

Estimation of a Nonvisible Field-of-View Mobile Target Incorporating Optical and Acoustic Sensors

Kuya Takami^{*1} · Tomonari Furukawa^{*1,3} · Makoto Kumon^{*2} · Daisuke Kimoto^{*2} · Gamini Dissanayake^{*3}

Received: date / Accepted: date

Abstract This paper presents a nonvisible field-of-view (NFOV) target estimation approach that incorporates optical and acoustic sensors. An optical sensor can accurately localize a target in its field-of-view (FOV) whereas the acoustic sensor could estimate the target location over a much larger space but only with limited accuracy. A recursive Bayesian estimation framework where observations of the optical and acoustic sensors are probabilistically treated and fused is proposed in this paper. A technique to construct the observation likelihood when two microphones are used as the acoustic sensor is also described. The proposed technique derives and stores the interaural level difference of observations from the two microphones for different target positions in advance and constructs the likelihood through correlation. A parametric study of the proposed acoustic sensing technique in a controlled test environment, and experiments with an NFOV target in an actual indoor environment are presented to demonstrate the capability of the proposed technique.

Keywords nonvisible field-of-view target estimation · recursive Bayesian estimation · interaural level difference · acoustic localization

^{*1} K. Takami and T. Furukawa
Department of Mechanical Engineering, Virginia Tech,
Blacksburg, VA, USA
E-mail: {kuya, furukawa}@vt.edu

^{*2} M. Kumon and D. Kimoto
Department of Mechanical System Engineering, Kumamoto
University, Kumamoto, Japan
E-mail: kumon@gpo.kumamoto-u.ac.jp

^{*3} T. Furukawa and G. Dissanayake
Center for Autonomous Systems, University of Technology,
Sydney, NSW, Australia
E-mail: gamini.dissanayake@uts.edu.au

1 Introduction

Mobile target estimation, or more strictly target localization and tracking, has been a research challenge over several decades due to the existence of a variety of applications in addition to the significance and the difficulty of each application. Indoor mobile target estimation, which is the focus of this paper, sees its significance in applications such as home security, home health care and urban search-and-rescue whilst its difficulty is primarily governed by the complexity of indoor structures [1, 14]. Complex indoor structures make estimation problems challenging as they could introduce largely unobservable regions when an optical sensor such as a camera is deployed [19]. This is because optical sensor's field-of-view (FOV) is determined by the range of the optical sensor and the line-of-sight (LOS) from the optical sensor, which could be small in highly constrained environments. This gives rise to need for nonvisible FOV (NFOV) mobile target estimation.

The coverage of a larger area can be achieved in four different ways. First is the most straightforward approach, which is to expand the optical FOV. It can be done by either implementing a single optical sensor with a larger FOV, such as an omnidirectional camera [17, 26], or multiple optical sensors each covering a different FOV, such as a spherical camera [12]. While such implementations are suitable for some applications, these approaches may still need to cope with a large unobservable area if environments are cluttered and complex since the FOV is determined by the LOS from the optical sensor(s). This approach, thus, is not essentially an adequate solution for the NFOV mobile target estimation.

The second approach deploys radio-frequency (RF) transmitters and receivers where an RF transmitter is

mounted on a target of concern. In one arrangement, RF receivers form a wireless sensor network (WSN), and numerical techniques are used to localize an NFOV target by processing information of received signals such as signal intensity [2, 5, 27, 37, 20, 10, 11]. For example, Wang *et al.* [36] developed and applied a Bayesian filter using the chirp-spread-spectrum ranging. The achievement of sub-meter accuracy has been reported, but the accuracy significantly depends on the settings of the networking infrastructure including the location and the number of wireless sensor nodes [13, 9]. A more accurate arrangement with minimal infrastructure uses “fingerprints” [1, 18]. In a static environment, there is a unique fingerprint at each location. A target can thus be localized by feature-matching the fingerprints. Whilst this arrangement could achieve high accuracy, the critical problem inherent in the RF-based approach is its applicability to only targets with an RF transmitter [4, 29, 33, 14].

In the third, acoustic sensors are used for target estimation. The acoustic approach can be advantageously used for any target that can make sound unlike the RF-based approach. Further, if the target is cooperative, it is possible to communicate with the target and estimate its location via sound. The approach most commonly utilizes the time-of-arrival (TOA) information of acoustic signals. Chan *et al.* [3] and Riba and Urruela [30] localized a target by positioning three or more LOS acoustic sensors. Furthermore, Nakadai *et al.* [24] and Sasaki *et al.* [31] demonstrated LOS mobile robot localization using microphone array. Mak and Furukawa [22] considered the diffraction characteristics of low-frequency sound and estimated target location. The former is an acoustic version of WSN, but the reflective nature of sound waves confines acoustic sensor locations to LOS, which is not possible in many occasions. The latter enables NFOV target estimation, but the time of sound generation, which is often unknown, must be informed beforehand. The majority of sound localization challenges have been focused on the direction of sound rather than its position due to the complexity of sound wave propagation [35, 34, 21]. Recently, acoustic physical models are used to improve the localization capability by investigating the reflection of sound in a known environment. Narang *et al.* [25] detected reflected sound by a combination of an image model and dynamic environment map. This approach estimates the direction of the sound source, however, it does not estimate the target position in NFOV, and does not fuse the vision information for localization. Even *et al.* [6] localized the sound source by ray tracing method based on the reflection signal arrival directions. However, their approach resulted in NFOV estimation

with meter order accuracy and some inconsistent variation in trials.

The last approach enhances the NFOV target estimation with a sensor having a limited FOV such as an optical sensor by using a numerical technique. Mauler [23] stated the NFOV estimation problem mathematically, and Furukawa, *et al.* [7, 8] developed the so-called Bayesian Search and Tracking (SAT) technique as a generalized numerical solution. In this technique, the event of “no detection” is converted into an observation likelihood and utilized to positively update probabilistic belief on the target that is dynamically maintained by the recursive Bayesian estimation (RBE). While search with no detection is made possible in addition to tracking with detection, the search has been found to fail unless the target is re-discovered within a short period after being lost. The unreliability of the belief significantly grows when there is no detection. It is, therefore, preferable to use any information on the target if available rather than purely relying on the information contained in no detection events.

This paper presents an NFOV target estimation approach, which incorporates an acoustic sensor in addition to an optical sensor. Whilst the optical sensor accurately localizes a target when the target is in the FOV, the acoustic sensor is used to estimate the target location even if the target is not in the optical FOV. The estimation is performed within the RBE framework where observations of the optical and acoustic sensors are each converted into an observation likelihood and making it possible to compute a joint observation likelihood via sensor fusion. Although the acoustic observation likelihood could be multi-modal with high uncertainty, the target belief updated in the past with sharply unimodal optical likelihoods when the target was in the optical FOV effectively acts as strong prior knowledge and allows accurate NFOV target estimation. A technique to construct an observation likelihood using an acoustic sensor composed of two microphones is also proposed. The proposed technique derives the interaural level difference (ILD) of observations from the two microphones for different target positions and stores the ILDs as fingerprints, or acoustic cues, *a priori*. Given a new acoustic observation, an acoustic observation likelihood is computed by quantifying the correlation of the ILD of the new observation to the stored ILDs. While the acoustic cues must be created in advance, the technique achieves the truly NFOV target estimation.

The paper is organized as follows. The following section reviews the conventional RBE that uses an optical sensor as well as the grid-based method. Section 3 presents the proposed target estimation approach incorporating an acoustic sensor. Section 4 demonstrates

the efficacy of the proposed target estimation through experimental analysis, and conclusions are summarized in the final section.

