

Person tracking using audio and depth cues

December 11, 2015

Abstract

In this paper, a novel probabilistic Bayesian tracking scheme is proposed and applied to bimodal measurements consisting of tracking results from the depth sensor and audio recordings collected using binaural microphones. We use random finite sets to cope with varying number of tracking targets. A measurement-driven birth process is integrated to quickly localize any emerging person. A new bimodal fusion method that prioritizes the most confident modality is employed. The approach was tested on real room recordings and experimental results show that the proposed combination of audio and depth outperforms individual modalities, particularly when there are multiple people talking simultaneously and when occlusions are frequent.

1 Introduction

Person tracking has been extensively studied in the field of computer vision, with various applications ranging from surveillance, video retrieval, remote meetings to computer-human interactive activities such as video games. Person tracking can be applied to different modalities, e.g. RGB images [3, 10, 7], acoustic recordings [23, 13, 14, 5], depth sensors [12, 21, 16], GPS and thermal sensors. There is a consensus that different modalities are complementary to each other, which has motivated an increasing interest in cross-modal tracking in the last decade. Most of these works are done in the audio-visual domain [24, 8, 9]. Combination of other modalities has recently started to become more popular. For instance, [15, 22] tracks person from both laser range and camera data; the work in [17] fuses RGB, depth and thermal features. Yet, there are some essential limitations associated

with the existing mono- or cross-modal person tracking methods. The mono-modal tracking is not robust enough, while the cross-modal methods often require a high hardware load.

To address the above limitations, we implemented a bimodal person tracking algorithm that combines depth and audio cues. A time-of-flight depth sensor, i.e. Kinect for Windows v2.0 [16], as well as a pair of binaural microphones, i.e. Cortex Manikin MK2 binaural head and torso simulator, are used for person tracking, which are respectively denoted as Kinect2 and Cortex MK2, as shown on the top right and top left of Fig. 1. Individually, both modalities have issues. Audio measurements from Cortex MK2 suffer from heavy background noise and the non-stationary nature of speech. Moreover, they are ambiguous between front and rear sound sources. On the other hand, depth cues from Kinect2 are affected by occlusions. Exploiting the complementary between these two modalities, more robust tracking are obtained from our proposed method. Based on random finite set (RFS) theory, we proposed a full-probabilistic model for multi-person tracking. Particle filters are implemented based on Bayesian filtering [4, 2].

We have several contributions in our proposed method. Firstly, the fusion of depth and audio measurements is not a widely explored territory, and our work has done the initial attempt in this field. Secondly, the proposed method balances the bimodal difference in their structures and robustness, which evaluates the validity of both streams and prioritizes the most confident modality. Thirdly, a measurement driven birth model is used to quickly localize any emerging person.

The remainder of the paper is organized as follows. Section 2 briefly introduces the RFS particle filters in person tracking. Section 3 presents the overall frame work

of our proposed algorithm, and describes in detail the fusion of depth and audio streams. Experimental results are shown and analyzed in Section 4. Finally, conclusions and insights for future research directions are raised in Section 5.

2 RFS-based particle filters

For single target tracking, the hidden state at time frame k , e.g. the position and velocity of the target, is often represented via a vector \mathbf{x}_k . To generalize this problem to the multi-object tracking problem, the hidden state is a finite-set-valued variable $X_k = \{\mathbf{x}_{k,1}, \dots, \mathbf{x}_{k,N_k}\}$ that contains N_k targets, with each $\mathbf{x}_{k,i}$ being the state vector associated with the i -th target. When $N_k = 0$, $X_k = \emptyset$ denotes no target being detected.

X_k can be estimated from a sequence of measurements $[Z_1, Z_2, \dots, Z_k]$ collected/extracted from the sensors, where $Z_k = \{\mathbf{z}_{k,1}, \dots, \mathbf{z}_{k,M_k}\}$ is also a finite-set-valued variable. Note that M_k does not necessarily equal N_k , and $\mathbf{x}_{k,i}$ is not necessarily associated with $\mathbf{z}_{k,i}$. Some measurements are clutters (false alarms) and some targets may fail to generate any measurement.

Bayesian filtering [4, 2] is often applied in target tracking, which propagates the posterior density over time with a recursive prediction and update process. It exploits the temporal involvement as well as the relationship between the underlying positions and the measurements, i.e. the state-space approach. However, this problem might be intractable if the state-space model does not satisfy certain restrictions. Sequential Monte Carlo (SMC) [4] methods can be devoted to its approximations, resulting in the so-called particle filters or bootstrap filters [1]. In multi-target tracking, random finite set (RFS) approach can be used, which takes into account the association uncertainty as well as spurious measurements. More details on RFS-based particle filters are available in [14].

3 Proposed method

Particle filters are applied to a sequence of measurements for target tracking. These measurements are often features extracted from the sensors, which are related to the underlying target state. In this paper, we exploit the

complementary relationship between the audio and depth streams, which are collected by Cortex MK2 and Kinect2 respectively. A novel bimodal person tracking scheme is proposed, whose main flow is shown in Figure 1.

