

# Online Multi-Object Tracking Framework with the GMPHD Filter and Occlusion Group Management

Young-min Song, Kwangjin Yoon, Young-Chul Yoon, Kin-Choong Yow, and Moongu Jeon\*

**Abstract**—In this paper, we propose an efficient online multi-object tracking framework based on the Gaussian mixture probability hypothesis density (GMPHD) filter and occlusion group management scheme where the GMPHD filter utilizes hierarchical data association to reduce the false negatives caused by miss detection. The hierarchical data association consists of two steps: detection-to-track and track-to-track associations, which can recover the lost tracks and their switched IDs. In addition, the proposed framework is equipped with an object grouping management scheme which handles occlusion problems with two main parts. The first part is “track merging” which can merge the false positive tracks caused by false positive detections from occlusions, where the false positive tracks are usually occluded with a measure. The measure is the occlusion ratio between visual objects, sum-of-intersection-over-area (SIOA) we defined instead of intersection-over-union (IOU) metric. The second part is “occlusion group energy minimization (OGEM)”, which prevents the occluded true positive tracks from false “track merging”. We define each group of the occluded objects as an energy function and find an optimal hypothesis which makes the energy minimal. We evaluate the proposed tracker in benchmark datasets such as MOT15 and MOT17 which are built for multi-person tracking. An ablation study in training dataset shows that not only “track merging” and “OGEM” complement each other but also the proposed tracking method has more robust performance and less sensitive to parameters than baseline methods. Also, SIOA works better than IOU for various sizes of false positives. Experimental results show that the proposed tracker efficiently handles occlusion situations and achieves competitive performance compared to the state-of-the-art methods. Especially, our method shows the best multi-object tracking accuracy among the online and real-time executable methods.

**Index Terms**—multiple object tracking, GMPHD filter, hierarchical data association, occlusion handling, energy minimization

## I. INTRODUCTION

MULTI-OBJECT Tracking (MOT) has become one of key techniques for intelligent video surveillance [5], [6] and autonomous vehicle systems [7] in the last decade.

In view of the processing pipelines, many state-of-the-art MOT methods [13]–[50] have exploited the tracking-by-detection paradigm. This phenomenon has been standing out

Y. Song, K. Yoon, and M. Jeon are with the School of Electrical Engineering and Computer Science, Gwangju Institute of Science and Technology, Gwangju 61005, South Korea (e-mail: sym@gist.ac.kr; yoon28@gist.ac.kr; mjeon@gist.ac.kr).

Y.-C. Yoon is with Convergence Biz. Development Office, LG Electronics, Seoul 07796, South Korea (email : youngchul.yoon@lge.com).

K.-C. Yow is with the Faculty of Engineering and Applied Science, University of Regina, Regina S4S 0A2, Canada (e-mail: Kin-Choong.Yow@uregina.ca).

while deep neural networks based detectors such as FRCNN [10], SDP [11], and EB [12] have shown breakthrough in object classification and detection.

Besides, MOT algorithms are categorized into two approaches: offline and online processes. The most different point between two approaches is that whereas the offline process can see the whole time sequences at once, the online process can see only the frames from initial time 1 to current processing time  $k$ . In other words, from the system user’s perspective, whereas the offline method is suitable for post-processing, the online process for real-time application.

Thus, many offline methods [28], [29], [31], [34], [36], [46] take advantage of the global optimization models. [29], [36], [46] exploit graphical models to solve MOT task. Pirsiavash *et al.* [36] designed a min-cost flow network where the nodes and the directed edges indicating observations and tracklets’ hypotheses, respectively form a directed acyclic graph (DAG). The DAG’s shortest (min-cost) path can be found with Dijkstra’s algorithm. Choi *et al.* [29] divided the tracking problem into subgraphs and solved each subgraph as conditional random field inference in parallel. Keuper *et al.* [46] applied vision-based perspective to the proposed graph optimization model. Feature points’ trajectories and bounding boxes build low-level and high-level graph models, respectively, and then, they find the optimal association results between the two levels graph models. Rezatofighi *et al.* [31] and Kim *et al.* [28] considered all possible hypotheses for data association. Because it involves the exponentially increasing complexity with a tree structure, [31] assumed  $m$ -best solutions and [28] pruned out invalid hypotheses using their own rule. Besides, Milan *et al.* [34] proposed a sophisticated energy minimization technique considering detection, appearance, dynamic model, mutual exclusion, and target persistence for MOT task in video. Those offline methods have strength to generate the accurate and refined tracking results but is not suitable for practical real-time application.

On the other hand, since the online approach cannot apply the global optimization models, intensive motion analysis and appearance feature learning have been popularly utilized with a hierarchical data association framework and the online Bayesian model [14], [15], [18], [21], [23], [25], [37], [42]. Yoon *et al.* [23] proposed a relative motion analysis between all objects in a frame, and then improved the work [23] by adding the cost optimization function using context constraints in [21]. Bae *et al.* [25] exploited the incremental linear discriminant analysis (LDA) for appearance learning and presented a tracklet confidence based data association framework. Also, in [14], they improved their previous work [25]

by using the deep neural network (DNN) based appearance learning instead of the incremental LDA. As we addressed in the previous paragraph, DNN has given breakthrough in appearance learning i.e., object classification and detection. So, some online MOT algorithms have focused on how to adopt deep appearance learning into their tracking frameworks. Yoon *et al.* [15] exploited the siamese convolutional neural networks (CNN) [51] to train appearance model. They train the deep appearance networks selectively where only the detection responses matched with high confidence between the historical object queues in the recent few frames. Then, they combine the trained networks to a simple Bayesian tracking model with the Kalman filter. Chen *et al.* [37] employed a re-identification (Re-ID) model [52] to their tracking framework. They measure the similarity between detection and track by calculating the distance between Re-ID feature vectors of them. Then, they associate the pairs of detections and tracks which make the sum of the distances minimal. Both approaches [15], [37] proposed online Bayesian tracking models with conventional DNN models to measure the similarity between the visual objects. Those online MOT methods have proposed successful solutions with excellent tracking accuracy but their intensive analysis and learning processes take heavy computing resource and time. Also, even if they just employ conventional DNN models through state-of-the-art GPU processing technique, the requirement for a lot of computing resource is inevitable and it makes the trackers difficult to achieve real-time speed.

Recently, the closed-form implementations [2], [3] of the probability hypothesis density (PHD) filtering have been employed as an emerging theory for many online MOT methods [16]–[18], [22], [38], [41]–[43]. That is because Vo *et al.* [2], [3] provided not only theoretically optimal approach to the online multi-target Bayes filtering but also approximate the original PHD recursions involving multiple integrals, which alleviate the computational intractability. Moreover, the PHD filter was originally designed for multi-target tracking in radar/sonar systems which receive uncountable false positive observations, i.e., clutters. So it is robust to deal with false positives errors but weak to handle false negatives. V. Eiselein *et al.* [43] combined the feature-based label tree to the Gaussian mixture PHD (GMPHD) filter, which use visual features to help the GMPHD filter work sensibly in video data system. Song *et al.* [16] extended the GMPHD filter based tracking with the two-stage hierarchical data association strategy and use simple motion estimation and appearance matching to recover lost tracks. T. Kutschbach *et al.* [42] joined the GMPHD filter with the kernelized correlation filters (KCF) [53] for online appearance update to overcome occlusion. Z. Fu *et al.* [18] adopted an adaptive gating technique and an online group-structured dictionary (appearance) learning strategy into the GMPHD filter. They make the GMPHD filter be sophisticated and fit to video based MOT. Besides, various tracking methods [22], [38], [41] utilizing the PHD filters have been proposed.

These latest MOT research trends motivate our work in terms of the three main contributions. Also, it reminds us of the requirements for the practical MOT applications. Thus, in this paper, we propose an online multi-object tracking framework

to resolve the practical tracking problems which are based on occlusion and the characteristics of video data system. First, we exploit the GMPHD filter for online MOT. To efficiently change the GMPHD filter's original domain, we define the tracking problems by miss detections in video data system. To deal with track loss by miss detection, we design a GMPHD filtering theory based hierarchical data association (HDA) strategy. Second, we assume that most of tracking problems are caused by occlusion in video data system. The occlusion between false positive tracks can cause ID-switch and the false positives, and the real occlusion between objects can make fragmented and miss tracks by miss detections. To handle these tracking problems, we propose a novel occlusion handling technique combined with HDA which is based on GMPHD filter tracking framework. Third, we consider that the proposed tracking framework should be implemented to run with real-time speed. That is because visual surveillance systems with higher intelligence require more immediate responses to the users with real-time speed. Also, immediate responses can help the systems' user and the machines to react abnormal situation rapidly. Finally, we evaluate the proposed method on the popular benchmark dataset. Our method shows the competitive performance against state-of-the-art methods in terms of “tracking accuracy versus speed”.

