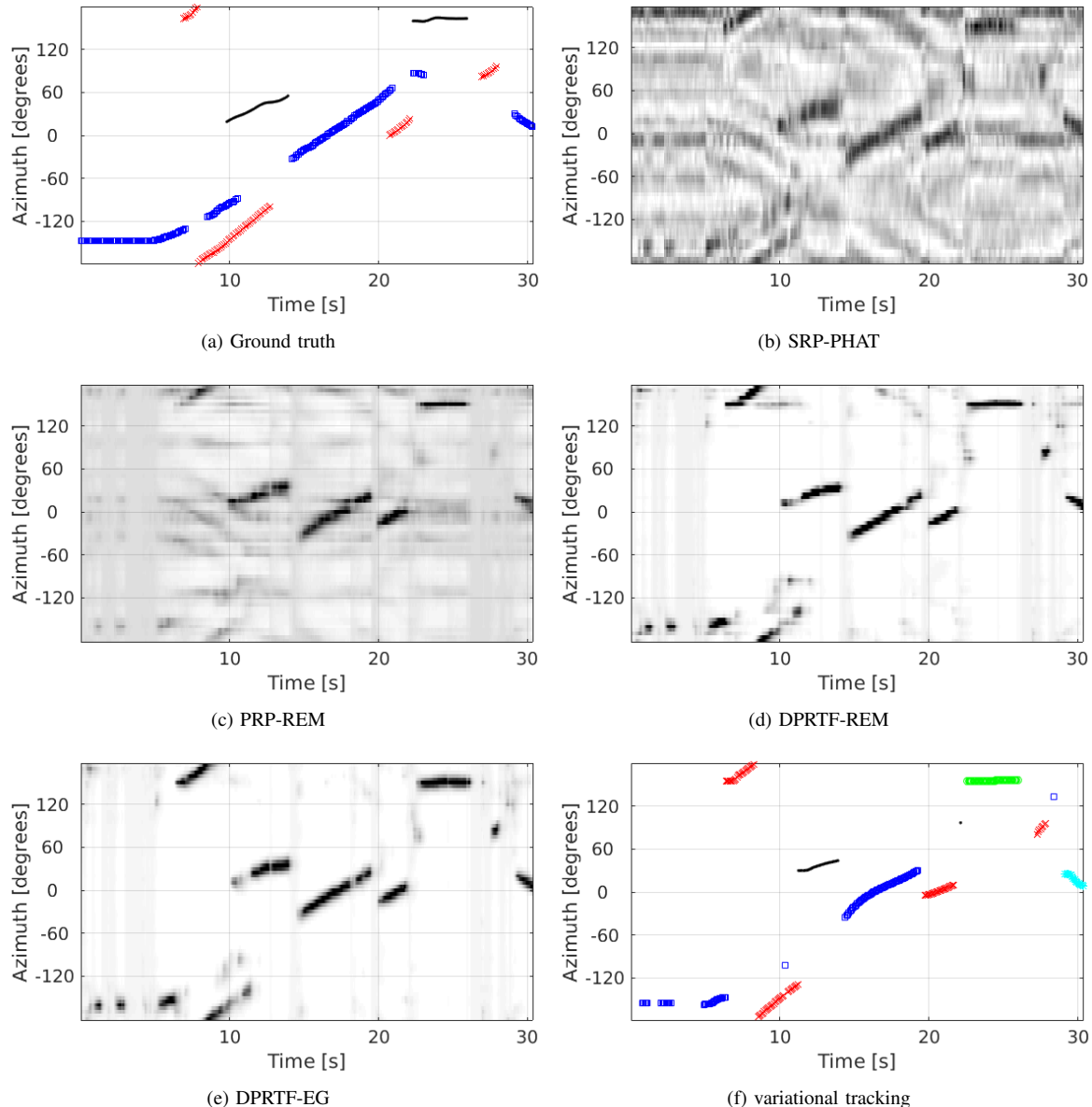


Fig. 3: Results of speaker localization and tracking for one sequence of the Kinovis dataset. (a) Ground truth trajectory and voice activity (red for speaker 1, black for speaker 2, blue for speaker 3). (b)-(e) One-dimensional heatmaps as a function of time for the four tested localization methods. (f) Results for the proposed VEM-based tracker.