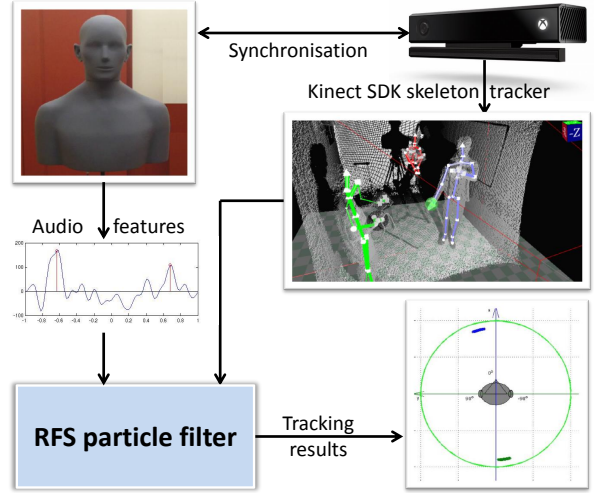


Figure 1: Flow of the proposed audio-depth person tracking method. Synchronized audio and depth measurements are collected from Cortex MK2 and Kinect2 respectively. A RFS particle filter is then employed to these synchronized measurements for person tracking.

We aim to find the relative angle of any person to Cortex MK2 in the horizontal plane, which can be considered as a 1D position or azimuth, as shown in Figure 2. The azimuth direction is not as informative as 3D position, but it is of great importance to attention switching in machine audition (e.g. for source separation) or computer vision (e.g. to handle occlusions). Azimuth estimation can be challenging and in this report we demonstrate the advantages of bimodal tracking over mono-modalities.

Following, we will introduce in detail what these measurements are and how they are fused together.

3.1 Audio-based likelihood function

The time delay of arrival (TDOA) cues are used as audio measurements in our method. The phase transform (PHAT)-GCC [11] method is applied to Cortex MK2 bin-aural recordings. Suppose $L_k(\omega)$ and $R_k(\omega)$ are the short

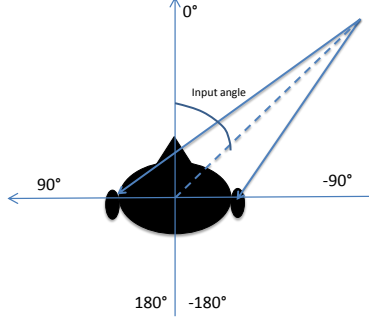


Figure 2: The input angle (azimuth) of a target source in the horizontal plane. A sound source arrives at the two ears via different paths, resulting an inter-aural time difference. The input angle increases from 0° from the nose anti-clockwisely.

time Fourier transforms (STFT) of the two audio segments at time k . The PHAT-GCC function can be applied as:

$$C(\tau) = \int_{-\infty}^{\infty} \frac{L_k(\omega)R_k^*(\omega)}{|L_k(\omega)R_k^*(\omega)|} e^{j\omega\tau} d\omega, \quad (1)$$

where the superscript $*$ denotes the conjugate operator and $|\cdot|$ is a modulus operator. By finding peak positions in PHAT-GCC, M_k^a TDOAs $Z_k^a = \{\tau_{k,1}, \dots, \tau_{k,M_k^a}\}$ can be obtained as the audio measurements¹.

Different positions (azimuths) yield different TDOAs. We need to model the relationship between the audio measurement with the azimuth, i.e. the audio likelihood function, which is complex due to reflections and diffraction of the head. From off-line training, we notice there exists a nonlinear relationship between the resultant TDOA τ with the azimuth α , as shown in Figure 3. Firstly, the curve is symmetric through the axis of 90° or -90° . This is quite understandable as TDOAs have some ambiguity between front and back. Secondly, TDOA from the front can be linearly fitted with the input azimuth (from -90° to 90°) using the polynomial fitting:

$$\tau = f(\alpha) = p_1\alpha + p_3\alpha^3, \quad (2)$$

and $p_1 = 2.405 \times 10^{-6}$ and $p_3 = 1.807 \times 10^{-2}$ are obtained in the off-line training process.

¹The superscript a indicates audio. Similarly, the superscript d stands for depth, and ad denotes audio-depth.

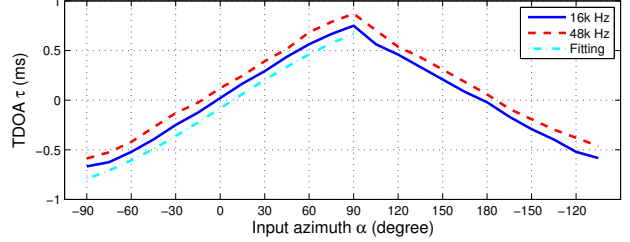


Figure 3: Illustration of the relationship between the resultant TDOA τ with the azimuth α . This was trained from off-line recordings at 16 kHz and 48 kHz. The third-order polynomial curve fitting is applied to model the audio likelihood function. We lifted the curve at 48 kHz, and lowered the fitted curve via a shift of 0.1 either way.

For an azimuth from the back, a mapping function $map(\cdot)$ can be applied to get its mirror reflection:

$$map(\alpha) = \begin{cases} \alpha, & \text{if } |\alpha| \leq 90^\circ, \\ sign(\alpha)(180^\circ - |\alpha|) & \text{otherwise.} \end{cases} \quad (3)$$

Considering zero-mean additive Gaussian noise with variance $\delta^{a,2}$, we can model the audio likelihood as:

$$g(\tau|\alpha) = \mathcal{N}(\tau - f(map(\alpha))|0, \delta^{a,2}), \quad (4)$$

where $\mathcal{N}(\cdot)$ denotes the Gaussian distribution. This noise term also relaxes the non-perfect fitting in Equation (2).