Our main contributions are described as follows:

1) To apply the GMPHD filter into video data system, we extended the conventional GMPHD filter based tracking process with a hierarchical data association (HDA) strategy. Also, we revised the equations of the GMPHD filter as a new cost function for HDA. HDA consists of detection-to-track association (D2TA) and track-to-track association (T2TA). Each cost matrix of each association stage is solved by the Hungarian method with the linear complexity  $O(n^3)$  (assignment problem). These D2TA and T2TA recovers lost tracks, while preserving real-time speed.

2) To handle occlusion in video-based tracking system, we devised “tracking merging” and “occlusion group energy minimization (OGEM)” which complement each other. “Tracking merging” relieves false positive tracks and “OGEM” recovers false “track merging” by using the occluded objects’ group energy minimization. “Tracking merging” runs in tracking-level so is different to detection-level merging such as non-maximum-suppression. To measure overlapping ratio between occluded objects, we devise a new metric named as sum-of-intersection-over-area (SIOA). We use the SIOA metric instead of intersection-over-union (IOU) which is an extensively used metric. For “OGEM”, we devise a new energy function to find the optimal state having the minimum energy in a group of occluded objects. “Tracking merging” and “OGEM” follow D2TA and T2TA, respectively. We name both techniques as occlusion group management (OGM).

3) Consequently, we propose an online multi-object tracking framework with the GMPHD filter and occlusion group management (GMPHD-OGM). In view of optimization techniques, the first and second contribution locally optimize tracking process which are the minimization of the association cost matrix and the occlusion group energy. We evaluate the proposed tracking framework on MOT15 [5] and MOT17 [6]

benchmarks. The ablation study on training set shows that our method is more robust than the given baselines. The qualitative and quantitative evaluation results shows that GMPHD-OGM efficiently handle the defined tracking problems by occlusion. Moreover, the proposed method achieves competitive tracking performance against state-of-the-art online MOT algorithms in terms of CLEAR-MOT metrics [54].

The related works are described in Section II. In Section III and IV, we introduce the GMPHD filter based tracking framework with HDA and OGM in detail, respectively. In Section V, our method is evaluated compared to baseline methods and state-of-the-art methods on the popular benchmarks MOT15 [5] and MOT17 [6]. We conclude this paper with future work in Section VI. Some preliminary results of this work was presented in Song *et al.* [16], [17].

## II. RELATED WORKS

Our proposed tracking framework is influenced from the PHD filter based online multi-object tracking, and grouping approach (topology and relative motion analysis).

The PHD filter [1]–[3] was originally designed to deal with radar/sonar data based multi-object tracking (MOT) systems. Mahler *et al.* [1] proposed a recursive Bayes filter equations for the PHD filter which optimizes MOT process in radar/sonar systems with the random-finite set (RFS) of state and observations. Following this PHD filtering theory, Vo *et al.* [3] implemented governing equations by using the Gaussian mixture model as closed-form recursions, named as the Gaussian mixture probability hypothesis density (GMPHD) filter. In the original domains the tracking algorithm should estimate true tracks (states) from a lot of observations as shown in Figure 1-(a). Whereas the radar/sonar sensors receive massive false positive but rarely missed observations, visual object detectors generate much less false positive and more missed observations than the radar/sonar sensors does as shown in Figure 1-(b). Thus, the GMPHD filter is efficient dealing with the false positive observations, but needs to be extended and improved by additional techniques for MOT in video data system.

As demand increases on online and real-time tracker in video-based tracking system, the PHD filter have been an emerging tracking model, recently. Song *et al.* [16] extended the GMPHD filter based tracking with the two-stage hierarchical data association strategy to recover fragmented and lost tracks. They defined the affinity in the track-to-track association step by using tracks' linear motion and color histogram appearance. This approach is an intuitive implementation of the GMPHD filter to handle tracking problems, but cannot correct the false associations already made in the detection-to-track association. T. Kutschbach *et al.* [42] added the kernelized correlation filters (KCF) [53] for online appearance update to overcome occlusion with the naive GMPHD filtering process. They showed a robust online appearance learning to re-find the IDs of the lost tracks. However, updating appearance information of all objects at every frame requires heavy computing resources. R. Sanchez-Matilla *et al.* [22] proposed a detection confidence based MOT model with the PHD filter.

Strong (high confidence) detections initiate and propagate tracks but weak (low confidence) detections only propagate existing tracks. This strategy works well when the detection results are reliable. However, the tracking performance is dependent on the detection performance, and especially weak to long-term missed detections. Z. Fu *et al.* [18] adopted an adaptive gating technique and an online group-structured dictionary (appearance) learning strategy into the GMPHD filter. They made the GMPHD filter have a sophisticated tracking process and fit to video based MOT.

Grouping approach e.g., relative motion and topological model, already have been exploited in [21], [23]. The key difference between their methods and ours is that [21], [23] consider the relations between all objects in a scene but we only consider topological information in the group of occluded objects. Grouping only the occluded objects exclude trivial solutions (associations) which focuses on solving sub-problems and reduces computing time.

## III. PROPOSED ONLINE MULTI-OBJECT TRACKING FRAMEWORK

In this section, we briefly introduce the general tracking process of the Gaussian mixture probability hypothesis density (GMPHD) filter in Subsection III-A. In III-B, we address how to extend the GMPHD filter with the hierarchical data association strategy in video-based online MOT systems.

### A. The GMPHD Filter

The Gaussian mixture model (GMM) of the GMPHD filter includes means, covariances, and weights which are propagated at every time stamp as follows; *Initialization*, *Prediction*, *Update*, and *Pruning* steps. We employ this basic process of the GMPHD filter but revise fit to the video-based MOT system.

$$X_k = \{x_k^1, \dots, x_k^{I_k}\}, \quad (1)$$

$$Z_k = \{z_k^1, \dots, z_k^{J_k}\}, \quad (2)$$

where  $X_k$  and  $I_k$  denote a set of objects' states and the number of them at time  $k$ , respectively. A state vector  $x_k$  is composed of  $(c_x, c_y, v_x, v_y)$ , where  $c_x$ ,  $c_y$ ,  $v_x$ , and  $v_y$  indicate the x-axis center point of the bounding box, the y-axis center point of the bounding box, the x-axis velocity, and the y-axis velocity, respectively. Likewise,  $Z_k$  and  $J_k$  denote a set of observations (detection responses) and the number of them at time  $k$ , respectively. An observation  $z_k$  is composed of  $(c_x, c_y)$ , where  $c_x$  and  $c_y$  indicate the x-axis and the y-axis center of the detection bounding box, respectively. Equation (3) and (4) describe the basic notations of state and observation.

$$x_k^i = \{c_{x,k}, c_{y,k}, v_{x,k}, v_{y,k}\}^T, \quad (3)$$

$$z_k^j = \{c_{x,k}, c_{y,k}\}^T. \quad (4)$$

The tracking process of the GM-PHD filter is composed of four steps: *Initialization*, *Prediction*, *Update*, and *Pruning* as follows.

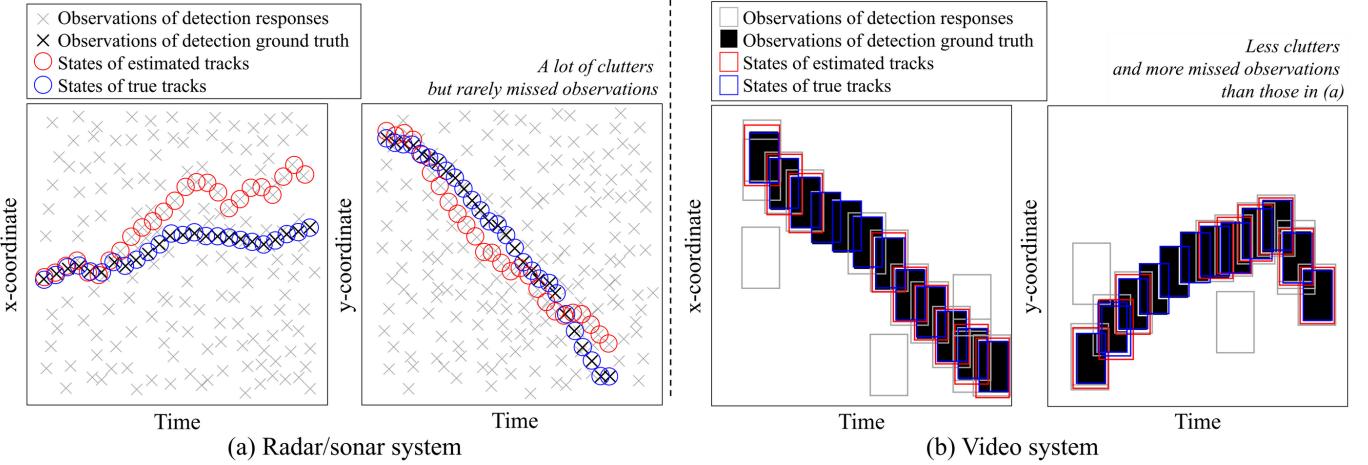


Fig. 1: Comparison between (a) radar/sonar and (b) video system in terms of input and output, i.e., observations (detection) and states (tracking). The radar/sonar sensors receive a lot of clutters (false positive error) but rarely miss objects (false negative error), whereas the detector in video data tends to receive a few clutters around the objects and misses more objects than the radar/sonar sensors do.