2 Optical Recursive Bayesian Estimation

2.1 Target Motion Model and Optical Sensor Model

Consider a target t of concern, the motion of which is given by

$$\mathbf{x}_{k+1}^t = \mathbf{f}^t(\mathbf{x}_k^t, \mathbf{u}_k^t, \mathbf{w}_k^t) \quad (1)$$

where $\mathbf{x}_k^t \in \mathcal{X}^t$ is the state of the target at time step k , $\mathbf{u}_k^t \in \mathcal{U}^t$ is the set of control inputs of the target, and $\mathbf{w}_k^t \in \mathcal{W}^t$ is the “system noise” of the target. For simplicity, the target state describes the two-dimensional position.

In order for the formulation of the NFOV target estimation problem, this moving target is observed by a sensor platform s . To focus on the estimation of a mobile target, let the sensor platform be stationary and its global state be accurately known as $\tilde{\mathbf{x}}^s \in \mathcal{X}^s$. Note that $()$ is an instance of $()$. The sensor platform carries an optical sensor to observe the target. The FOV, or more precisely the “observable region”, of the optical sensor s_c can be expressed with the probability of detecting the target $P_d(\mathbf{x}_k^t|\tilde{\mathbf{x}}^s)$ as

$${}^{s_c}\mathcal{X}_o^t = \{\mathbf{x}_k^t | 0 < P_d(\mathbf{x}_k^t|\tilde{\mathbf{x}}^s) \leq 1\}.$$

Accordingly, the target position observed from the optical sensor, ${}^{s_c}\mathbf{z}_k^t \in \mathcal{X}^t$, is given by

$${}^{s_c}\mathbf{z}_k^t = \begin{cases} {}^{s_c}\mathbf{h}^t(\mathbf{x}_k^t, \tilde{\mathbf{x}}^s, {}^{s_c}\mathbf{v}_k^t) & \mathbf{x}_k^t \in {}^{s_c}\mathcal{X}_o^t \\ \emptyset & \mathbf{x}_k^t \notin {}^{s_c}\mathcal{X}_o^t \end{cases} \quad (2)$$

where ${}^{s_c}\mathbf{v}_k^t$ represents the observation noise, and \emptyset represents an “empty element”, indicating that the optical observation contains no information on the target or that the target is unobservable.

2.2 Recursive Bayesian Estimation

The RBE updates belief on a dynamical system, given by a probability density, in both time and observation. Let a sequence of observations of a moving target t by a stationary sensor platform s from time step 1 to time step k be ${}^s\tilde{\mathbf{z}}_{1:k}^t \equiv \{{}^s\tilde{\mathbf{z}}_k^t | \forall k \in \{1, \dots, k\}\}$. Given the initial belief $p(\mathbf{x}_0^t)$, the sensor platform state $\tilde{\mathbf{x}}^s$ and a sequence of observations ${}^s\tilde{\mathbf{z}}_{1:k}^t$, the belief on the target at any time step k , $p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k}^t, \tilde{\mathbf{x}}^s)$, can be estimated recursively through the two stage equations; update and prediction.

In the prediction process, the target belief $p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s)$ is updated from that in the previous time step $p(\mathbf{x}_{k-1}^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s)$ by Chapman-Kolmogorov equation as

$$\begin{aligned} & p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s) \\ &= \int_{\mathcal{X}^t} p(\mathbf{x}_k^t | \mathbf{x}_{k-1}^t) p(\mathbf{x}_{k-1}^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s) d\mathbf{x}_{k-1}^t, \end{aligned} \quad (3)$$

where $p(\mathbf{x}_k^t | \mathbf{x}_{k-1}^t)$ is a probabilistic form of the motion model (1). Note that $p(\mathbf{x}_{k-1}^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s) = p(\mathbf{x}_0^t)$ when $k = 1$. The correction process, on the other hand, updates the belief using information available in the observations. The target belief $p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k}^t, \tilde{\mathbf{x}}^s)$ is corrected from the corresponding state estimated with the observations up to the previous time step $p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s)$ and a new observation ${}^s\tilde{\mathbf{z}}_k^t$ as

$$p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k}^t, \tilde{\mathbf{x}}^s) = \frac{q(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}^s)}{\int_{\mathcal{X}^t} q(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}^s) d\mathbf{x}_{k-1}^t}, \quad (4)$$

where

$$q(\cdot) = l(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}^s) p(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_{1:k-1}^t, \tilde{\mathbf{x}}^s), \quad (5)$$

and $l(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}^s)$ represents the likelihood of \mathbf{x}_k^t given ${}^s\tilde{\mathbf{z}}_k^t$ and $\tilde{\mathbf{x}}^s$, which is a probabilistic version of the sensor model; i.e., Equation (2) if the sensor is optical. It is to be noted that the likelihood does not need to be a probability density since the normalization in Equation (4) makes the output belief be a probability density regardless of the formulation of the likelihood.

2.3 Modeling of Optical Observation Likelihood

The optical observation likelihood is modeled by first defining the “detectable region”. Due to the existence of uncertainty, the observation of a no-empty element does not necessarily indicate that the target has been reliably detected. The detectable region of the optical sensor s_c that describes the region within which the sensor confidently detects the target is thus defined as:

$${}^{s_c}\mathcal{X}_d^t = \{\mathbf{x}_k^t | \epsilon^t < P_d(\mathbf{x}_k^t|\tilde{\mathbf{x}}^s) \leq 1\} \subset {}^{s_c}\mathcal{X}_o^t,$$

where ϵ^t is a positive threshold value which judges the detection of the target. Given the observation ${}^s\tilde{\mathbf{z}}_k^t$, the optical observation likelihood is resultantly stated as

$$l^c(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}^s) = \begin{cases} p({}^s\tilde{\mathbf{z}}_k^t | \mathbf{x}_k^t, \tilde{\mathbf{x}}^s) & \exists {}^s\tilde{\mathbf{z}}_k^t \in {}^{s_c}\mathcal{X}_d^t \\ 1 - P_d(\mathbf{x}_k^t|\tilde{\mathbf{x}}^s) & \nexists {}^s\tilde{\mathbf{z}}_k^t \in {}^{s_c}\mathcal{X}_d^t \end{cases} \quad (6)$$

where, depending on whether there exists a target within the detectable region, the upper and lower formulas return likelihoods with detection and no-detection events, respectively.

Figure 1 illustrates the configuration of the optical observation likelihood when a sensor platform is in a one-dimensional target space. When a target is not detected without having it in the detectable region, the likelihood tells where the target is unlikely to be and is thus represented as a heavily non-Gaussian distribution. When the target is detected, the likelihood becomes near-Gaussian with its peak located at the observed location. The closer the target to the sensor platform, the more accurate the estimation.

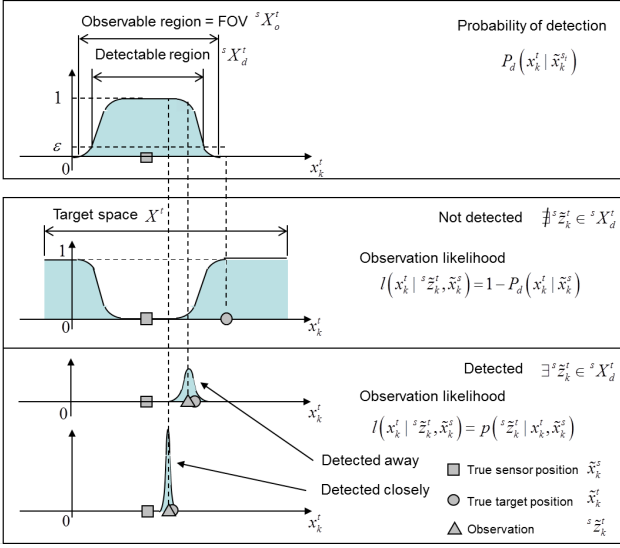


Fig. 1 Optical observation likelihood

2.4 Grid-based Method

Handling the heavily non-Gaussian no-detection likelihood necessitates the grid-based method for RBE. As the grid-based method represents the belief space in terms of regularly aligned grid cells, let the cell of concern be positioned at l th and m th partitions in x and y directions. At the grid cell $[l, m]$, the prediction and the correction are processed independently. The prediction requires the numerical evaluation of the Chapman-Kolmogorov equation in Equation (3) at each grid cell. Given the belief $p_{\mathbf{x}_{k-1}^t}^{l,m}(s\mathbf{z}_{1:k-1}^t)$ at time step k as well as the motion model $p_{\mathbf{x}_k^t|\mathbf{x}_{k-1}^t}^{l,m}$ constructed in the matrix form as the convolution kernel, the target belief at the grid cell $[l, m]$ can be predicted as

$$p_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_{1:k-1}^t) = \sum_{\alpha=0}^{I_x^t} \sum_{\beta=0}^{I_y^t} p_{\mathbf{x}_k^t|\mathbf{x}_{k-1}^t}^{\alpha,\beta} p_{\mathbf{x}_{k-1}^t}^{l-\alpha,m-\beta}(s\mathbf{z}_{1:k-1}^t) \quad (7)$$

where \otimes indicates the convolution of the last belief with the motion model.