3.2 Depth-based likelihood function

As mentioned before, to get depth measurements, we used the Kinect for Windows v2.0 time of flight sensor, dubbed as Kinect2 in this paper. This sensor emits near infra-red pulses and a fast infrared camera is used to estimate depth based on phase difference. The SDK provided by Microsoft [16] includes a method that detects up to six people and estimates their pose based on a skeleton model with 25 joints. It uses a body part detector based on random decision forests [21]. Each point in a point cloud is classified using simple depth comparison features and a random decision forest (RDF). This RDF is trained with millions of samples of humans, combining real and synthetic samples, at a wide range of poses, with ground

truth labels annotated for each body part (hand, arm, elbow, forearm, etc.). This generates a point cloud where each point is labeled as a body part or as background, when their features do not match a body part. The result is spatially filtered and a post-processing method fits up to six skeletons to the resulting labeled point cloud.

The center right sub-plot in Figure 1 shows detected skeletons in a sample point cloud. Since our goal is to objectify speakers, we are interested in the location of their mouth, which is close to the center of their head. We thus use the position of the heads detected by Kinect2 SDK the 3D position of the sound sources.

A number of methods have been proposed to detect and track people in depth images [25, 18], particularly those generated using sensors based on structured light projection, such as the first version of Kinect. Although the full pipeline implemented in Kinect2 SDK has not been published, we have performed a set of preliminary experiments comparing this method with other state of the art implementations available for 3D head tracking from depth measurements, such as the method of Fanelli et al. [6] and RGB methods, such as that of [20]. Our qualitative observations indicate that the method implemented in Kinect2 SDK robustly achieves state of the art accuracy in head position estimation. Since it has been designed to work on living rooms, the range of distances where it operates is optimal for our application, whereas other implementations available off-the-shelf have been optimized to be used on web-cam scenarios, with a much smaller working distance range.

Despite its robustness, this method has some drawbacks. Since it is based on a *tracking as detection* framework, it does not incorporate a mechanism to handle occlusions based on inference from tracking results. Occlusions cause this implementation to lose measurements or to generate noise outliers and to swap the identity of people being tracked, as shown in Figure 4. It can also fail due to limitations of the sensor itself, such as its working range. If a person is closer than 1.2 meters or further than 3.5 meters, the detector fails. It also fails in cluttered scenes or when people are close to each other.

Further to detecting people, the Kinect2 skeleton detector also locates the dummy automatically as a static person sitting at the center of the room. By having the prior knowledge that the dummy is the audio recording device and that it remains static, we can easily detect it by

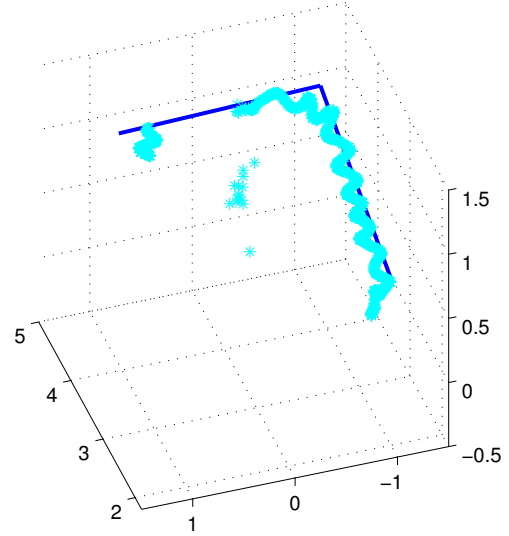


Figure 4: Kinect2 tracking results for a person (Actor B) following an L-shaped trajectory in the room of Figure 5. The target path that this person followed is highlighted in blue lines. The cyan stars show detected positions, which wiggles because the head actually swings from side to side as this person walked. There is a cluster of mis-detected positions, i.e. this subject’s head was detected around the dummy head position when he was occluded by the dummy. Shortly after that, there was a period of consecutive frames where the target is not detected because of this occlusion.

analyzing a sequence of recordings in a pre-processing step. This enables us to label the dummy and distinguish it from moving people. It also enables us to project the 3D position of detected humans to the polar coordinate system centered at the dummy head. Since the head position is estimated in 3D from Kinect2’s viewpoint, there is no front/back ambiguity w.r.t. the dummy head and the mapping of Equation 3 is not necessary for depth-based cues. The obtained azimuth angle measurements are used as depth-based observations, denoted as $Z_k^d = \{\theta_{k,1}, \dots, \theta_{k,M_k^d}\}$, i.e., in the remainder of this paper, we assume that the pipeline that maps from depth images to azimuth angles relative to the dummy head is part of the measurement process.

As mentioned earlier, the head tracker is usually re-

liable, but occlusions introduce severe noise or missing depth-based measurements, which we approximate using the zero-mean additive Gaussian noise with variance, defined as $\delta^{d,2}$. Therefore, the likelihood of the associated Kinect detection given an input angle follows

$$g(\theta|\alpha) = \mathcal{N}(\theta - \alpha|0, \delta^{d,2}). \quad (5)$$

3.3 Audio-depth fusion

As introduced in Section 2, RFS particle filters are applied to the audio and depth measurements. Their state space model contains two essential parts: the dynamic model and the measurement model.