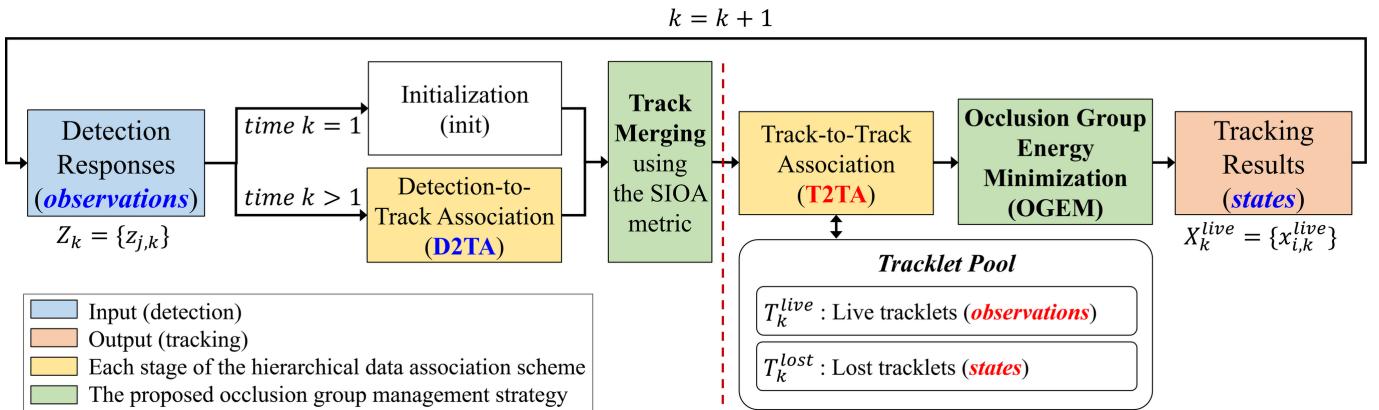


Fig. 2: Flow chart of the proposed online multi-object tracking framework. The red dotted line divides the proposed hierarchical data association into two stages. Each stage and its states and observations are marked as **blue** and **red**, i.e., D2TA and T2TA, respectively. The key components of this chart, such as init, D2TA, Merge, T2TA, OGEM, live and lost tracklets are used in Figure 3 and Figure 9, also.

### **Initialization:**

$$\sum_{i=1}^{I_0} w_0^i \mathcal{N}(x; m_0^i, P_0^i), \quad (5)$$

where the GMM is initialized by the initial observations from the detection responses. Besides, when an observation fails to find the association pair, i.e., updating object state, the observation initializes a new Gaussian model (a new state). Gaussian probability function  $\mathcal{N}$  represents tracking objects with weight  $w$ , mean vector  $m$ , object state vector  $x$ , and covariance matrix  $P$ . At this step, we set the initial velocities of mean vector to zeros. Each weight is set to the normalized confidence value of the corresponding detection response.

### *Prediction:*

$$\sum_{i=1}^{I_{k-1}} w_k^i \mathcal{N}(x; m_{k-1}^i, P_{k-1}^i), \quad (6)$$

$$m_{k|k-1}^i = F m_{k-1}^i, \quad (7)$$

$$P_{k|k-1}^i = Q + F P_{k-1}^i (F)^T, \quad (8)$$

where we assume that the GMM representing the objects' states was initialized or active at the previous frame  $k-1$  in (6). In (7) and (8),  $F$  is the state transition matrix and  $Q$  is the process noise covariance matrix.  $F$  and  $Q$  are constants in our tracker. Then, we can predict the state at time  $k$  using the Kalman filtering. In (7),  $m_{k|k-1}^i$  is derived by using the velocity of  $m_{k-1}^i$ . Covariance  $P_{k|k-1}^i$  is also predicted by the Kalman filtering method in (8).

### Update:

$$\sum_{i=1}^{I_{k|k-1}} w_k^i(z) \mathcal{N}(x; m_{k|k}^i, P_{k|k}^i), \quad (9)$$

$$q_k^i(z) = \mathcal{N}(z; Hm_{k|k-1}^i, R + HP_{k|k-1}^i(H)^T), \quad (10)$$

$$w_k^i(z) = \frac{w_{k|k-1}^i q_k^i(z)}{\sum_{l=1}^{I_{k|k-1}} w_{k|k-1}^l q_k^l(z)}, \quad (11)$$

$$m_{k|k}^i(z) = m_{k|k-1}^i + K_k^i(z - Hm_{k|k-1}^i), \quad (12)$$

$$P_{k|k}^i = [I - K_k^i H] P_{k|k-1}^i, \quad (13)$$

$$K_k^i = P_{k|k-1}^i(H)^T (HP_{k|k-1}^i(H)^T + R)^{-1}, \quad (14)$$

where the goal of update step is deriving (9). First, we should find an optimal observation  $z_k$  at time  $k$  to update a Gaussian model. The optimal  $z$  makes  $q_k$  be the maximum in (10).  $R$  denotes the observation noise covariance.  $H$  denotes the observation matrix to transit a state vector to an observation vector. Both  $R$  and  $H$  are constants in our application. In the perspective of application, the update step involves data association. Updating the Gaussian state models follows finding the optimal observations updating the states through the data association. After finding the optimal  $z$ , the GMM is updated to (9) through (10), (11), and (14), (13), (12).

### Pruning:

$$\tilde{X}_k = \{m_k^i : w_k^i \geq \theta_w, i = 1, \dots, I_k\}, \quad (15)$$

$$\tilde{W}_k = \{w_k^i : m_k^i \in \tilde{X}_k, i = 1, \dots, I_k\}, \quad (16)$$

$$\tilde{W}_k = \{\tilde{w}_{k,1}, \dots, \tilde{w}_{k,\tilde{I}_k}\}, \tilde{I}_k = |\tilde{W}_k|, \quad (17)$$

$$w_k^i = \frac{\tilde{w}_k^i}{\sum_{l=1}^{\tilde{I}_k} \tilde{w}_k^l}, \quad (18)$$

$$X_k = \tilde{X}_k, \quad (19)$$

where the states with the weight under threshold  $\theta_w$  are pruned as in (15). We experimentally set  $\theta_w$  to 0.1. Then the weights of the surviving states are normalized as shown in (18). The pruning step handles the false positive tracks by the false positive detections.

The GMPHD filter [3] is specialized in handling false positives e.g., clutters and noise. However, tracking systems have the different problems, depending on their domains as shown in Figure 1, where input and output indicate detection results (observations) and tracking results (states), respectively. As presented in [4] and Figure 1-(a), at radar/sonar systems, the sensors receive uncountable detection responses with a lot of clutters but objects are rarely missed. On the other hand, as shown in Figure 1-(b), the video data based detectors observe less clutters and miss more objects than the radar/sonar sensors do. The conventional GMPHD filter is effective to handle the clutters (false positive) but missed detections cause the new tracking problems in video data system (false negative). Thus, we propose the GMPHD filtering based tracker with a hierarchical data association strategy.

### B. Hierarchical Data Association

Video-based tracking systems have inherent problems as shown in Figure 1-(b). Generally, when objects are not

detected, the objects' IDs are frequently changed and the tracks are fragmented if only detection-to-track association is employed. To prevent these problems by missing objects, we take advantage of a hierarchical data association (HDA) strategy which has been widely used in many online multi-object tracking methods [14], [16], [17], [22], [25]. Thus, in this paper, we propose a simple HDA scheme with just two stages. The proposed HDA includes detection-to-track (D2T) and track-to-track (T2T) associations. We implement the both association methods with the GMPHD filtering process as given in III-A. Also, we derive a cost function from (11) of the GMPHD filtering process as follows:

$$Cost(x_{k|k-1}^i, z_k^j) = -\ln w_k^i(z_k^j), \quad (20)$$

where  $w_k^i$  indicate the weight value, assuming that observation  $z_k^j$  updates state  $x_{k|k-1}^i$ . We use  $-\ln w_k^i(z_k^j)$  as a cost between  $x_{k|k-1}^i$  and  $z_k^j$ . Then, cost matrix  $C$  can be built by every pair between state set  $X_{k|k-1}$  and observation set  $Z_k$  as follows:

$$C[i, j] = Cost(X_{k|k-1}[i], Z_k[j]). \quad (21)$$

When the cost matrix  $C$  is built, the Hungarian algorithm is used to solve it. Then, the optimal pairs between observations and states are found, and consequently state  $x_{k|k-1}$  is updated to  $x_k$  in D2T and T2T associations. In III-B1 and III-B2, we introduce the definition of observations and states in each association stage with more detail usage of the cost function.