The correction requires the computation of Equation (4) at each grid cell. Given the predicted belief $p_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_{1:k-1}^t)$ and the observation likelihood $l_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_k^t)$, the target belief at the grid cell $[l, m]$ can be corrected as

$$p_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_{1:k}^t) = \frac{q_{\mathbf{x}_k^t}^{l,m}(\cdot)}{\Delta x_r \Delta y_r \sum_{\alpha} \sum_{\beta} q_{\mathbf{x}_k^t}^{\alpha,\beta}(\cdot)}, \quad (8)$$

where

$$q_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_{1:k}^t) = l_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_k^t) p_{\mathbf{x}_k^t}^{l,m}(s\mathbf{z}_{1:k-1}^t). \quad (9)$$

and $[\Delta x_r, \Delta y_r]$ is the size of the grid.

Whilst the generalized optical observation likelihood allows belief update and maintenance regardless of whether the target has been detected, the RBE with the optical observation likelihood does not update and maintain the belief effectively. Figure 2 illustratively depicts this problem where the FOV and the NFOV are given by the light blue and the white colors respectively. When the configuration of the target space is constrained complicatedly, the FOV becomes significantly limited compared to the target space. This makes the belief dominantly updated by the observation likelihood with no detection. If no detection continuously takes place, the belief keeps spread out with predictions and becomes highly uncertain and unreliable. The next section will describe the proposed target estimation incorporating an acoustic sensor to solve this problem.

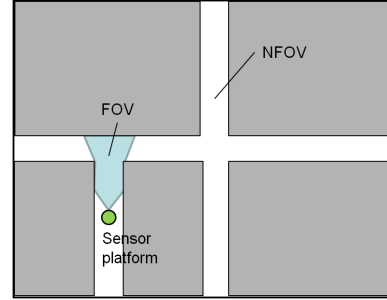


Fig. 2 Optical observation likelihood

3 Target Estimation Incorporating Optical and Acoustic Sensors

3.1 Acoustic Sensor Model and Observation Likelihood

The acoustic sensor incorporated in the proposed approach can observe a target on the non-line-of-sight (NLOS) or even in the NFOV though accuracy is limited due to the complex behavior of sound signals including reflection, refraction, and diffraction. Because

of its broad range, the observable region of the acoustic sensor could be considered unlimited when compared to that of the optical sensor. The acoustic sensor model s_a can be therefore constructed without defining an observable region unlike the optical sensor model:

$${}^{s_a}\mathbf{z}_k^t = {}^{s_a}\mathbf{h}^t(\mathbf{x}_k^t, \tilde{\mathbf{x}}_k^s, {}^{s_a}\mathbf{v}_k^t), \quad (10)$$

which is probabilistically equivalent to the likelihood given by

$$l^a(\mathbf{x}_k^t | {}^{s_a}\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^s) = p({}^{s_a}\tilde{\mathbf{z}}_k^t | \mathbf{x}_k^t, \tilde{\mathbf{x}}_k^s). \quad (11)$$

Figure 3 illustrates the observation likelihood of the acoustic sensor in comparison to that of the optical sensor in Figure 1. The observation likelihood could be heavily non-Gaussian with multiple peaks if the target is on the NLOS though it is highly likely that one of the peaks is found near the location of the target as shown in the figure. The likelihood with a target on the LOS could still be multi-modal if there are adjacent structures that create reflective, refractive and/or diffractive sound signals, but it captures the target location more confidently with a sharper peak. The likelihood with a target on the LOS without adjacent structures will be a sharp near-Gaussian distribution since the direct sound dominates the observation.

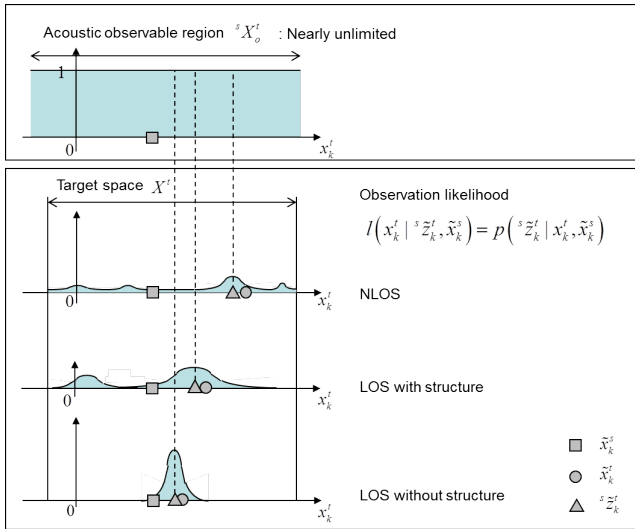


Fig. 3 Acoustic observation likelihood

3.2 RBE Using Joint Optical/Acoustic Observation Likelihood

Given the observation likelihood of the optical sensor $l^c(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^s)$ and that of the acoustic sensor

$l^a(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^s)$, the proposed approach derives a joint likelihood by multiplying the two likelihoods:

$$l(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^s) = l^c(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^s) l^a(\mathbf{x}_k^t | {}^s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^s) \quad (12)$$

in accordance to the canonical data fusion formula. The joint optical/acoustic likelihood may not be a probability density similarly to the optical and acoustic observation likelihoods.

Figure 4 illustratively shows the resulting joint optical/acoustic observation likelihood when the target is in the NFOV. The possible locations of the target are narrowed down since the optical likelihood with no detection clears out likelihood in the detectable region and dropped some peak(s) as shown in the figure. However, the joint likelihood could still remain heavily non-Gaussian with multiple peaks and thus may not solely make a good estimation about where the target is.

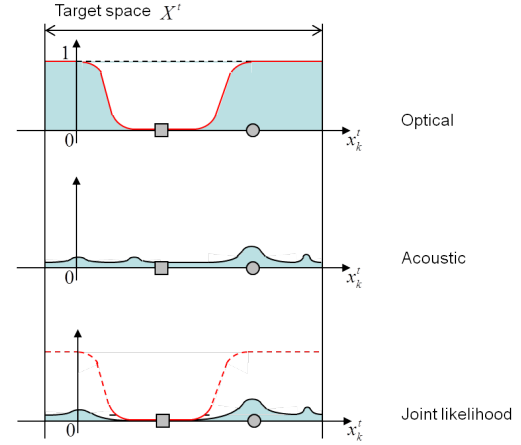


Fig. 4 Joint optical/acoustic observation likelihood

Figure 5 shows the further ability of the proposed approach in reliable target estimation. The proposed approach performs RBE by using the joint likelihood in correction (substituting Equation (12) into Equation (5)) within the standard RBE framework. Because sharpest and most Gaussian is the optical observation likelihood with detection, the prior belief is most determined by the last optical observation and remains a sharp Gaussian distribution. The posterior belief with the joint observation likelihood inherits this characteristics since the joint likelihood most likely captures the target location with a peak and magnifies the confidence of the prior belief with the joint likelihood.