3.3.1 Dynamic model

The dynamic model describes the temporal evolution of target states. For multi-targets, each state vector $\mathbf{x}_k \in \mathcal{X}_k$ at frame k can either survive with probability P_s or die with probability $1 - P_s$ at the next frame. Let \mathbf{x}_k contain the input angle α and the angular velocity $\dot{\alpha}$; the Largeven model can be utilized to model the relationship between a survived target \mathbf{x}_{k+1} and its previous state:

$$\mathbf{x}_{k+1} = \begin{bmatrix} 1 & T \\ 0 & e^{-\beta T} \end{bmatrix} \mathbf{x}_k + \begin{bmatrix} 0 \\ \nu \sqrt{1 - e^{-2\beta T}} \mathcal{N}(\cdot|0, 1) \end{bmatrix}, \quad (6)$$

where T is the time duration between two consecutive frames; β and ν parametrize the motion model.

Moreover, a new target may be born in the searching field with probability P_b . To quickly localize any appearing target, we propose a measurements-driven target birth model as follows.

The current measurements \mathcal{Z}_k can be mapped to a group of azimuths. We assume the birth model as a mixture of Gaussian kernels, whose mean and standard deviation are these mapped azimuths and 0.1 m. The velocity of newborn targets is zero. Following that distribution, newborn targets are enforced to those potential positions yielding the current measurements. The proposed method can therefore quickly localize any emerging target. Similar idea of adaptive target birth intensity is used in [19].

3.3.2 Measurement model

The measurement or observation model describes the relationship between the target state and the measurement. From sections 3.1 and 3.2, we know the relationship between a single observed mono-modal feature and its associated single-target state. However, for cross-modal multi-person tracking, we need $g(Z_k|X_k)$, where both the bimodal feature Z_k and the multiple-target state X_k are set-valued variables. From empirical study, we notice the azimuth estimates based on depth data alone has far fewer outlier as compared to the audio stream. As a result, a depth-dominant fusion scheme is proposed.

We assume there are up-to-two people in the searching field. As a result, the hidden target state X_k can either be \emptyset , $\{\mathbf{x}_{k,1}\}$ or $\{\mathbf{x}_{k,1}, \mathbf{x}_{k,2}\}$.

When there is no detected target, i.e. $X_k = \emptyset$,

$$g(Z_k|\emptyset) = \left(\frac{P_c^a}{2\tau_{max}}\right)^{|Z_k^a|_0} \left(\frac{P_c^d}{360}\right)^{|Z_k^d|_0}, \quad (7)$$

where P_c^a and P_c^d are the expected number of false alarms at each frame, and $|\cdot|_0$ counts the number.

When there is one detected target, i.e. $X_k = \{\mathbf{x}_{k,1}\}$,

$$g(Z_k|\{\mathbf{x}_{k,1}\}) = p(Z_k|\emptyset)((1 - P_d) + P_d g^{ad}(\mathbf{x}_{k,1})), \quad (8)$$

where $g^{ad}(\mathbf{x}_{k,1}) = \max(g^a(\mathbf{x}_{k,1}), g^d(\mathbf{x}_{k,1}))$ with $g^a(\mathbf{x}_{k,1}) = \max_{\mathbf{z} \in Z_k^a} \frac{g(\mathbf{z}|\mathbf{x}_{k,1})^{2\tau_{max}}}{P_c^a}$ using Equation (4), and $g^d(\mathbf{x}_{k,1})$ similarly using Equation (5). P_d is the chance a target being detected.

When there are two detected targets,

$$\begin{aligned} g(Z_k|\{\mathbf{x}_{k,1}, \mathbf{x}_{k,2}\}) &= p(Z_k|\emptyset)((1 - P_d)^2 \\ &\quad + P_d(1 - P_d)g^{ad}(\mathbf{x}_{k,1}) \\ &\quad + P_d(1 - P_d)g^{ad}(\mathbf{x}_{k,2}) \\ &\quad + P_d^2 g^{ad}(\mathbf{x}_{k,1})g^{ad}(\mathbf{x}_{k,2})). \end{aligned} \quad (9)$$

Computational complexity of the above full-probabilistic model becomes much bigger with an increasing number of targets. As a result, we constrain our algorithm to up-to-two person. To prioritize the depth stream, we make P_c^a larger than P_c^d . We also evaluate the validity of the audio stream in each frame via a straightforward energy

thresholding. If the audio frame is invalid, i.e., speech energy is very low, we then degenerate the proposed model to depth-only mode. This is implemented by keeping only the depth term in Equations (7-9).

4 Experiments

4.1 Recording setup

Our testbed is a TV/film studio set built following professional media production standards, with furniture and features of a relatively large hallway whose dimensions are very similar to those of a typical living room: $244 \times 396 \times 242$ cm. As with typical TV/film production sets, its ceiling and one of the walls are missing, though this set was assembled inside a larger room. The reverberation time of this room is about 430 ms. In the recordings for our experiments, the binaural microphone (Cortex MK2) stood in the center of the room with ear height of 165 cm. The depth sensor was placed around the center at the height of 170 cm, 329 cm away from the depth sensor, as shown in Figure 5. The sampling rate for audio is $F_s^a = 44.1$ kHz. The depth-based head tracker has a sampling rate of $F_s^d = 27.43$ Hz. We used hand clapping at the beginning and end of each recording session to synchronize these two streams. The hand claps can be detected from the audio stream via energy thresholding, and arm pose detection using skeletal tracker from the depth stream.