1) *Detection-to-Track Association (D2TA, Stage 1):* In D2TA, observation set  $Z_k$  is filled with detection responses at time  $k$ . We assume that state set  $X_{k|k-1}$  already exists from time  $k-1$ , and then  $X_{k|k-1}$  is predicted by using the Kalman filtering as shown in (6)-(8). Thus, the cost matrix  $C_{D2T}$  is easily calculated with these sets  $X_{k|k-1}$  and  $Z_k$ .

2) *Track-to-Track Association (T2TA, Stage 2):* In T2TA, a simple temporal analysis of tracklet is conducted. A tracklet means a fragment of the track, and becomes a calculation unit. Before T2TA, all tracklets are categorized into two types, according to success or failure of tracking at the present time  $k$  as follows:

$$T_k^{lost} \cup T_k^{live} = T_k^{all}, \quad (22)$$

$$T_k^{lost} \cap T_k^{live} = \phi, \quad (23)$$

$$T_k^{lost} = \{\tau_{1,k}^{lost}, \dots, \tau_{i,k}^{lost}\}, \quad (24)$$

$$\tau_{i,k}^{lost} = \{a_s^i, \dots, a_t^i\}, \quad 0 \leq s < t < k, \quad (25)$$

$$T_k^{live} = \{\tau_{1,k}^{live}, \dots, \tau_{j,k}^{live}\}, \quad (26)$$

$$\tau_{j,k}^{live} = \{a_s^j, \dots, a_t^j\}, \quad 0 \leq s < t, t = k, \quad (27)$$

where "live" indicates that tracking succeeds at time  $k$ . "lost" indicates that tracking fails at time  $k$ . Then, for the T2TA, observation set  $Z_k$  is filled with the first (oldest) elements  $a_s^j$ s of "live" tracklets. However, the state set  $X_{k|k-1}$  is not filled with the last (most recent) elements  $a_t^j$ s of "lost" tracklets.

## Track State Machine

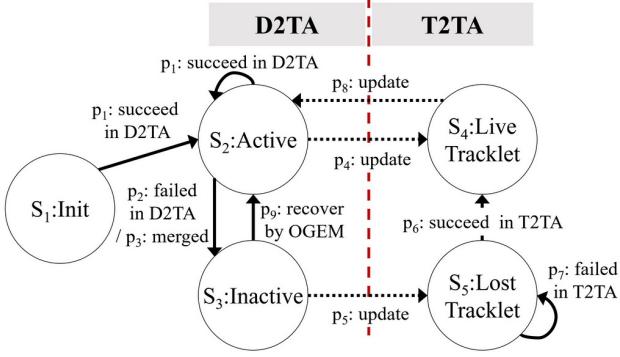


Fig. 3: In the proposed tracking framework, an object with state  $S_1$  : *Init* is transited within the defined states  $\{S_2, S_3, S_4, S_5\}$  by the state-transition functions  $\{p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8, p_9\}$ .

One prediction step is needed as follows:

$$a_s^i = \{c_{x,s}, c_{y,s}, v_{x,s}, v_{y,s}\}^T, \quad (28)$$

$$a_t^i = \{c_{x,t}, c_{y,t}, v_{x,t}, v_{y,t}\}^T, \quad (29)$$

$$x_t^i = \{c_{x,t}, c_{y,t}, \frac{c_{x,t} - c_{x,s}}{t-s}, \frac{c_{y,t} - c_{y,s}}{t-s}\}^T, \quad (30)$$

$$x_{k|k-1}^i = F^{T2T} x_t^i, \quad (31)$$

$$F^{T2T} = \begin{pmatrix} 1 & 0 & d_f & 0 \\ 0 & 1 & 0 & d_f \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, \quad (32)$$

$$d_f(i, j) = \text{frame difference between } a_t^i \text{ and } a_s^j. \quad (33)$$

In (30)  $\frac{c_{x,t} - c_{x,s}}{t-s}$  and  $\frac{c_{y,t} - c_{y,s}}{t-s}$  are the averaged velocities in terms of x-axis and y-axis, respectively. The velocities are calculated by subtracting the center position of the first object state  $a_s^i$  from that of the last state  $a_t^i$ , and dividing it by the frame difference  $t - s$  which is equivalent to the length of “lost” tracklet  $\tau_{i,k}^{lost}$ . D2TA has the identical time interval “1” between states and observations in transition matrix  $F$ , whereas in T2TA , each cost of matrix  $C_{T2T}$  has different time interval (frame difference) between states and observations. Variable  $d_f$  depends on which state of “lost” tracklet and observation of “live” tracklet are paired. (31) means the prediction process of state with linear motion analysis. Finally, the cost matrix  $C_{T2T}$  is filled by (31) and the oldest element  $a_s^j$  of live tracklet  $\tau_{j,k}^{live}$ .

The pseudo-code in Algorithm 1 includes the procedures presented in this section. *Initialization*, *Prediction*, Cost-minimization, *Update*, and *Pruning* in D2TA correspond to each of line 5-10, 13-15, 16-21, 22-24, and 25-27 in Algorithm 1. Tracklet-categorization, Cost-minimization, *Update* in T2TA correspond to line 36-44, 49-54, and 55-66 in Algorithm 1, respectively.

## IV. OCCLUSION GROUP MANAGEMENT SCHEME

In Section III, we addressed that the proposed online multi-object tracking framework is based on the GMPHD filter

## Algorithm 1 Proposed Online MOT Algorithm

---

$\triangleright k$  : the current frame number  
 $\triangleright X_{k-1}$  : a set of states at time  $k-1$   
 $\triangleright Z_k$  : a set of observations at time  $k$   
 $\triangleright \sigma_m$  : threshold for track merging  
 $\triangleright \tau_{T2T}$  : the minimum track length for T2TA  
 $\triangleright \theta_{T2T}$  : the maximum frame interval for T2TA  
 $\triangleright T^{live}$  : a {key:id,value:tracklet} set of live tracklets  
 $\triangleright T^{lost}$  : a {key:id,value:tracklet} set of lost tracklets

```

1: procedure GMPHD_OGM( $k, X_{k-1}, Z_k, \sigma_m, \tau_{T2T}, \theta_{T2T}, T^{live}, T^{lost}$ )
2:    $l = |X_{k-1}|;$  // the number of states
3:    $m = |Z_k|;$  // the number of observations
4:    $G_{k-1}, G_k;$  // a set of occlusion groups at time  $k-1$  and  $k$ .
5:   if  $k = 1$  or  $l = 0$  then
6:     Initialize states  $X'_k$  with  $Z_k$ ;
7:      $G_{k-1} = G_k$ ;
8:      $X_k = \text{MERGE}(X'_k, \sigma_m, G_k)$ ;
9:     return  $X_k$ ;
10:    end if
/* 1. Detection-to-Track Association (D2TA) */
11:     $C_{D2T}[1 \dots l][1 \dots m];$  // for cost matrix
12:     $P_{D2T}[1 \dots l];$  // for pairing observations' indices
/*predict states  $X_{k-1}$  to be  $X_{k|k-1}$ */
13:    for  $i = 1$  to  $l$  do
14:       $X_{k|k-1}[i] = \text{PREDICT}(X_{k-1}[i]);$ 
15:    end for
/*calculate the GMPHD filter cost matrix  $C_{D2T}$ */
16:    for  $i = 1$  to  $l$  do
17:      for  $j = 1$  to  $m$  do
18:         $C_{D2T}[i][j] = \text{COST}_{D2T}(X_{k|k-1}[i], Z_k[j]);$ 
19:      end for
20:    end for
/*find min-cost pairs by the Hungarian method*/
21:     $P_{D2T} = \text{HungrianMethod}(C_{D2T});$ 
/*update and birth states*/
/*update  $X_{k|k-1}$  with the min-costly observations*/
22:    for  $i = 1$  to  $l$  do
23:       $X'_k[i] = \text{UPDATE}(X_{k|k-1}[i], Z_k[P_{D2T}[i]]);$ 
24:    end for
/*prune  $X_{k|k-1}$  with the weight under 0.1*/
25:    for  $i = 1$  to  $l$  do
26:       $X'_k[i] = \text{PRUNE}(X_{k|k-1}[i]);$ 
27:    end for
28:    for  $j = 1$  to  $m$  do
29:      if  $Z_k[j]$  is not assigned to update any state then
30:        Initialize newly birth state  $x$  with  $Z_k[j]$ ;
31:         $X'_k = X'_k \cup \{x\};$ 
32:      end if
33:    end for
/* 2. Merge States and Find Occlusion Groups */
34:     $G_{k-1} = G_k;$ 
35:     $X_k = \text{MERGE}(X'_k, \sigma_m, G_k);$ 
/*manage tracklet pool after D2TA and MERGE*/
36:    for  $i = 1$  to  $|X_k|$  do
37:      if  $X_k[i]$  is active then
38:        update  $T^{live}[X_k[i].id]$  with  $X_k[i]$ ;
39:        delete  $T^{lost}[X_k[i].id]$ ;
40:      else
41:        update  $T^{lost}[X_k[i].id]$  with  $X_k[i]$ ;
42:        delete  $T^{live}[X_k[i].id]$ ;
43:      end if
44:    end for