3.3 Modeling of Acoustic Observation Likelihood

The technique proposed in this paper to model an acoustic observation likelihood uses two microphones as an

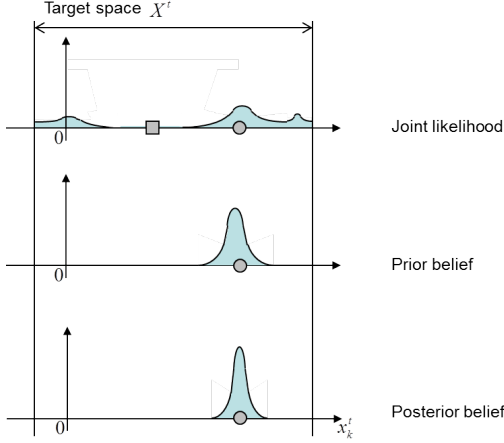


Fig. 5 RBE incorporating an acoustic sensor

acoustic sensor and constructs acoustic cues of the target in the environment of concern *a priori*. This is because of the potential of the proposed technique for NLOS target estimation through the preliminary investigations of the authors [16, 15, 28] and an inability of the aforementioned existing techniques. Figure 6 shows a schematic diagram of the proposed technique to model an acoustic observation likelihood. We assume that the target emits sound with white noise, and we indeed use it to create the acoustic observation likelihood. A white noise sound emitted at a specific position by the target for a certain time period is first recorded by two microphones. After applying fast Fourier transform (FFT), the difference between the frequency domain amplitude responses, known as the ILD, is then derived and further sampled to form an ILD vector within the frequency range of interest. The ILD vector is created with various target positions and each saved as an acoustic cue. The acoustic observation likelihood modeling essentially corresponds to creating the set of ILD vectors. When a target emitted a white noise sound, the ILD vector of the sound observation is compared to all the acoustic cues. The degree of similarity is then used to develop a correlation map indicating where the target is likely to be. The correlation map is the acoustic likelihood of the particular sound observation.

Mathematically, let the frequency-domain sound level of the target at the i th position $(\mathbf{x}_k^t)_i$, which is observed by the left and right microphones, be $s_l(\omega | (\mathbf{x}_k^t)_i)$ and $s_r(\omega | (\mathbf{x}_k^t)_i)$ where ω is a frequency of sound. The ILD for the i th position $(\mathbf{x}_k^t)_i$, $\Delta S(\omega | (\mathbf{x}_k^t)_i)$, is then given by

$$\Delta S(\omega | (\mathbf{x}_k^t)_i) = 20 \log |s_l(\omega | (\mathbf{x}_k^t)_i)| - 20 \log |s_r(\omega | (\mathbf{x}_k^t)_i)|. \quad (13)$$

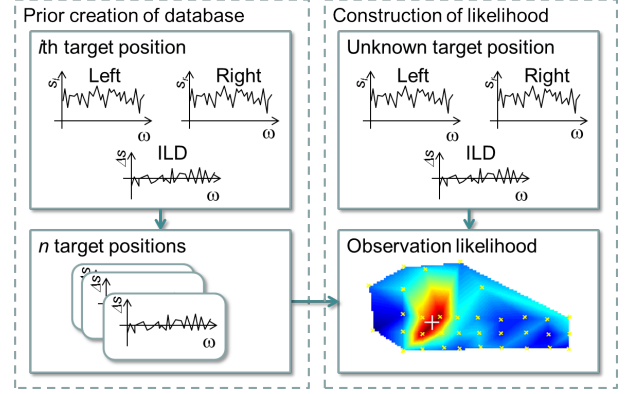


Fig. 6 Acoustic sensor

If the ILD is sampled at N frequencies $\boldsymbol{\Omega} = [\omega_1, \dots, \omega_N]^\top$, the ILD vector can be described as

$$\mathbf{S}(\boldsymbol{\Omega} | (\mathbf{x}_k^t)_i) = [a_1 \Delta S(\omega_1 | (\mathbf{x}_k^t)_i), \dots, a_N \Delta S(\omega_N | (\mathbf{x}_k^t)_i)]^\top, \quad (14)$$

where

$$a_j = \langle \min \{ |s_l(\omega_j | (\mathbf{x}_k^t)_i)|, |s_r(\omega_j | (\mathbf{x}_k^t)_i)| \} - \epsilon \rangle. \quad (15)$$

In the equation, $\langle \cdot \rangle$ is Macaulay brackets, and $\min \{ \cdot, \cdot \}$ returns the smaller value of the two entities. Thus, a_j effectively filters noise components of the signal. When the ILD vector is created with a set of frequencies $\tilde{\boldsymbol{\Omega}}$ at n known target positions, i.e., $(\tilde{\mathbf{x}}_k^t)_i, \forall i \in \{1, \dots, n\}$, the acoustic cues to be prepared in advance and used to create the acoustic observation likelihood become $\mathbf{S}(\tilde{\boldsymbol{\Omega}} | (\tilde{\mathbf{x}}_k^t)_i), \forall i \in \{1, \dots, n\}$.

Given the ILD vector $\mathbf{S}(\tilde{\boldsymbol{\Omega}} | (\tilde{\mathbf{x}}_k^t)_i)$ with observation $s\tilde{\mathbf{z}}_k^t$ at unknown target position \mathbf{x}_k^t , the proposed technique quantifies the degree of correlation of the i th ILD vector to that of the unknown target position as

$$X((\tilde{\mathbf{x}}_k^t)_i | s\tilde{\mathbf{z}}_k^t) = \frac{1}{2} \left\{ \frac{\mathbf{S}(\tilde{\boldsymbol{\Omega}} | s\tilde{\mathbf{z}}_k^t)^\top \mathbf{S}(\tilde{\boldsymbol{\Omega}} | (\tilde{\mathbf{x}}_k^t)_i)}{|\mathbf{S}(\tilde{\boldsymbol{\Omega}} | s\tilde{\mathbf{z}}_k^t)| |\mathbf{S}(\tilde{\boldsymbol{\Omega}} | (\tilde{\mathbf{x}}_k^t)_i)|} - 1 \right\}. \quad (16)$$

where $0 \leq X(\cdot) \leq 1$. The acoustic observation likelihood with the unknown target position \mathbf{x}_k^t can be finally calculated as

$$l^a(\mathbf{x}_k^t | s\tilde{\mathbf{z}}_k^t, \tilde{\mathbf{x}}_k^t) = \sum_{i=1}^n \mu_i(\boldsymbol{\xi}_k^t) X((\tilde{\mathbf{x}}_k^t)_i | s\tilde{\mathbf{z}}_k^t), \quad (17)$$

where $\mu_i(\boldsymbol{\xi}_k^t)$ is a basis function of natural coordinates $\boldsymbol{\xi}_k^t$, which are transformed from \mathbf{x}_k^t . Based on the characteristics of the likelihood, the proposed technique uses

the T-spline basis function such that $\mu_i(\xi_k^t)$ in a T-mesh, for the two-dimensional parameter space case $\xi_k^t = [\xi_k^t, \eta_k^t]^\top$, can be represented as

$$\mu_i(\xi_k^t) = \mu_i^1(\xi_k^t) \mu_i^2(\eta_k^t) \quad (18)$$

where $\mu_i^{(\cdot)}$ is a cubic B-spline basis function. Further detailed formulations are found in [32]. It is to be noted here that other basis functions are also possible while the proposed technique uses T-spline basis functions.

4 NUMERICAL AND EXPERIMENTAL ANALYSIS

The efficacy of the proposed approach was examined experimentally in two steps. The first step was aimed at studying the capability and limitation of the proposed acoustic sensing technique by parametrically changing the complexity of the environment where the experimental system with a speaker array and a movable/replaceable wall was developed specifically for this study. After verifying the feasibility of the acoustic sensing for NLOS target localization, the applicability of the proposed approach to the estimation of an NFOV target in a complex practical environment was investigated. The investigation looked into the performance of both the joint optical/acoustic observation likelihood and the RBE with the joint likelihood.