Three sequences were recorded in total about 7.5 minutes, involving two actors: Actor A is a male, with height of 1.82 m and Actor B is a female, 1.58 m. In the first sequence, Actor A started at Position 1 (labeled with a yellow circle in Figure 5), facing the center, walking slowly along the gray circular trajectory anti-clock-wisely, reading randomly-selected sentences from the TIMIT database. He walked back clock-wisely along the gray circle when reaching Position 24. Actress B repeated this process with a higher speed, and this was recorded in the second sequence. In the third sequence, Actor A started at Position A, walking along the path $A \rightarrow B \rightarrow C \rightarrow B \rightarrow A$, facing forward. At the same time Actress B started at Position C, walking along the path $C \rightarrow D \rightarrow A \rightarrow D \rightarrow C$, facing forward. Therefore, both actors followed L-shaped paths (symmetric to

each other, relative to the room), moving independently from each other, each walking at his/her preferred pace while reading the material mentioned earlier.

4.2 Implementation details

To obtain audio measurements, 8192-point (approximately 186 ms) Hamming windowed STFT with 0.75 overlap is applied. The time length between two neighboring frames is therefore $T = 139$ ms. The candidate τ is linearly sampled in the range of -1 ms to 1 ms ($\tau_{max} = 1$ ms) with the resolution of $1/F_s^a$. At each time frame, at most two TDOAs are extracted as audio measurements.

Figure 6 shows the extracted audio measurements from Sequences 1 and 3. Sequence 1 has only one speaker facing the binaural microphone, while Sequence 3 has two speakers and they do not face the microphone most of the time. In addition, Sequence 3 contains heavy background noise.

To implement RFS-particle filters, the following parameters are used. The target survive chance $P_s = 0.99$, and a target birth $P_b = 0.02$. Parameters in the Largeiven model are set as $\beta = 10, \nu = 10$. Chance a target being detected is $P_d = 0.75$, and the false alarm expectations are $P_c^a = 0.5$ and $P_c^d = 0.1$. The mono-modal likelihood functions in Equations 4 and 5 have the standard variance of $\delta^a = 1/16$ ms and $\delta^d = 5^\circ$.

4.3 Results and analysis

4.3.1 Single person, audio-only features

Firstly, we tested the proposed algorithm on the first two sequences, using only audio features. In Sequences 1 and 2, only less than 10 seconds occlusions is observed. In addition, when there is no occlusion, very accurate depth-based tracking results are obtained except for only a few frames of outlier. As a result, we manually corrected these outlier and labeled the misdetected frames from the depth images when occlusions happened to obtain the ground-truth, which was down-sampled to be synchronized with the audio measurements on a frame basis.

Note that, the TDOA audio cues cannot distinguish a signal from front or back. For instance, the signal from 45° and 135° yields the same TDOA features. To address

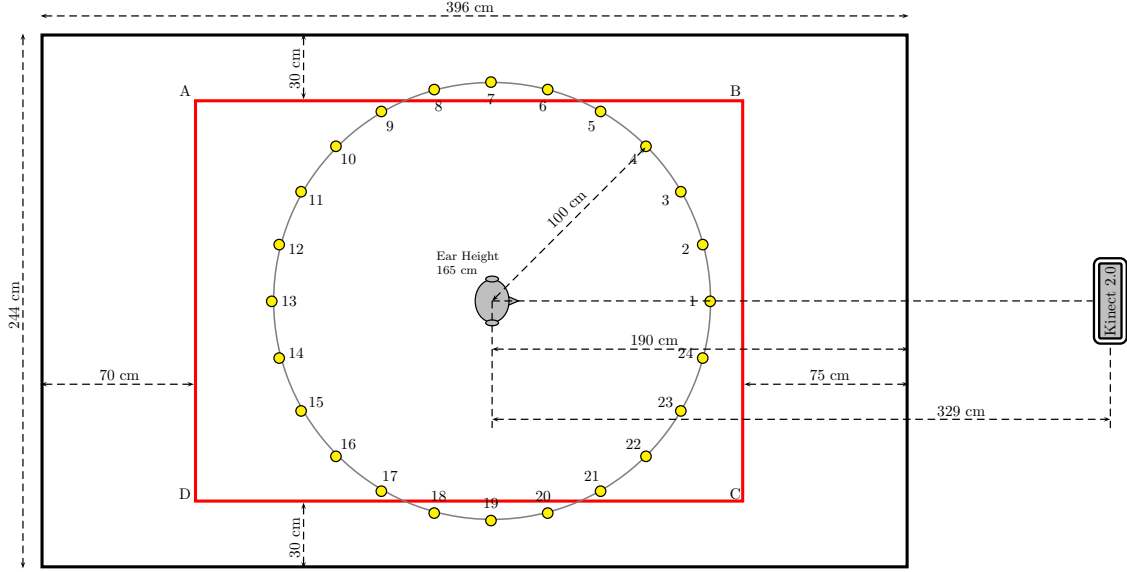


Figure 5: Setup for data recordings. The 24 highlighted dots in a circle labels the positions used to model the relationship between different input angles and the exhibited audio features.

this ambiguity, an audio range assumption of $[-90^\circ, 90^\circ]$ was imposed. In other words, we assumed the signal comes in front of the dummy head. The tracking results for Sequence 1 is shown in Figure 7.