```

---

ing theory with the two-stage hierarchical data association. However, the tracking results from that framework still give uncertainty to us, even if we effectively extend the conventional GMPHD filter to be suitable for video-based tracking system. Thus, to handle it, we define two types of tracking problems and provide a solution. One is an intrinsic occlusion and the other is an extrinsic occlusion. The intrinsic occlusion is defined when the number of detection responses on one

---

```

/* 3. Track-to-Track Association (T2TA) */
45:    $t_1 = |T^{lost}|;$                                 // the number of lost tracklets
46:    $t_2 = |T^{live}|;$                                 // the number of live tracklets
47:    $C_{T2T}[1 \dots t_1][1 \dots t_2];$                 // for cost matrix
48:    $P_{T2T}[1 \dots t_1];$                             // for pairing observations' indices
49:   /*calculate the GMPHD filter cost matrix  $C_{T2T}*/$ 
50:   for  $i = 1$  to  $t_1$  do
51:     for  $j = 1$  to  $t_2$  do
52:        $C_{T2T}[i][j] = COST_{T2T}(T^{lost}[i], T^{live}[j], \tau_{T2T}, \theta_{T2T});$ 
53:     end for
54:   end for
55:   /*find min-cost pairs by the Hungarian method*/
56:    $P_{T2T} = HungarianMethod(C_{T2T});$ 
57:   /*update tracklets and manage tracklet pool after T2TA*/
58:   for  $i = 1$  to  $l$  do
59:      $X'_k[i] = UPDATE(X_{k|k-1}[i], Z_k[P_{D2T}]);$ 
60:   end for
61:   for  $i = 1$  to  $|X_k|$  do
62:     if  $X_k[i]$  is active then
63:       update  $T^{live}[X_k[i].id]$  with  $X_k[i];$ 
64:       delete  $T^{lost}[X_k[i].id];$ 
65:     else
66:       update  $T^{lost}[X_k[i].id]$  with  $X_k[i];$ 
67:       delete  $T^{live}[X_k[i].id];$ 
68:     end if
69:   end for
70:   /* 4. Occlusion Group Energy Minimization (OGEM) */
71:   if  $k > 1$  and  $|G_{k-1}| > 0$  then
72:     OGEM( $k, G_{k-1}, X_k$ );
73:   /*manage tracklet pool after OGEM*/
74:   for  $i = 1$  to  $|X_k|$  do
75:     if  $X_k[i]$  is active then
76:       update  $T^{live}[X_k[i].id]$  with  $X_k[i];$ 
77:       delete  $T^{lost}[X_k[i].id];$ 
78:     else
79:       update  $T^{lost}[X_k[i].id]$  with  $X_k[i];$ 
80:       delete  $T^{live}[X_k[i].id];$ 
81:     end if
82:   end for
83:   return  $X_k;$                                 // return final states  $X_k$ 
84: end procedure

```

---

object is more than one. Generally, it makes the object ID switched and false positive tracks. The extrinsic occlusion is defined when the number of detection responses on objects, occluded each other, is less or more than the number of the occluded objects. That can cause false negative and positive tracks, respectively. Figure 8 and 9 show the defined tracking issues well. The false positive detections made by intrinsic and extrinsic occlusions inevitably generate false tracks as shown in the second row Figure 9 (D2TA), if appropriate techniques do not handle it. To resolve the two types of problems, we design a new occlusion group management (OGM) scheme. OGM consists of “*Track Merging*” and “*Occlusion Group Energy Minimization (OGEM)*” routines which execute just after D2TA and T2TA, respectively. Figure 2 briefly shows the tracking pipeline with those two components of OGM. Consequently, our occlusion group management technique not only decreases false positive tracking results but also prevents occluded tracks from false “track merging”. The effectiveness of the proposed OGM method is discussed in Section V in more detail.

#### A. Track Merging

Merging the neighboring objects’ states with the distances under a threshold is proposed in [3] already. However, it

#### Algorithm 2 Track Merging using the SIOA Metric

---

```

function MERGE( $X_k, \sigma_m, G_k$ )
2:    $l = |X_k|;$                                 //  $l$  : the number of states  $X_k$ 
3:   Let  $M[1 \dots l][1 \dots l]$  be the array set to all false;
4:   /* measure occlusion ratio between all states
5:    * by using the SIOA metric */
6:   for  $i = 1$  to  $l$  do
7:     for  $j = i + 1$  to  $l$  do
8:        $r_{occ} = SIOA_{X_k[i], X_k[j]};$            // SIOA occlusion ratio.
9:       if  $r_{occ} > \sigma_m$  then
10:         $M[i][j] = true;$                       // check to be merged
11:         $M[j][i] = true;$                       // double check
12:        else if  $r_{occ} \leq \sigma_m$  and  $r_{occ} > 0$  then
13:           $id_i = X_k[i].id, id_j = X_k[j].id;$ 
14:          if  $id_i < id_j$  then
15:             $G_k[id_i] = G_k[id_i] \cup \{X_k[i], X_k[j]\};$ 
16:          else
17:             $G_k[id_j] = G_k[id_j] \cup \{X_k[i], X_k[j]\};$ 
18:          end if
19:        end if
20:      end for
21:    end for
22:    /* merge the states where SIOA value >  $\sigma_m$  */
23:    for  $i = 1$  to  $l$  do
24:      for  $j = i + 1$  to  $l$  do
25:        if  $M[i][j] = true$  then
26:          if  $X_k[i].id < X_k[j].id$  then
27:             $X_k[i] = 0.9 * X_k[i] + 0.1 * X_k[j];$ 
28:            Deactivate state  $X_k[j];$ 
29:          else
30:             $X_k[j] = 0.9 * X_k[j] + 0.1 * X_k[i];$ 
31:            Deactivate state  $X_k[i];$ 
32:          end if
33:        end if
34:      end for
35:    end function

```

---

can only reflect point-to-point distance without considering regional information e.g., overlapping ratio between visual objects (bounding boxes). To measure the overlapping ratio, the intersection-over-union (IOU) metric has been widely used which was originally designed to measure mAP in object detection research fields [55], [56]. However, the IOU metric is nice to refine the detection bounding boxes but not adjustable to measure overlapping ratio for merging the objects. Figure 8 explains that reason by a case study. The case study mainly assumes that the number of detection responses (observations) is larger than the number of real objects. When the observations most likely include false positive detections, the object states by those observations also most likely become the false positive states. So, to handle and consider the characteristic of those observations with the false positive errors, we propose a new metric named as sum-of-intersection-over-area (SIOA). The IOU and SIOA metrics are formulated as follows:

$$IOU_{AB} = \frac{\text{area}(A) \cap \text{area}(B)}{\text{area}(A) \cup \text{area}(B)}, \quad (34)$$

$$\begin{aligned} SIOA_{AB} = \\ 0.5 * \left( \frac{\text{area}(A) \cap \text{area}(B)}{\text{area}(A)} + \frac{\text{area}(A) \cap \text{area}(B)}{\text{area}(B)} \right), \end{aligned} \quad (35)$$