also miss some speaking frames, e.g. at the beginning of Speaker 3’s trajectory (in blue). The possible reasons are i) the NAO robot (v5) has a relative strong ego-noise [39], and thus the signal-to-noise ratio of the recorded signals is relative low, and ii) the speakers are moving with a varying source-to-robot distance and the direct-path speech is contaminated by more reverberations when the speakers are distant. Overall, DPRTF-REM and DPRTF-EG are able to monitor the moving, appearance, and disappearance of active speakers for most of the time, with a small time lag due to the temporal smoothing.

This kind of recording/scenario is very challenging for the tracking method, especially for speaker identity preservation, since the participants are intermittently speaking and moving. In a general manner, the proposed tracker achieves relatively

good results, as illustrated in Fig. 3 (f). The tracked trajectories are smooth and clean. If the true trajectory of one speaker has an approximately constant direction, the tracker is able to re-identify the speaker even after a long silence thanks to the above-mentioned combination of observations and dynamics in (35), e.g. Speaker 1’s trajectory in red. In the case that the speaker changes his/her movement when he/she is silent, the track can be lost. When the person speaks again, it is indeed difficult to re-identify him/her based on the dynamics estimated before the silence period. The tracker may then prefer to give birth to a new speaker. This is illustrated by the black trajectory turning into green, and the blue trajectory turning into cyan in Fig. 3. Note that the silence periods are here much longer than in the LOCATA example of Fig. 1.

Fig 4 show the ROC curves for the Kinovis dataset. Compared to the ROC curves for the LOCATA dataset, all the four localization methods have a worse ROC curve, especially along the MD rate axis, for the reasons mentioned above. Table II summarizes the localization and tracking results for the optimal bottom-left point of the ROC curves. It can be seen that, for the four localization methods, MAEs are quite close, namely the heatmap peaks have similar biases. Compared with the results for the LOCATA dataset, the advantage of the proposed tracker is more significant for the Kinovis dataset. The MD and FA rates are more largely reduced compared with DPRTF-EG (respectively by 4.3% and 2.7%), and MAE is reduced by 0.4° . The identity switches are mainly caused by speakers changing their movement while being silent, as discussed above.

VI. CONCLUSION

In this paper, we proposed and combined i) a recursive DP-RTF feature estimation method, ii) an online multiple-speaker localization method, and iii) an multiple-speaker tracking method. The resulting framework provides online speaker counting, localization and consistent tracking (i.e. preserving speaker identity over a track in spite of intermittent speech production). The three algorithms are computationally efficient. In particular the tracking algorithm implemented in variational Bayesian framework yields a tractable solver under the form of VEM. Experiments with two datasets, recorded in realistic environment, verify that the proposed method is robust against reverberation and noise. Moreover, the tracker is able to efficiently track multiple moving speakers, detect whether there are speech or they are silent, as long as the motion associated with silent people is smooth. However, the tracking of the person from silent to active remains a difficult task. The combination of the proposed method with speaker identification techniques will be addressed in future works to attempt to solve this problem.

The proposed VEM tracker can be easily adapted to work in tandem with any frame-wise localizer providing source location estimates and/or corresponding weights (and if no weights are provided by the localizer, the tracker can be applied with all weights set to one). This makes the proposed tracker very flexible, and easily reusable by the audio processing research community.