We then quantitatively evaluated the proposed method for single speakers using only audio cues. In Sequence 1, the ground truth have 1282 frames in total, and the proposed method results show 1273 frames have one speaker, and 9 frames have no speaker, caused by long silent periods. The standard deviation for the error is 7.8° . In Sequence 2, the ground truth have 1282 frames in total, and the proposed method results show 599 frames have one speaker, and 26 frames have no speaker. Of the 599 frames, 109 frames has the error more than 30° , in the beginning and the end of the recording session, when the person is not silent, and an interfering speaker outside the recording field is talking. The standard deviation for the error of the remaining 480 frames is 10.6° .

4.3.2 Single person, audio and depth features

Secondly, we tested the proposed algorithm on the first two sequences, using both audio and depth features. The

tracking results for Sequence 1 is shown in Figure 8, from which very good tracking results were observed. The detected trajectory is almost overlapped with the ground-truth.

We then did some quantitative evaluations. In Sequence 1, in all of the 1282 frames, one person was successfully detected, with the deviation of 2.4° . In Sequence 2, in all but 2 frames one person was detected, with the deviation of 3.8° . Compared with the results using audio-only features, the combination of audio and depth greatly reduced the error.

4.3.3 Two people, audio and depth features

Finally, we tested our algorithm on the two people scenario. Using audio-only features, the proposed method did not converge since the audio measurements are too noisy. Using depth-only features, the outliers were removed. However, the occluded person was not tracked. Using both audio and depth features, we successfully tracked both speakers, as shown for Sequence 3 in Figure 9. However, note that the identity of the speakers got swapped occasionally. This problem can be solved by ap-

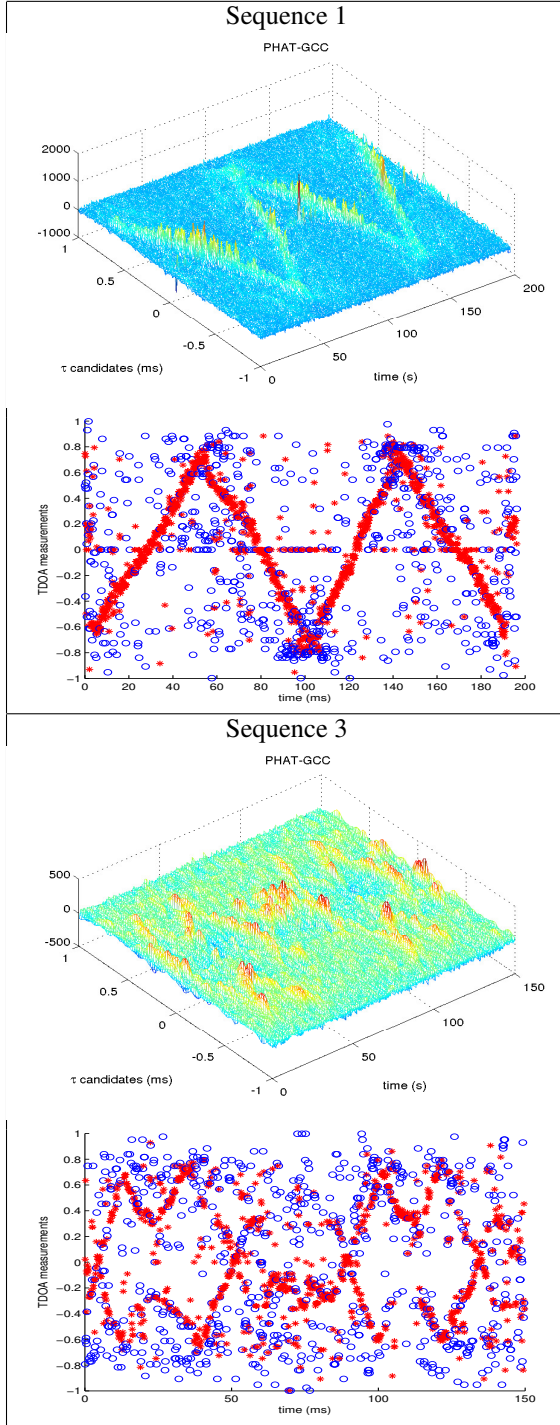


Figure 6: PHAT-GCC results and detected TDOAs from Sequence 1 and Sequence 3. The red-starred points denote the first peak-related TDOA while the blue circles represent the second one. The peak in Sequence 1 is very smooth, which clearly exhibits the speaker’s trajectory. However, despite some peaks related to the real positions in Sequence 3, much more false alarms are obtained.