**Algorithm 3** Occlusion Group Energy Minimization

```

     $\triangleright k$  : the current frame number
     $\triangleright G_{k-1}$  : a set of occlusion groups at time  $k-1$ 
     $\triangleright X_k$  : a set of states at time  $k$ 

1: function OGEM( $k, G_{k-1}, X_k$ )
2:    $l = |G_{k-1}|$ ;           //  $l$  : the number of the groups  $G_{k-1}$ .
3:    $n = |X_k|$ ;           //  $n$  : the number of the states  $X_k$ .
   /*build the GMMs for all occlusion groups at time  $k-1$ */
   /*a GMM is used for the defined energy function in (36)*/
4:   for  $i = 1$  to  $l$  do
5:      $p_i = |G_{k-1}[i]| P_2$            //the number of topological vectors.
6:      $GMM[1 \dots p_i]$ ;           //the Gaussian mixture for a group.
7:     for  $j = 1$  to  $p_i$  do           //iterate topologies in a group.
8:       Initialize a Gaussian mixture  $GMM[j]$  with
9:         the mean vectors  $m$  having topological info and
10:        the covariance matrix  $R$  as defined in (36)
11:    end for
12:     $E_{min} = DBL\_MAX$ ;           //variable for the min-energy.
13:     $h_{min} = 0$ ;           //index to the optimal hypothesis.
14:    for  $h = 1$  to  $|H|$  do           //iterate topological hypotheses.
15:      if  $E(h) < E_{min}$  then           //find the optimal hypothesis.
16:         $h_{min} = h$ ;
17:         $E_{min} = E(h)$ ;
18:      end if
19:    end for
20:    Update  $G_{k-1}[i]$  with  $h_{min}$ ;
21:    for  $g$  in  $G_{k-1}[i]$  do           //iterate group  $G_{k-1}[i]$ .
22:      Find the state  $x$  with  $g.id$  in  $X_k$ ;
23:      if  $x$  is in  $X_k$  then
24:         $X_k[g.id] = g$ ;
25:      else
26:         $X_k = X_k \cup \{g\}$ ;
27:      end if
28:    end for
29:  end for
30:  return  $X_k$ ;
31: end function

```

where A and B indicate two different objects. *area* represent a bounding box ( $x, y, width, height$ ). Algorithm 2 describes the proposed track merging method. Track merging with the SIOA metric follows after the D2T association as presented in Subsection III-B1 and Figure 2.

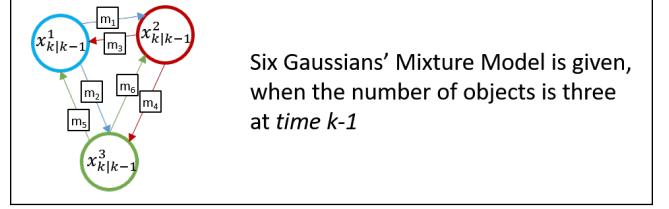
### B. Occlusion Group Energy Minimization

The occlusion group energy minimization method is devised to prevent the true objects which are occluded to others from false merging. In other words, track merging after detection-to-track association (D2TA) may merge occluded objects with correct number of observations into the states with less number of real objects. That can cause tracking errors such as false negative and fragmented tracks.

Thus, we propose a new energy minimization model to prevent false merging, named as “*Occlusion Group Energy Minimization (OGEM)*”. Each group of occluded objects has an energy function represented by a Gaussian mixture model (GMM) as follows:

$$E(h) = -\ln \sum \mathcal{N}(t|m, R), \quad (36)$$

where  $h$ ,  $t$ ,  $m$ , and  $R$  indicate hypothesis, topological vector, mean vector, and Gaussian covariance marix (noise), respectively. A Gaussian probability function  $\mathcal{N}$ , i.e., component of the GMM, indicates a topological position vector between two objects in an occlusion group, which is given at time  $k-1$ . The Gaussian function has a mean vector  $m$  which denotes



Six Gaussians' Mixture Model is given, when the number of objects is three at time  $k-1$

Ordered Arrangements of Three Objects : Six Hypotheses

$$= \{h1, \dots, h6\} = \{x^1x^2x^3, x^1x^3x^2, x^2x^1x^3, x^2x^3x^1, x^3x^1x^2, x^3x^2x^1\}$$

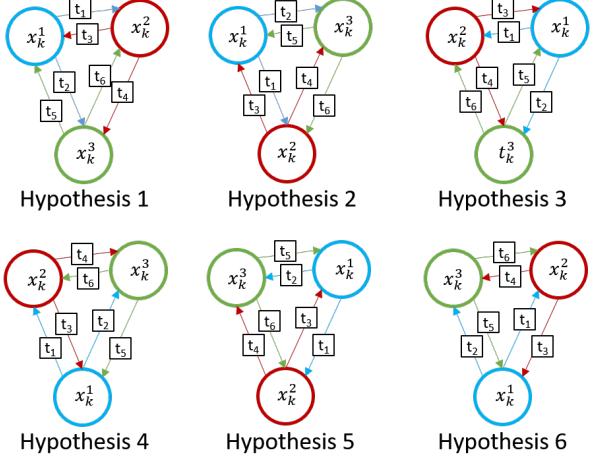


Fig. 4: Illustration of the proposed occlusion group energy minimization represented by the Gaussian mixture model. Six hypotheses exist in the case of three occluded objects.

the topological position, i.e., relative position between the predicted center positions at time  $k$  of the objects in the group. Those objects are denoted by  $x_{k|k-1}$  and the notation  $k|k-1$  indicates the prediction at time  $k$  from position and velocity at time  $k-1$ . If there are three occluded objects in a group, six hypotheses exists as shown in Figure 4. One hypothesis is a set of six topological vectors  $\{t_1, t_2, t_3, t_4, t_5, t_6\}$ . For example,  $m_1$  is calculated by  $x_{k|k-1}^2 - x_{k|k-1}^1$  and  $t_1$  is calculated by  $x_{k|k-1}^2 - x_{k|k-1}^1$ . In the case that an object state  $x_{k|k-1}^d$  becomes inactive by false merging in occlusion, we build a new hypothesis using a  $x_{k|k-1}^d$  as a dummy. Then the dummy added hypothesis recovers the false merged object. If there are  $n$  occluded objects in a group,  $n(n-1)$  hypotheses exists with the condition  $1 < n < 4$ . Then with these topological models, we can find an optimal one among all hypotheses making the group cost minimal.

Whereas track merging runs after the D2TA, the proposed occlusion group energy minimization follows Track-to-Track Association (T2TA) as described in Figure 2. Figure 9 includes some examples to explain that the proposed group energy minimization complements track merging step. The tracking stage at frame 42 explains it.

In summary, both “*Track Merging*” and “*Occlusion Group Energy Minimization*” procedures assume occlusion situations, and the GMHD filtering is adopted as the main framework. The pseudo-code examples of proposed occlusion group management scheme are described in Algorithm 2 and 3. Also, both methods correspond to line 8, 35 and 67-78 in

Algorithm 1 which represents the whole tracking framework. From now on we use GMPHD-OGM as the abbreviation for the proposed algorithm, online multi-object tracking with the GMPHD filter and occlusion group management.

## V. EXPERIMENTS

In this section, we present development environment including parameter settings, and also discuss evaluation results of the GMPHD-OGM tracker which include an ablation study with baselines and comparisons to state-of-the-art methods. The GMPHD-OGM tracker is implemented by Visual C++ with OpenCV3.4.1 and boost1.61.0 libraries, and without any GPU-accelerated libraries such as CUDA. All experiments are conducted on Windows 10 with Intel i7-7700K CPU @ 4.20GHz and DDR4 32.0GB RAM.

### A. Parameter Setting

Our proposed tracking framework involves several parameter settings. Parameter  $\sigma_m$  indicates the threshold for “Track Merging” which is set to 0.5 in terms of the SIOA metric. 0.5 is set not only empirically set but by considering the occlusion cases between size-variant bounding boxes as shown in Case 4 and 5 of Figure 8.  $\tau_{T2T}$  and  $\theta_{T2T}$  are related to track-to-track association (T2TA) of the hierarchical data association, whose parameters are selected adaptively, scene-by-scene. The optimal values of  $\tau_{T2T}$  and  $\theta_{T2T}$  are gained from the ablation study presented in Figure 7. We use the optimal parameter settings from training to test sequences. These three parameters are summarized in Table I.

The GMPHD filtering process has a set of static parameters. The matrices  $F$ ,  $Q$ ,  $P$ ,  $R$ , and  $H$  are used in *Prediction Step* and *Update Step*. Also,  $\theta_w$  is used in *Pruning Step*. Experimentally, we set the parameters for the GMPHD filter’s tracking process as follows:

$$F = \begin{pmatrix} 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{pmatrix}, Q = \frac{1}{2} \begin{pmatrix} 5^2 & 0 & 0 & 0 \\ 0 & 10^2 & 0 & 0 \\ 0 & 0 & 5^2 & 0 \\ 0 & 0 & 0 & 10^2 \end{pmatrix},$$

$$P = \begin{pmatrix} 5^2 & 0 & 0 & 0 \\ 0 & 10^2 & 0 & 0 \\ 0 & 0 & 5^2 & 0 \\ 0 & 0 & 0 & 10^2 \end{pmatrix}, R = \begin{pmatrix} 5^2 & 0 \\ 0 & 10^2 \end{pmatrix},$$

$$H = \begin{pmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{pmatrix}, \theta_w = 0.1,$$

TABLE I: Parameter Settings for “Track Merging” and track-to-track association (T2TA).