4.1 Acoustic Observation of NLOS Target

Figure 7(a) shows the design of the experimental system that changes the complexity of the environment for the evaluation of the proposed acoustic sensing technique. An acoustic sensor consisting of two microphones is fixed next to an outer wall and faces open space where a speaker array and movable/replaceable wall(s) are placed. The complexity of the environment can be changed by varying the two parameters of the movable/replaceable wall: the distance of the wall to the edge of speaker array L_d and the length of the wall L_w . The longer the distance and/or the larger the length, the more complex the environment since sound from speakers result in more reflections.

Shown in the figure as blue crosses are speaker locations. A microcontroller controls speakers so that each speaker sequentially emits white noise sound for a programmed period. A set of ILDs for a wall setting can be thus collected automatically. Once the ILDs are collected, the ability of the proposed acoustic sensing technique is evaluated by emitting sound from a speaker at some location within the area of the speaker array and

identifying the location in the form of observation likelihood. The location is not where one of the speakers of the speaker array is located to demonstrate the ability of the proposed technique in identifying the target at an arbitrary position.

Figure 7(b) shows the developed experimental system

whereas the dimensions and other parameters used in the experiments are listed in Table 1. The sound was sampled and represented at 8,192 frequency bins within the audible range to capture its behavior accurately. 54 speakers were aligned to cover the open space. The distance and the length of the wall were varied to introduce both lightly NLOS and heavily NLOS environments. The case of two walls ($n_w = 2$) was tested in addition to the single wall case to make the environment more complex. The distance of only the wall closer to the acoustic sensor was varied.

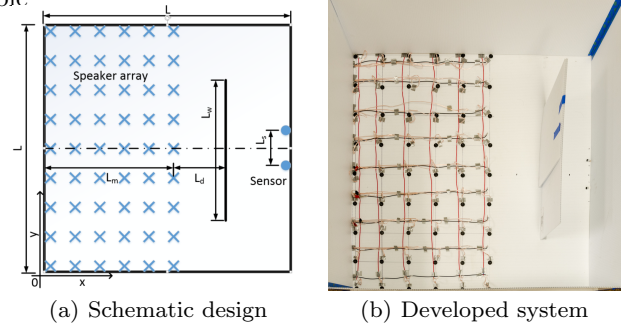


Fig. 7 Experimental system for investigating influence of environmental complexity

Table 1 Dimensions and other parameters in the experiments

Parameter	Value	Parameter	Value
$\tilde{\mathbf{x}}^t$	[42, 34] cm	L	90 cm
ω_1	0 Hz	L_m	50 cm
ω_N	22 kHz	L_s	10 cm
N	8,192	L_d	{0, 10, 20, 30} cm
ϵ	0.01	L_w	{50, 60, 70} cm
n	54	n_w	{1, 2}

Figure 8 shows the four ILDs each observed with a target at one of the 54 positions when $[L_d, L_w, n_w] = [50, 40, 1]$. Two positions are on the LOS, and four are on the NLOS. It is first seen that the configuration of the ILD varies depending on the target position. The configuration is different even when the target is on the NLOS. This indicates that the ILD contains informa-

tion on the target position no matter whether the target is on the LOS.

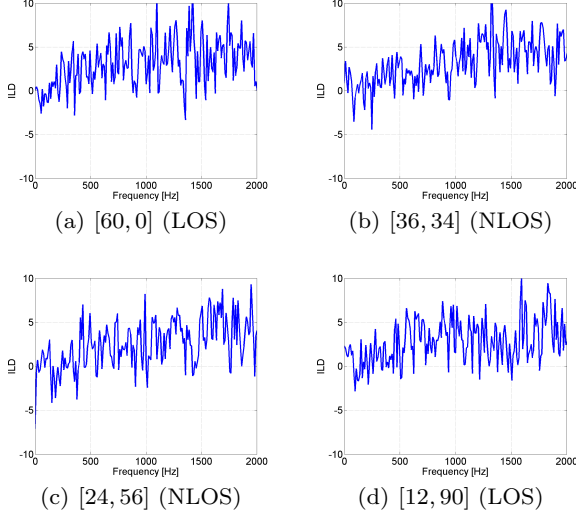


Fig. 8 ILDs at four of the 54 positions

Figure 9 shows four resulting acoustic observation likelihoods each with a different condition. The former two cases are with a single wall having different distances and the same length whereas the latter two cases are with two walls having different lengths and the same distance. The result first indicates that the target location is well estimated when the distance is short (Figure 9(a)) or when the length is small (Figure 9(c)). The target is closer to LOS in these conditions since sound reaches the acoustic sensor with a small number of reflections. The identification of the target location in the remaining two cases (Figures 9(b) and 9(d)) is hard due to a number of sound reflections. The identification with two walls (Figures 9(c) and 9(d)) is seen to be harder than that with a single wall (Figures 9(a) and 9(b)) for the same reason. While the acoustic observation likelihood is heavily multi-modal in these cases, the target location is still captured by the highest peak (Figures 9(a)-9(c)) or at least by one of the peaks (Figure 9(d)). This demonstrates the ability of the proposed acoustic sensing technique for identifying the location of the NFOV target though with limited accuracy.

Figure 10(a) and 10(b) show the mean error of the acoustic observation likelihood when the distance and the length were varied for single and double wall cases. The mean error is a distance of the nearest peak of the acoustic observation likelihood to the true target location. The result of the mean error shows that the proposed technique could locate the target within 2 cm error in most (11) of the 12 cases for the single wall

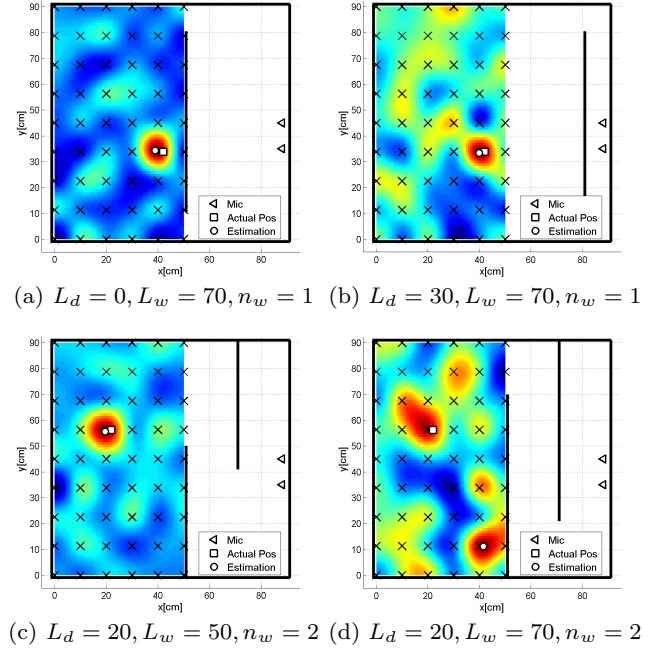


Fig. 9 Acoustic observation likelihoods for different environmental complexity

case. The estimation is particularly good when the wall length was small. The mean error with two walls, meanwhile, is over 20 cm in four of nine cases. This indicates that the proposed acoustic sensing technique is not sufficient enough to obtain less than the data points (10 cm) accuracy for NLOS localization. Figure 10(c) shows the variation of the differential entropy with respect to the number of walls and the wall length. It is seen that the addition of a wall dominantly increases uncertainty. While the proposed acoustic sensing technique is effective enough in relatively simple NLOS environments, enhancement is necessary when target estimation in more complex environments is pursued.