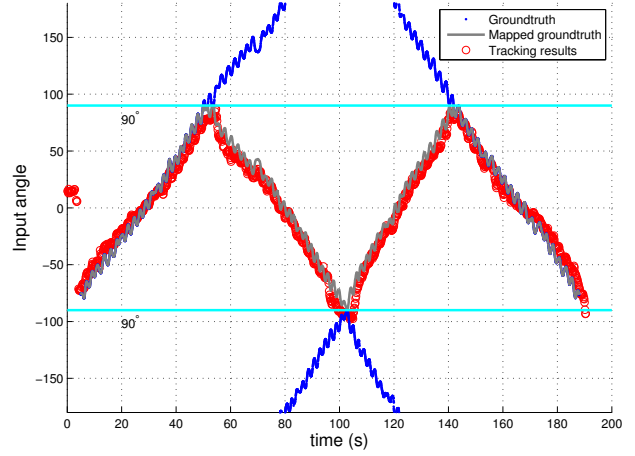


Figure 7: The proposed method applied to Sequence 1 using only the audio features. Since the audio features have the front and back confusion, we imposed the input angle range of $[-90, 90]$. The blue dots represent the ground-truth input angle. We symmetrically mapped the angles at the back of the dummy head, i.e. $[-180, -90)$ and $(90, 180]$, to the front. The mapped ground-truth is the gray curve. The tracking results are represented via the red circles. Comparison between the tracked results and the mapped ground-truth demonstrates the practicability of the proposed method, and the previously-set parameters are well defined.

plying a simple filter in space-time, e.g. by calculating the distance between detected person in two consecutive frames. Our depth-audio results on Sequence 2 were also consistent with the walking trajectory described earlier, demonstrating success with the fusion of depth and audio cues.

5 Conclusions

We presented a method for multimodal tracking using audio and depth features. TDOA features are extracted from binaural recording (Cortex MK2); 3D positions from the depth sensor (Kinect2) are mapped into 1D azimuth relative to Cortex MK2 as the depth features. The measurements from both modalities were fused in a particle filtering framework that enables birth and death of multiple

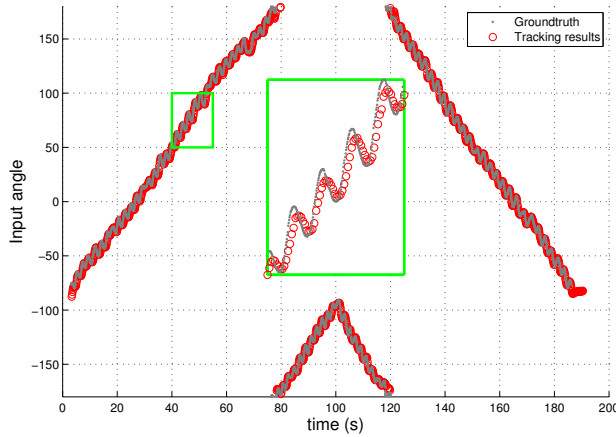


Figure 8: Application of the proposed method to Sequence 1 using both the audio and depth features. The dots represent the ground-truth input angle. The tracking results are represented via the circles. We noticed the tracking results almost overlapped with the ground-truth. We have zoomed in a short segment highlighted in the rectangle.

tracks using Random Finite Sets (RFS). These two modalities are obviously very different and have very different levels of confidence. We showed how to take that into account and how they can complement each other. Our results show that this combination clearly outperforms individual modalities, particularly when there are multiple people talking simultaneously and when there is a significant amount of occlusion.

As future work, we plan to perform experiments on more datasets, aiming to highlight the method’s potential to handle birth and death of targets. We also intend to compare our results against other tracking and fusion methods. The RFS tracking framework is a principled way to simultaneously track a varying number of targets, but its complexity grows as the number of targets increase. We suggest that depth-based tracking results, including the detected targets identities, should help us to design a modified version of RFS, with lower complexity w.r.t. the number of targets. We also plan to use the most confident modality to provide strong priors on the birth and death of tracking targets.

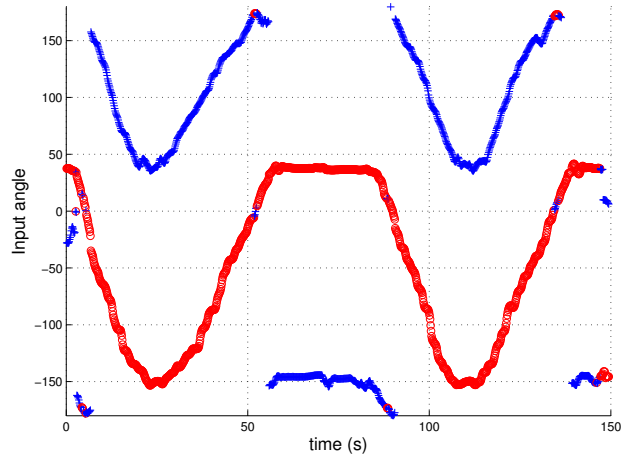


Figure 9: Application of the proposed method to Sequence 3 using both the audio and depth features. The blue dots represent Actor B, and the red dots represent Actress A.