Symbol	Description	Value
$\sigma_m$	threshold for track merging.	0.5
$\tau_{T2T}$	the minimum track length for T2TA	1, 2, 3
$\theta_{T2T}$	the maximum frame interval for T2TA	5 to 100

### B. Evaluation Results

In this section, we evaluated the propose method with state-of-the-art online [14]–[16], [18]–[25], [37]–[43] and offline [26]–[36], [44]–[50] MOT methods in terms of the CLEAR-MOT metrics [54]. The CLEAR-MOT metrics gracefully measure multi-object tracking performance from the detailed perspectives such as multi-object tracking accuracy (MOTA), multi-object tracking precision (MOTP), mostly tracked objects (MT), mostly lost objects (ML), the total number of false positives (FP), the total number of false negatives (missed tracks, FN), the total number of identity switches (IDS), the total number of times that a trajectory is fragmented (Frag), and processing speed (frames per second, FPS). Among these metrics, MOTA is normally proposed as the key metric, because it considers three error sources including FP, FN, and IDS, comprehensively. The evaluation results contain not only the tracking results on the MOT15 and MOT17 test datasets but also an ablation study on the MOT15 training dataset.

First, in the ablation study, we employ the two baseline methods, to find optimal parameters settings and to prove the effectiveness of the proposed method. One is the GMPHD filter based tracker with the hierarchical data association (HDA) and without the occlusion group management (OGM). The other is the GMPHD filter based tracker with HDA and OGM by using the IOU metric for measuring occlusion ratio. We name these three methods as GMPHD-HDA, GMPHD-OGM (/w IOU), and GMPHD-OGM (/w SIOA) as shown in Table II and Figure 7. The GMPHD-OGM (/w SIOA) method is our final tracking model. The scene-by-scene optimal parameter settings of those three methods are obtained by another ablation study as shown in Figure 7. The same  $\tau_{T2T}$  and  $\theta_{T2T}$  settings are applied to the whole training sequences with the range {1, 2, 3} and {5, 10, 20, 30, 50, 70, 100}, respectively. GMPHD-OGM (/w IOU) is improved over GMPHD-HDA in terms of the upper bound of tracking accuracy (the maximum MOTA). GMPHD-OGM (/w SIOA) shows that upper bound and lower bound of tracking accuracy increases. Besides, with  $\theta_{T2T}$  over 20, the maximum and minimum values of MOTA increase on average. Figure 8 shows the comparison results of “Track Merging” between using the IOU metric and the SIOA metric when the detection results with a lot of false positives are given. Also, through the case study on occlusion, we observe that the IOU metric cannot consider size-variant detections with false positives and too sensitive to be used for merging as shown in Figure 8. On the other hand, the SIOA metric can consider the size-variant detection and the optimal value of merging threshold  $\sigma_m$  is decided to be 0.5 by the occlusion cases 4 and 5, empirically. Table II provides the quantitative results on the MOT15 training dataset with the best performance results on each sequence and  $\sigma_m = 0.5$ . GMPHD-OGM (w/ IOU) does not show outstanding improvement compared to GMPHD-HDA even though GMPHD-OGM (w/ IOU) takes more processing time since the OGM scheme runs whereas it does not in GMPHD-HDA. GMPHD-OGM (w/ SIOA) shows meaningful improvements in terms of MOTA, The ablation study proves that GMPHD-OGM is not

only overall improved but also more robust and less sensitive in parameters than baseline methods.

Figure 9 demonstrates some qualitative results of our tracking framework in view of overall process. Detection results (observations) initialize tracking objects (states). In the sequential tracking process, the states are associated with the proper observations by the detection-to-track association (D2TA) using the GMPHD filtering process. From false positive detections, false positive tracks can be generated and then “Track Merging” handles it. If objects are occluded and their IDs are switched, track-to-track association (T2TA) can recover their IDs. In the case that “Track Merging” merges true tracks (false merging), the occlusion group energy minimization (OGEM) process can recover it which optimizes energy of a group of occluded objects at current time by calculating the probability of the Gaussian mixture model as described in Subsection IV-B.

Table III and IV show the quantitative evaluations results on MOT15 and MOT17 test dataset, respectively. Those two benchmark datasets have crucial different characteristics. First, provided public detection results are different, MOT15 provides ACF [8] detector based detections, and MOT17 provides three types of detections such as DPM [9], FRCNN [10], and SDP [11]. Compared to the DNN based detectors FRCNN and SDP, ACF and DPM exploit hand-crafted features learning and models, and thus show relatively poor performance. DPM generates more false positives than FRCNN and SDP do, and especially ACF misses much more objects than others do. Thus, in MOT15, state-of-the-art trackers shows wider range of MOTA distribution than that in MOT17. Among online methods, the trackers with DNN [14], [37] shows the top MOTA scores in MOT15 and MOT17, respectively. Our method achieves the second best MOTA 30.7 vs. the best speed 169.5 fps in MOT15, but we think that the performance is competitive and enough to consider real-time application. In Figure 5-(a), the proposed method is located in a distinguished spot in terms of tracking accuracy (MOTA) vs. speed (fps). That also proves effectiveness of our occlusion group based object analysis (OGM), compared to other relation analysis between all objects in the scene [21], [23]. However, in MOT17, the speed of the proposed method decreases to 30.7 fps. That speed still belongs to real-time processing time but is not outstanding compared to other online methods. That is caused by the second different point between two datatsets. MOT15 includes 5783 frames with 721 tracks, 61440 bounding boxes, and 10.6 density i.e., the average number of objects a frame, whereas MOT17 includes 17757 frames with 2355 tracks, 564228 bounding boxes, and 31.8 density. Because MOT17 has the scenes not only with much higher density but also accurate detection results, those points increase tracking accuracy and processing time. Figure 5-(b) proves those facts where the performances of state-of-the-art methods are saturated on the spot with MOTA around 50 and speed under 5 fps. Even though our method achieves the second best MOTA and speed among online approaches in MOT17, the speed is drastically decreased compared to MOT15. Figure 6 explains that reason. In MOT17-03, the speed is around 10 fps since many objects appear in the

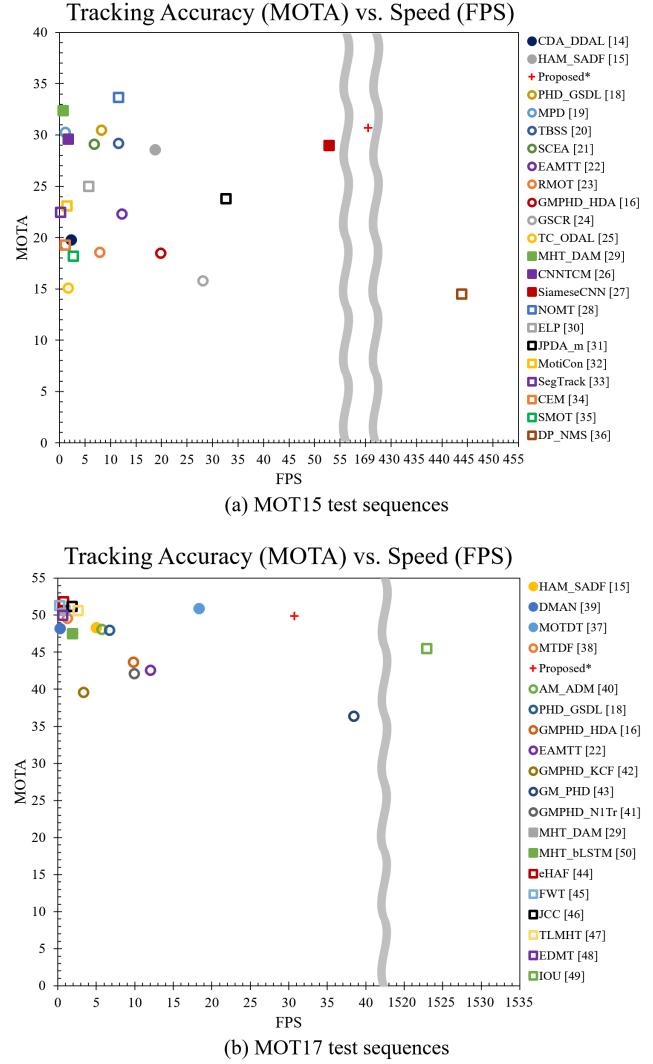


Fig. 5: Comparisons of tracking accuracy against speed with the state-of-the-art methods on the (a) MOT15 and (b) MOT17 test sequences.