4.2 NFOV Target Estimation in Complex Practical Environments

This subsection investigates the applicability of the proposed approach to NFOV target estimation in complex practical environments by first enhancing the acoustic sensing with the optical sensing and then executing the RBE with the joint optical/acoustic observation likelihood. Figure 11 shows the map of the indoor environment used to demonstrate the practical applicability of the proposed approach as well as the details of the demonstration. A sensor platform with a camera and two microphones for optical and acoustic observations was located in a corridor and faced with the open space the target could move around. The environment

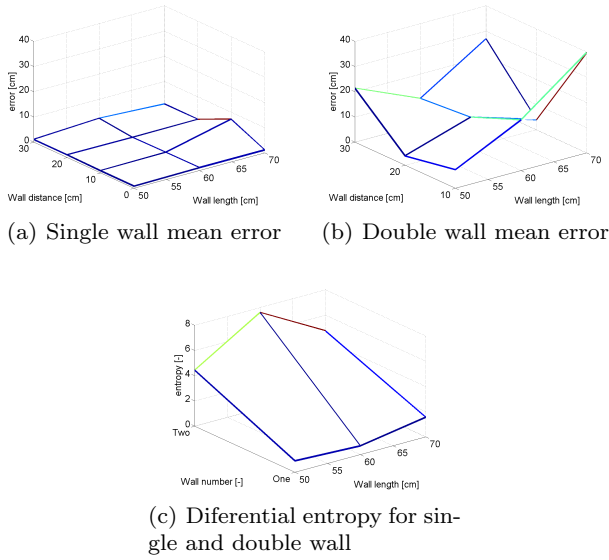


Fig. 10 Mean error and differential entropy of the acoustic observation likelihood with a single and double wall

is complex with the existence of walls and structures so that the FOV of the camera is significantly limited compared to the target space. Shown in the figure as yellow crosses are the target locations at which sound was emitted to collect ILDs. After the collection, the target was then moved along the lines indicated in the figure and emitted sound. The observation and estimation were examined at the four positions marked by red circles.

Figure 12 shows the target and the sensor platform. The sensor platform is with a camera and two microphones as aforementioned whereas the target is a wireless speaker. The same speaker was used to construct the acoustic observation likelihood. White noise was emitted from the speaker, and the parameters used to construct the acoustic observation likelihood are listed in Table 2.

Table 2 Parameters for acoustic observation likelihood

Parameter	Value
ω_1	0 [Hz]
ω_N	386 [Hz]
N	100
ϵ	0.01
n	65

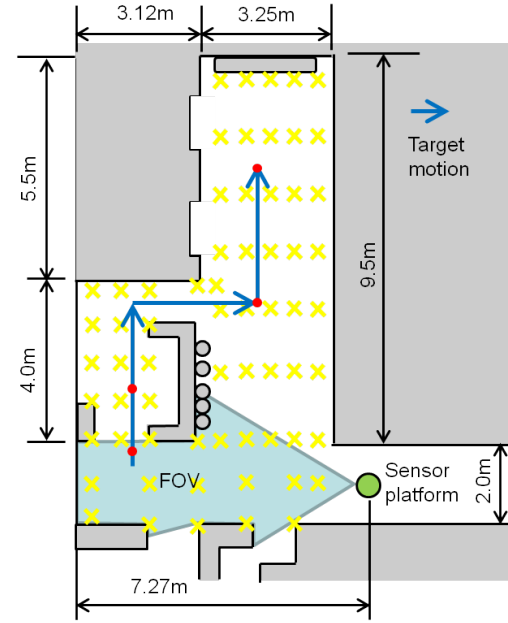


Fig. 11 Map of the test environment and demonstration



Fig. 12 Target and sensor platform

4.2.1 Joint Optical/Acoustic Observation of NFOV Target

Figure 13 shows the acoustic observation likelihoods when the target moved and emitted sound at the four marked target positions, which are at the 1st, 7th, 35th and 51th steps. The target is in the FOV only at the 1st step. The acoustic observation likelihood is seen to be multi-modal due to sound reflection even when the target is within the FOV and thus on the LOS. The proposed acoustic sensing technique has, however, been able to accurately identify the true target position at one of the peaks and successfully detect it except for the 7th step. Failure in the 7th step is a result of the limitation of acoustic sensing shown and concluded in the last subsection, but it is to be importantly noted that the target position is captured near the second highest peak.

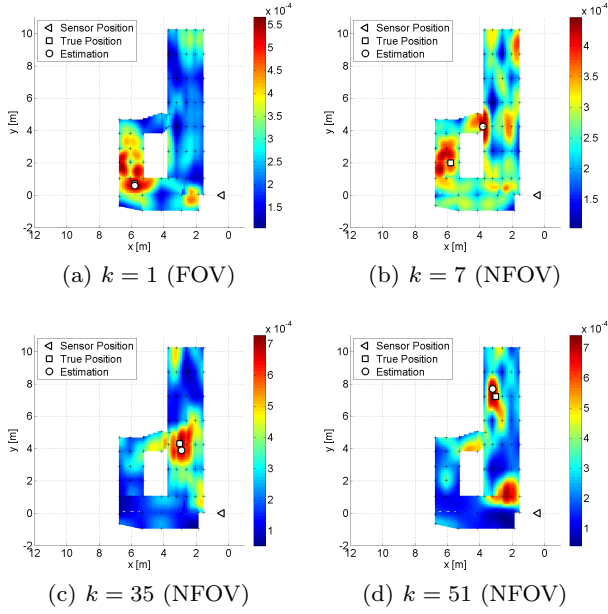


Fig. 13 Acoustic observation likelihoods

The effectiveness of the proposed acoustic sensing technique is further understood comparatively by seeing optical observation likelihoods in Figure 14. The optical sensor can identify the target accurately when it is within the FOV in Fig. 14(a). The observation likelihood with the target outside the FOV in Fig. 14(b)-14(d) can however provide no localization capability on the target. The likelihood of the target anywhere in NFOV is equally distributed with no specific estimation of the target position. Finally, Figure 15 shows the joint optical/ acoustic observation likelihoods. It is seen that the joint observation likelihoods most narrow down the possible target locations by detecting the target dominantly with the optical observation likelihood when the target is within the FOV and with the acoustic observation likelihood when the target is outside the FOV. The wrong computation at the 7th step, however, remains and necessitates RBE for target estimation.

4.2.2 RBE with Joint Optical/Acoustic Observation Likelihoods

Having understood the limitation of the target detection with observations only in the last section, the effectiveness of the proposed RBE with the joint optical/acoustic observation likelihoods was investigated using the same test data. Without knowing the target motion well, the target motion model was given by a random walk model assuming that the target is a human who could move to any direction with equal probability.