References

- [1] M. Arulampalam, S. Maskell, N. Gordon, and T. Clapp. A tutorial on particle filters for online nonlinear/non-gaussian bayesian tracking. *IEEE Transactions on Signal Processing*, 50(2):174–188, February 2002. 2
- [2] Z. Chen. Bayesian filtering: From Kalman filters to particle filters, and beyond. Technical report, McMaster University, 2003. 1, 2
- [3] D. Comaniciu, V. Ramesh, and P. Meer. Kernel-based object tracking. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 25(5):564–575, May 2003. 1
- [4] A. Doucet, N. de Freitas, and N. Gordon. An introduction to sequential monte carlo methods. In *Sequential Monte Carlo Methods in Practice*, Statistics for Engineering and Information Science, pages 3–14. Springer New York, 2001. 1, 2
- [5] M. Fallon and S. Godsill. Acoustic source localization and tracking of a time-varying number of speakers. *IEEE Transactions on Audio, Speech, and Language Processing*, 20(4):1409–1415, May 2012. 1
- [6] G. Fanelli, M. Dantone, J. Gall, A. Fossati, and L. Van Gool. Random forests for real time 3d face analysis. *Int. J. Comput. Vision*, 101(3):437–458, February 2013. 4
- [7] F. Fleuret, J. Berclaz, R. Lengagne, and P. Fua. Multicamera people tracking with a probabilistic occupancy map.

- IEEE Transactions on Pattern Analysis and Machine Intelligence*, 30(2):267–282, Feb 2008. 1
- [8] D. Gatica-Perez, G. Lathoud, J. Odobez, and I. McCowan. Audiovisual probabilistic tracking of multiple speakers in meetings. *IEEE Transactions on Audio, Speech, and Language Processing*, 15(2):601–616, Feb 2007. 1
 - [9] V. Kilic, M. Barnard, W. Wang, and J. Kittler. Audio assisted robust visual tracking with adaptive particle filtering. *IEEE Transactions on Multimedia*, 17(2):186–200, Feb 2015. 1
 - [10] K. Kim and L. S. Davis. Multi-camera tracking and segmentation of occluded people on ground plane using search-guided particle filtering. In *Proceedings of the 9th European Conference on Computer Vision*, pages 98–109, 2006. 1
 - [11] C. Knapp and G. Carter. The generalized correlation method for estimation of time delay. *IEEE Transactions on Acoustics, Speech and Signal Processing*, 24(4):320–327, Aug 1976. 2
 - [12] S. J. Krotosky and M. M. Trivedi. Mutual information based registration of multimodal stereo videos for person tracking. *Computer Vision and Image Understanding*, 106(23):270–287, 2007. 1
 - [13] E. A. Lehmann and R. C. Williamson. Particle filter design using importance sampling for acoustic source localisation and tracking in reverberant environments. *EURASIP Journal on Advances in Signal Processing*, 2006(1):1–9, 2006. 1
 - [14] 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*, 54(9):3291–3304, 2006. 1, 2
 - [15] M. P. Michalowski and R. Simmons. Multimodal person tracking and attention classification. In *Proceedings of the 1st Annual Conference on Human-Robot Interaction*, pages 347–348. ACM, March 2006. 1
 - [16] Microsoft. Kinect for Windows. Online, Retrieved in September 2015. <https://dev.windows.com/en-us/kinect/>. 1, 3
 - [17] A. Mogelmose, C. Bahnsen, T. Moeslund, A. Clapes, and S. Escalera. Tri-modal person re-identification with RGB, depth and thermal features. In *Proceedings of CVPR Workshops*, pages 301–307, June 2013. 1
 - [18] C. Redondo-Cabrera, R. Lopez-Sastre, and T. Tuytelaars. All together now: Simultaneous detection and continuous pose estimation using a hough forest with probabilistic locally enhanced voting. In *Proc 25th British Machine Vision Conf, Nottingham, Sept 1-5*. BMVA Press, 2014. 4
 - [19] B. Ristic, D. Clark, B.-N. Vo, and B.-T. Vo. Adaptive target birth intensity for phd and cphd filters. *IEEE Transactions on Aerospace and Electronic Systems*, 48(2):1656–1668, 2012. 5
 - [20] J. M. Saragih, S. Lucey, and J. F. Cohn. Face alignment through subspace constrained mean-shifts. In *12th International Conference on Computer Vision*, pages 1034–1041. IEEE, 2009. 4
 - [21] J. Shotton, A. Fitzgibbon, M. Cook, T. Sharp, M. Finocchio, R. Moore, A. Kipman, and A. Blake. Real-time human pose recognition in parts from single depth images. In *Proc of the IEEE Conf on Computer Vision and Pattern Recognition*, 2011. 1, 3
 - [22] L. Spinello, R. Triebel, and R. Siegwart. Multimodal people detection and tracking in crowded scenes. In *Proceedings of the 23rd National Conference on Artificial Intelligence*, pages 1409–1414, 2008. 1
 - [23] J. Vermaak and A. Blake. Nonlinear filtering for speaker tracking in noisy and reverberant environments. In *IEEE International Conference on Acoustics, Speech, and Signal Processing*, volume 5, pages 3021–3024, 2001. 1
 - [24] J. Vermaak, M. Gangnet, A. Blake, and P. Perez. Sequential monte carlo fusion of sound and vision for speaker tracking. In *IEEE International Conference on Computer Vision*, volume 1, pages 741–746. IEEE, 2001. 1
 - [25] M. Ye, Q. Zhang, L. Wang, J. Zhu, R. Yang, and J. Gall. A survey on human motion analysis from depth data. In *Time-of-Flight and Depth Imaging. Sensors, Algorithms, and Applications*, volume 8200 of LNCS, pages 149–187. Springer, 2013. 4