REFERENCES

- [1] J. Chen, J. Benesty, and Y. Huang, "Time delay estimation in room acoustic environments: an overview," *EURASIP Journal on applied signal processing*, vol. 2006, pp. 170–170, 2006.
- [2] C. Knapp and G. C. Carter, "The generalized correlation method for estimation of time delay," *IEEE Transactions on Acoustics, Speech and Signal Processing*, vol. 24, no. 4, pp. 320–327, 1976.
- [3] J. H. DiBiase, H. F. Silverman, and M. S. Brandstein, "Robust localization in reverberant rooms," in *Microphone Arrays* (M. S. Brandstein and D. Ward, eds.), pp. 157–180, Springer, 2001.
- [4] C. T. Ishi, O. Chatot, H. Ishiguro, and N. Hagita, "Evaluation of a music-based real-time sound localization of multiple sound sources in real noisy environments," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2027–2032, 2009.
- [5] O. Yilmaz and S. Rickard, "Blind separation of speech mixtures via time-frequency masking," *IEEE Transactions on Signal Processing*, vol. 52, no. 7, pp. 1830–1847, 2004.
- [6] M. I. Mandel, R. J. Weiss, and D. P. Ellis, "Model-based expectation-maximization source separation and localization," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, no. 2, pp. 382–394, 2010.
- [7] Y. Dorfan and S. Gannot, "Tree-based recursive expectation-maximization algorithm for localization of acoustic sources," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 23, no. 10, pp. 1692–1703, 2015.
- [8] Y. Huang and J. Benesty, "Adaptive multichannel time delay estimation based on blind system identification for acoustic source localization," in *Adaptive Signal Processing*, pp. 227–247, Springer, 2003.
- [9] S. Doclo and M. Moonen, "Robust adaptive time delay estimation for speaker localization in noisy and reverberant acoustic environments," *EURASIP Journal on Applied Signal Processing*, vol. 2003, pp. 1110–1124, 2003.
- [10] T. G. Dvorkind and S. Gannot, "Time difference of arrival estimation of speech source in a noisy and reverberant environment," *Signal Processing*, vol. 85, no. 1, pp. 177–204, 2005.
- [11] K. Kowalczyk, E. A. Habets, W. Kellermann, and P. A. Naylor, "Blind system identification using sparse learning for TDOA estimation of room reflections," *IEEE Signal Processing Letters*, vol. 20, no. 7, pp. 653–656, 2013.
- [12] X. Li, L. Girin, R. Horaud, and S. Gannot, "Multiple-speaker localization based on direct-path features and likelihood maximization with spatial sparsity regularization," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 25, no. 10, pp. 1997–2012, 2017.
- [13] Y. Avargel and I. Cohen, "System identification in the short-time Fourier transform domain with crossband filtering," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 4, pp. 1305–1319, 2007.
- [14] R. Talmon, I. Cohen, and S. Gannot, "Relative transfer function identification using convolutive transfer function approximation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 17, no. 4, pp. 546–555, 2009.
- [15] D. Pavlidi, A. Griffin, M. Puigt, and A. Mouchtaris, "Real-time multiple sound source localization and counting using a circular microphone array," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 21, no. 10, pp. 2193–2206, 2013.
- [16] O. Schwartz and S. Gannot, "Speaker tracking using recursive EM algorithms," *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, vol. 22, no. 2, pp. 392–402, 2014.
- [17] N. Roman and D. Wang, "Binaural tracking of multiple moving sources," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 16, no. 4, pp. 728–739, 2008.
- [18] C. Evers, A. H. Moore, P. A. Naylor, J. Sheffer, and B. Rafaely, "Bearing-only acoustic tracking of moving speakers for robot audition," in *IEEE International Conference on Digital Signal Processing (DSP)*, pp. 1206–1210, 2015.
- [19] Y. Ban, L. Girin, X. Alameda-Pineda, and R. Horaud, "Exploiting the complementarity of audio and visual data in multi-speaker tracking," in *ICCV Workshop on Computer Vision for Audio-Visual Media*, vol. 3, 2017.
- [20] Z. Liang, X. Ma, and X. Dai, "Robust tracking of moving sound source using multiple model kalman filter," *Applied acoustics*, vol. 69, no. 12, pp. 1350–1355, 2008.
- [21] J. Vermaak and A. Blake, "Nonlinear filtering for speaker tracking in noisy and reverberant environments," in *Acoustics, Speech, and Signal Processing, 2001. Proceedings.(ICASSP'01). 2001 IEEE International Conference on*, vol. 5, pp. 3021–3024, IEEE, 2001.
- [22] S. Ba, X. Alameda-Pineda, A. Xompero, and R. Horaud, "An online variational bayesian model for multi-person tracking from cluttered scenes," *Computer Vision and Image Understanding*, vol. 153, pp. 64–76, 2016.
- [23] I. Gebru, S. Ba, X. Li, and R. Horaud, "Audio-visual speaker diarization based on spatiotemporal Bayesian fusion," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2017.
- [24] M. F. Fallon and S. J. Godsill, "Acoustic source localization and tracking of a time-varying number of speakers," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 20, no. 4, pp. 1409–1415, 2012.
- [25] J.-M. 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.