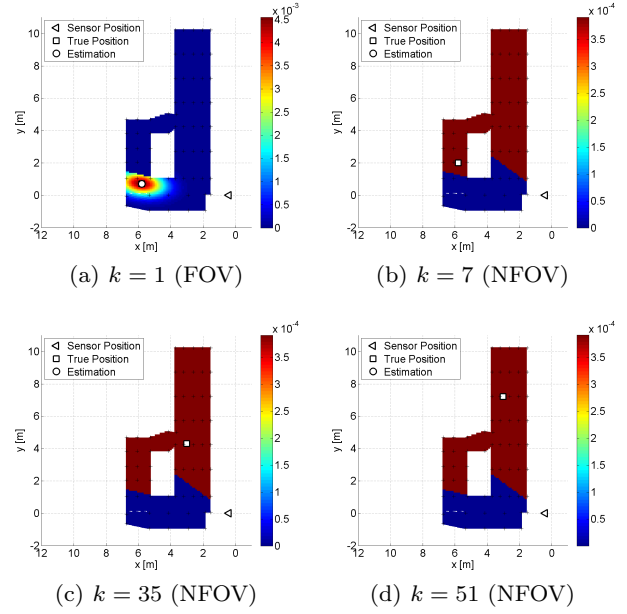


Fig. 14 Optical observation likelihoods

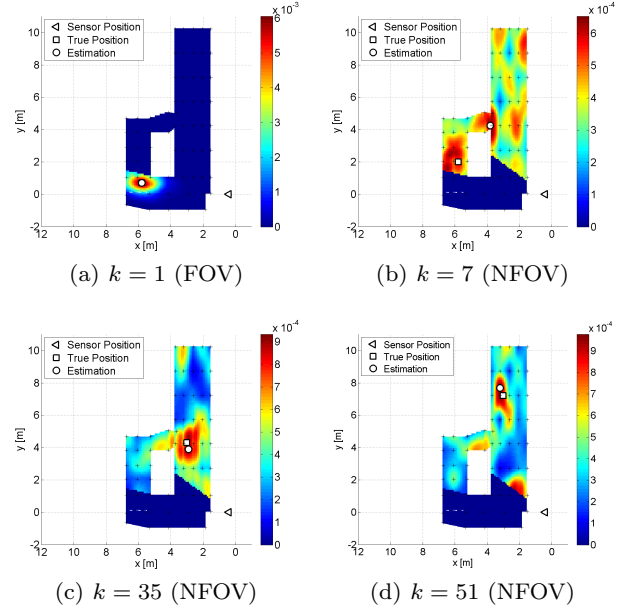


Fig. 15 Joint optical/acoustic observation likelihoods

Figure 16 shows the target belief estimated via RBE with the joint optical/acoustic observation likelihoods. The result shows that the target position is well estimated with all the time steps including the 7th step where the joint observation likelihood did not detect the target with the highest peak. Because the target was initially observed by the optical sensor, strong and accurate prior belief was constructed. Since the prior knowledge is updated by prediction with the random walk model and correction with the joint observation

likelihood, the proposed approach can eliminate wrong detections and estimate the target position near the true position.

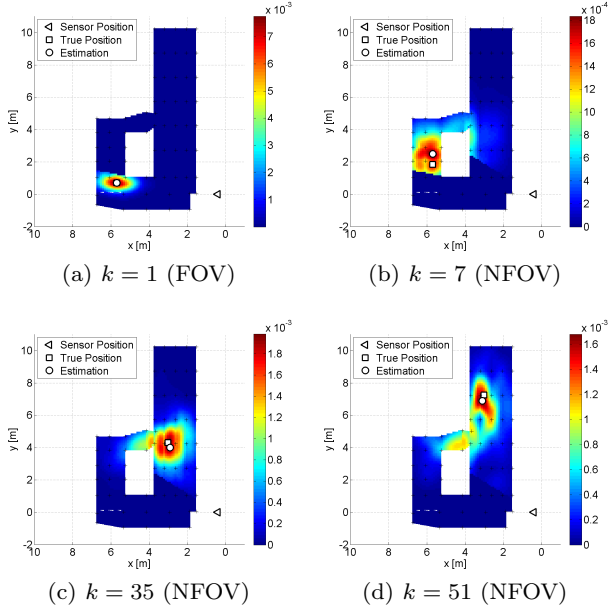


Fig. 16 Proposed optical/acoustic target estimation

Figure 17 shows the target belief estimated conventionally with the optical observation likelihoods only to comparatively verify the effectiveness of the proposed approach. The result with the optical observation likelihoods is seen to estimate the target position wrongly when the target is in the NFOV. Because an inaccurate random-walk motion model is used, the target is estimated continuously at the location where it was lost. Finally, Figure 18 shows the results quantitatively evaluating the performance of the proposed approach. Figure 18(a) shows the transition of the mean error of the estimated target position from the true position whereas the transition of the Kullback-Leibler (KL) divergence is exhibited in Figure 18(b). The error transition indicates that the proposed approach maintains error within 1 m even when the target has not been lost from the FOV for some time whilst the conventional RBE with optical observation likelihoods increases the error with high gradient. The KL divergence transition also shows this behavior, indicating that the proposed approach maintains target information with the use of the acoustic sensor.

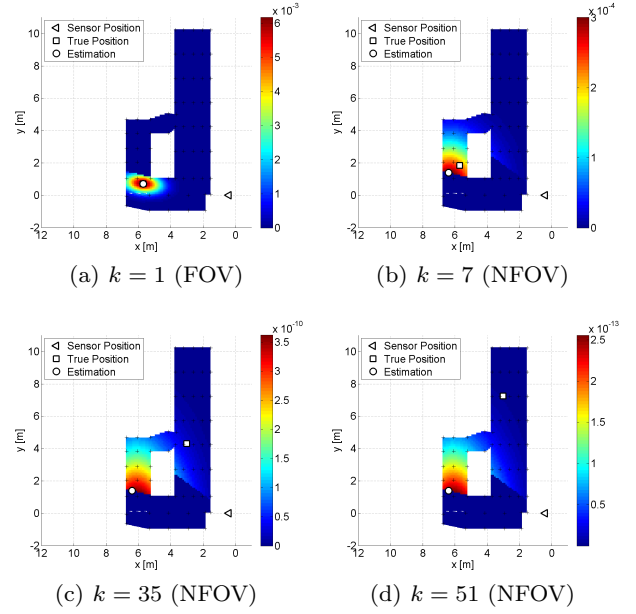


Fig. 17 Conventional optical target estimation

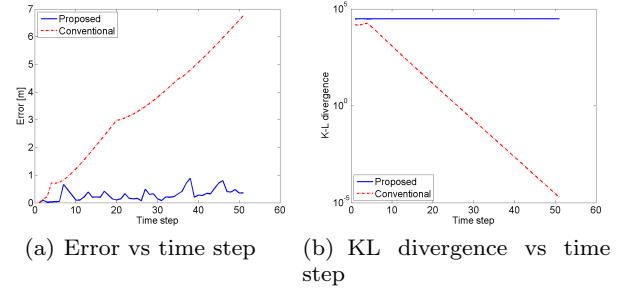


Fig. 18 Quantitative analysis

5 Conclusions

This paper has presented an NFOV target estimation approach which incorporates optical and acoustic sensors. The proposed approach performs RBE with joint optical/acoustic observation likelihoods. Although the acoustic observation likelihood could be multi-modal with high uncertainty, the target belief updated in the past with sharply unimodal optical likelihoods effectively acts as strong prior knowledge and enables accurate NFOV target estimation. A technique to construct an observation likelihood using an acoustic sensor composed of two microphones has also been proposed. The acoustic observation likelihood is created by correlating the ILD of the new acoustic observation with the ILDs collected in advance.

The first experiment studied the capability and limitation of the proposed acoustic sensing technique by parametrically changing the complexity of the environ-

ment. It has been found that the proposed technique can identify the location of an NLOS target but with limited accuracy particularly when the complexity of the environment is severe. The applicability of the proposed RBE approach with joint optical/acoustic observation likelihoods to the estimation of an NFOV target in a complex practical environment was investigated as the second experiment. The result shows that the target position is well estimated at all the time steps even when the joint observation likelihood does not identify the target location well due to the use of prediction with a motion model and strong prior belief constructed with the past optical observation.

The paper has demonstrated the new concept for NFOV target estimation, and many challenges are still open for future study. One of the improvements to the approach is the enhancement of acoustic sensing by incorporating the interaural time difference (ITD) and the interaural phase difference (IPD) as well as the use of non-white noise sound so that the approach could be used for various applications. Other ongoing work includes the implementations of the proposed approach to the mobile sensor platform and to the infrastructure. The proposed approach can be used for various practical applications including home security, home health care, and urban search-and-rescue by the extension.

References

1. P. Bahl and V. N. Padmanabhan. Radar: An in-building rf-based user location and tracking system. In *INFOCOM 2000. Nineteenth Annual Joint Conference of the IEEE Computer and Communications Societies. Proceedings. IEEE*, volume 2, pages 775–784. Ieee, 2000. 1, 2
2. M. Bertinato, G. Ortolan, F. Maran, R. Marcon, A. Marcassa, F. Zanella, M. Zambotto, L. Schenato, and A. Cenedese. Rf localization and tracking of mobile nodes in wireless sensors networks: Architectures, algorithms and experiments. 2008. 2
3. Y. Chan, W. Tsui, H. So, and P. Ching. Time-of-arrival based localization under nlos conditions. *Vehicular Technology, IEEE Transactions on*, 55(1):17–24, 2006. 2
4. P. Chen. A non-line-of-sight error mitigation algorithm in location estimation. In *Wireless Communications and Networking Conference, 1999. WCNC. 1999 IEEE*, pages 316–320. IEEE, 1999. 2
5. H. Dai, Z. Zhu, and X. Gu. Multi-target indoor localization and tracking on video monitoring system in a wireless sensor network. *Journal of Network and Computer Applications*, 2012. 2
6. J. Even, Y. Morales, N. Kallakuri, C. Ishi, and N. Hagita. Audio ray tracing for position estimation of entities in blind regions. In *Intelligent Robots and Systems (IROS 2014), 2014 IEEE/RSJ International Conference on*, pages 1920–1925. IEEE, 2014. 2
7. T. Furukawa, F. Bourgault, B. Lavis, and H. Durrant-Whyte. Recursive bayesian search-and-tracking using coordinated uavs for lost targets. In *Robotics and Automation, 2006. ICRA 2006. Proceedings 2006 IEEE International Conference on*, pages 2521–2526. IEEE, 2006. 2
8. T. Furukawa, L. C. Mak, H. DurrantWhyte, and R. Madhavan. Autonomous bayesian search and tracking, and its experimental validation. *Advanced Robotics*, 26(5-6):461–485, 2012. 2
9. P. Gao, W. Shi, W. Zhou, H. Li, and X. Wang. A location predicting method for indoor mobile target localization in wireless sensor networks. *International Journal of Distributed Sensor Networks*, 2013, 2013. 2
10. S. Gezici. A survey on wireless position estimation. *Wireless Personal Communications*, 44(3):263–282, 2008. 2
11. I. Guvenc and C. Chong. A survey on toa based wireless localization and nlos mitigation techniques. *Communications Surveys & Tutorials, IEEE*, 11(3):107–124, 2009. 2
12. N. D. Jankovic and M. D. Naish. Developing a modular active spherical vision system. In *Robotics and Automation, 2005. ICRA 2005. Proceedings of the 2005 IEEE International Conference on*, pages 1234–1239. IEEE, 2005. 1
13. J. Jung and H. Myung. Indoor localization using particle filter and map-based nlos ranging model. In *Robotics and Automation (ICRA), 2011 IEEE International Conference on*, pages 5185–5190, 2011. 2
14. H. M. Khoury and V. R. Kamat. Evaluation of position tracking technologies for user localization in indoor construction environments. *Automation in Construction*, 18(4):444–457, 2009. 1, 2
15. D. Kimoto and M. Kumon. On sound direction estimation by binaural auditory robots with pinnae. 35th Meeting of Special Interest Group on AI Challenges, 2011. 6
16. D. Kimoto and M. Kumon. Optimization of the ear canal position for sound localization using interaural level difference. 36th Meeting of Special Interest Group on AI Challenges, 2012. 6
17. M. Kobilarov, G. Sukhatme, J. Hyams, and P. Batavia. People tracking and following with mobile robot using an omnidirectional camera and a laser. In *Robotics and Automation, 2006. ICRA 2006. Proceedings 2006 IEEE International Conference on*, pages 557–562. IEEE, 2006. 1
18. A. M. Ladd, K. E. Bekris, A. P. Rudys, D. S. Wallach, and L. E. Kavradi. On the feasibility of using wireless ethernet for indoor localization. *IEEE Transactions on Robotics and Automation*, 20(3):555–559, 2004. 2
19. L. Ledwich and S. Williams. Reduced sift features for image retrieval and indoor localisation. Citeseer. 1
20. H. Liu, H. Darabi, P. Banerjee, and J. Liu. Survey of wireless indoor positioning techniques and systems. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 37(6):1067–1080, 2007. 2
21. Y. Lu and M. Cooke. Binaural estimation of sound source distance via the direct-to-reverberant energy ratio for static and moving sources. *Audio, Speech, and Language Processing, IEEE Transactions on*, 18(7):1793–1805, 2010. 2
22. L. C. Mak and T. Furukawa. Non-line-of-sight localization of a controlled sound source. In *Advanced Intelligent Mechatronics, 2009. AIM 2009. IEEE/ASME International Conference on*, pages 475–480. IEEE, 2009. 2
23. R. Mauler. *Recent Developments in Cooperative Control and Optimizatio*, chapter Objective Functions for Bayesian Control-Theoretic Sensor Management, II:

- MHC-Like Approximation, pages 273–316. Kluwer Academic Publishers, Norwell, MA, 2003. 2
24. K. Nakadai, H. Nakajima, M. Murase, S. Kaijiri, K. Yamada, T. Nakamura, Y. Hasegawa, H. G. Okuno, and H. Tsujino. Robust tracking of multiple sound sources by spatial integration of room and robot microphone arrays. In *Acoustics, Speech and Signal Processing, 2006. ICASSP 2006 Proceedings. 2006 IEEE International Conference on*, volume 4, pages IV–IV. IEEE, 2006. 2
 25. G. Narang, K. Nakamura, and K. Nakadai. Auditory-aware navigation for mobile robots based on reflection-robust sound source localization and visual slam. In *Systems, Man and Cybernetics (SMC), 2014 IEEE International Conference on*, pages 4021–4026. IEEE, 2014. 2
 26. S. K. Nayar. Catadioptric omnidirectional camera. In *Computer Vision and Pattern Recognition, 1997. Proceedings., 1997 IEEE Computer Society Conference on*, pages 482–488. IEEE, 1997. 1
 27. L. M. Ni, Y. Liu, Y. C. Lau, and A. P. Patil. Landmarc: indoor location sensing using active rfid. *Wireless networks*, 10(6):701–710, 2004. 2
 28. Y. Noda and M. Kumon. Sound source direction estimation in the median plane by two active pinnae. 13th SICE System Integration Division Annual Conference, 2012. 6
 29. E. A. Prigge. *A positioning system with no line-of-sight restrictions for cluttered environments*. PhD thesis, Stanford University, 2004. 2
 30. J. Riba and A. Urruela. A non-line-of-sight mitigation technique based on ml-detection. In *Proceedings of IEEE International Conference on Acoustic, Speech and Signal Processing*, volume 2, pages 153–156, May 2005. 2
 31. Y. Sasaki, S. Kagami, and H. Mizoguchi. Online short-term multiple sound source mapping for a mobile robot by robust motion triangulation. *Advanced Robotics*, 23(1-2):145–164, 2009. 2
 32. T. W. Sederberg, J. Zheng, A. Bakenov, and A. Nasri. T-splines and t-nurccs. In *ACM transactions on graphics (TOG)*, volume 22, pages 477–484. ACM, 2003. 7
 33. C. K. Seow and S. Y. Tan. Non-line-of-sight localization in multipath environments. *Mobile Computing, IEEE Transactions on*, 7(5):647–660, 2008. 2
 34. P. Svaizer, A. Brutti, and M. Omologo. Environment aware estimation of the orientation of acoustic sources using a line array. In *Signal Processing Conference (EU-SIPCO), 2012 Proceedings of the 20th European*, pages 1024–1028. IEEE, 2012. 2
 35. J. Valin, F. Michaud, J. Rouat, and D. Létourneau. Robust sound source localization using a microphone array on a mobile robot. In *Intelligent Robots and Systems, 2003.(IROS 2003). Proceedings. 2003 IEEE/RSJ International Conference on*, volume 2, pages 1228–1233. IEEE, 2003. 2
 36. J. Wang, Q. Gao, Y. Yu, H. Wang, and M. Jin. Toward robust indoor localization based on bayesian filter using chirp-spread-spectrum ranging. *Industrial Electronics, IEEE Transactions on*, 59(3):1622–1629, 2012. 2
 37. D. Zhang, Y. Yang, D. Cheng, S. Liu, and L. M. Ni. Cocktail: an rf-based hybrid approach for indoor localization. In *Communications (ICC), 2010 IEEE International Conference on*, pages 1–5. IEEE, 2010. 2