- [26] V. Cevher, R. Velmurugan, and J. H. McClellan, "Acoustic multitarget tracking using direction-of-arrival batches," *IEEE Transactions on Signal Processing*, vol. 55, no. 6, pp. 2810–2825, 2007.
- [27] B.-N. Vo, S. Singh, and W. K. Ma, "Tracking multiple speakers using random sets," in *Acoustics, Speech, and Signal Processing, 2004. Proceedings.(ICASSP'04). IEEE International Conference on*, vol. 2, pp. ii–357, IEEE, 2004.
- [28] W.-K. Ma, B.-N. Vo, S. S. Singh, and A. Baddeley, "Tracking an unknown time-varying number of speakers using tdoa measurements: A random finite set approach," *IEEE Transactions on Signal Processing*, vol. 54, no. 9, pp. 3291–3304, 2006.
- [29] B.-N. Vo and W.-K. Ma, "The gaussian mixture probability hypothesis density filter," *IEEE Transactions on signal processing*, vol. 54, no. 11, pp. 4091–4104, 2006.
- [30] X. Li, B. Mourgue, L. Girin, S. Gannot, and R. Horaud, "Online localization of multiple moving speakers in reverberant environments," in *The Tenth IEEE Workshop on Sensor Array and Multichannel Signal Processing*, 2018.
- [31] J. Kivinen and M. K. Warmuth, "Exponentiated gradient versus gradient descent for linear predictors," *Information and Computation*, vol. 132, no. 1, pp. 1–63, 1997.
- [32] G. Xu, H. Liu, L. Tong, and T. Kailath, "A least-squares approach to blind channel identification," *IEEE Transactions on signal processing*, vol. 43, no. 12, pp. 2982–2993, 1995.
- [33] X. Li, L. Girin, R. Horaud, and S. Gannot, "Estimation of relative transfer function in the presence of stationary noise based on segmental power spectral density matrix subtraction," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 320–324, 2015.
- [34] A. L. Yuille and A. Rangarajan, "The concave-convex procedure," *Neural computation*, vol. 15, no. 4, pp. 915–936, 2003.
- [35] I. D. Gebru, X. Alameda-Pineda, F. Forbes, and R. Horaud, "Em algorithms for weighted-data clustering with application to audio-visual scene analysis," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 12, pp. 2402–2415, 2016.
- [36] H. W. Löllmann, C. Evers, A. Schmidt, H. Mellmann, H. Barfuss, P. A. Naylor, and W. Kellermann, "The LOCATA challenge data corpus for acoustic source localization and tracking," in *IEEE Sensor Array and Multichannel Signal Processing Workshop (SAM)*, (Sheffield, UK), July 2018.
- [37] X. Li, R. Horaud, L. Girin, and S. Gannot, "Voice activity detection based on statistical likelihood ratio with adaptive thresholding," in *IEEE International Workshop on Acoustic Signal Enhancement (IWAENC)*, pp. 1–5, 2016.
- [38] <https://kinovis.inria.fr/inria-platform>.
- [39] X. Li, L. Girin, F. Bader, and R. Horaud, "Reverberant sound localization with a robot head based on direct-path relative transfer function," in *IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 2819–2826, IEEE, 2016.