scene with 69.8 density which means the average number of objects per frame. That makes the number of track-to-track association greatly increase. The proposed method is still comparative and positioned at meaningful spot for real-time application as shown in Figure 5-(b). In addition to our tracking algorithm (GMPHD-OGM), many PHD filter based online approaches [16], [18], [22], [38], [41]–[43] have been proposed in the past decade. Against them, GMPHD-OGM achieves not only the best MOTA, MOTP, MT, ML, FN, and speed scores on MOT15 but also the second best MOTA, speed, and best MT, FN, and Frag scores in MOT17. Especially, against to state-of-the-art online approaches, the proposed method is distinguished in terms of tracking accuracy (MOTA) vs. speed (fps), even though we did not utilize any complex visual features except bounding boxes. Also, the proposed tracker (GMPHD-OGM) against state-of-the-art algorithms including online with DNN even including offline, GMPHD-OGM shows the competitive MOTA versus speed as

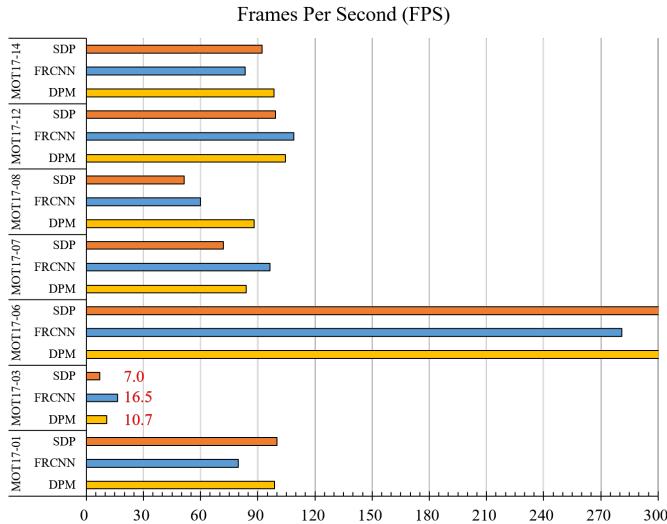


Fig. 6: Speed comparison of the proposed tracking method on MOT17 test dataset which provides three types of detection results for each scene, including DPM [9], FRCNN [10], and SDP [11].

described in Figure 5, Table III, and Table IV.

## VI. CONCLUSION AND FUTURE WORK

In this paper, we proposed an efficient online multi-object tracking framework with the GMPHD filter and the occlusion group management (OGM) named as GMPHD-OGM. In the proposed framework, our first contribution is that the Gaussian mixture probability hypothesis density (GMPHD) filter [3] is exploited to resolve MOT task. Since the GMPHD filter is originally designed to handle MOT in radar/sonar system, we should revise the filter to fit to video data system. To resolve missed tracks problem in the difference domain, we extended the conventional GMPHD filtering process with the hierarchical data association (HDA) strategy as explained in Figure 1. The second contribution is that to solve the occlusion problems, we proposed an occlusion group management (OGM) scheme. OGM is composed of “Track Merging” and “Occlusion Group Energy Minimization (OGEM)”. “Track Merging” reduced the number of false positives by merging them. The OGEM prevents false merging between true tracks. Both modules complement each other, and also instead of the IOU metric, we designed a new metric named as sum-of-intersection-over-area (SIOA) to measure the occlusion ratio between visual objects. The third is that the effectiveness of our tracker (GMPHD-OGM) was introduced by the ablation study with the baselines and the evaluation results on MOT15 [5] and MOT17 [6] benchmarks with state-of-the-art MOT methods. The ablation study proves that GMPHD-OGM (w/ SIOA) is more efficient to solve the defined problems than the given baseline methods such as GMPHD-HDA and GMPHD-OGM (w/ IOU). GMPHD-OGM achieves the best MOTA scores in MOT15 and MOT17 datasets, respectively, in comparison with the PHD filter based online trackers [16], [18], [22], [38], [41]–[43]. Finally, by the comprehensive

evaluation, we conclude the proposed tracker (GMPHD-OGM) against state-of-the-art algorithms including online with DNN even including offline, GMPHD-OGM shows the competitive value in “MOTA versus speed”. As a future work, we will develop an efficient real-time tracker even with the number of objects over hundred, simultaneously, achieving the state-of-the-art level tracking accuracy.

## REFERENCES

- [1] R. P. S. Mahler, “Multitarget Bayes Filtering via First-Order Multitarget Moments,” *IEEE Trans. Aerosp. Electron. Syst.*, vol. 39, no. 4, pp. 1152–1178, Oct. 2003.
- [2] B.-N. Vo, S. Singh, and A. Doucet, “Sequential Monte Carlo implementation of the PHD filter for multi-target tracking,” in Proc. Int. Conf. Information Fusion (ICIF), pp. 792–799, Jul. 2003.
- [3] B.-N. Vo and W.-K. Ma, “The Gaussian mixture probability hypothesis density filter,” *IEEE Trans. Signal Processing*, vol. 54, no. 11, pp. 4091–4104, Oct. 2006.
- [4] B.-T. Vo, “Random finite sets in multi-object filtering,” PhD Thesis, School of Electrical, Electronic and Computer Engineering, The Univ. of Western Australia, Perth, Australia, 2008.
- [5] A. Milan, L. Leal-Taixé, I. Reid, S. Roth, and K. Schindler, “MOTChallenge 2015: Towards a benchmark for multi-target tracking,” [Online]. Available: <http://arxiv.org/abs/1504.01942>, Apr. 2015.
- [6] L. Leal-Taixé, A. Milan, I. Reid, S. Roth, and K. Schindler, “MOT16: A Benchmark for Multi-Object Tracking,” [Online]. Available: <https://arxiv.org/abs/1603.00831>, May. 2016.
- [7] A. Geiger, P. Lenz, and R. Urtasun, “Are we ready for Autonomous Driving? The KITTI Vision Benchmark Suite,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 3354–3361, Jun. 2012.
- [8] P. Dollár, R. Appel, S. Belongie, and P. Perona, “Fast feature pyramids for object detection,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 36, no. 8, pp. 1532–1545, Aug. 2014.
- [9] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan, “Object detection with discriminatively trained part-based models,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 32, no. 9, pp. 1627–1645, Sep. 2010.
- [10] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in Proc. 28th Int. Conf. Neural Inf. Process. Syst. (NIPS), Montréal, QC, Canada, vol. 1, pp. 91–99, Dec. 2015.
- [11] F. Yang, W. Choi, and Y. Lin, “Exploit all the layers: Fast and accurate cnn object detector with scale dependent pooling and cascaded rejection classifiers,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 2129–2137, Jun. 2016.
- [12] L. Wang, Y. Lu, H. Wang, Y. Zheng, H. Ye, and X. Xue, “Evolving boxes for fast vehicle detection,” in Proc. IEEE Conf. Multi. Expo (ICME), pp. 1135–1140, Jul. 2017.
- [13] S. Murray, “Real-Time Multiple Object Tracking - A Study on the Importance of Speed,” [Online]. Available: <https://arxiv.org/abs/1709.03572>, Oct. 2017.
- [14] S.-H. Bae and K.-J. Yoon, “Confidence-Based Data Association and Discriminative Deep Appearance Learning for Robust Online Multi-Object Tracking,” *IEEE Trans. Pattern Anal. Mach. Intell.*, vol. 40, no. 3, pp. 595–610, Mar. 2018.
- [15] Y.-C. Yoon, A. Boragule, Y. Song, K. Yoon, and M. Jeon, “Online Multi-Object Tracking with Historical Appearance Matching and Scene Adaptive Detection Filtering,” in Proc. IEEE Int. Workshop Traffic Street Survill. Safety Secur. (AVSS), pp. 1–6, Nov. 2018.
- [16] Y. Song and M. Jeon, “Online Multiple Object Tracking with the Hierarchically Adopted GM-PHD Filter using Motion and Appearance,” in Proc. IEEE Int. Conf. Consumer Electronics-Asia (ICCE-Asia), pp. 1–4, Oct. 2016.
- [17] Y. Song, Y.-C. Yoon, K. Yoon, and M. Jeon, “Online and Real-Time Tracking with the GMPHD Filter using Group Management and Relative Motion Analysis,” in Proc. IEEE Int. Workshop Traffic Street Survill. Safety Secur. (AVSS), pp. 1–6, Nov. 2018.
- [18] Z. Fu, P. Feng, F. Angelini, J. Chambers, and S. M. Naqvi, “Particle phd filter based multiple human tracking using online group-structured dictionary learning,” *IEEE Access*, vol. 6, pp. 14764–14778, Mar. 2018.
- [19] Y. Xiang, A. Alahi, and S. Savarese, “Learning to Track: Online Multi-Object Tracking by Decision Making,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), pp. 4705–4713, Dec. 2015.

TABLE II: Quantitative evaluation results on MOT15 training dataset. The proposed method namely GMPHD-OGM (w/ SIOA) is compared to two baseline methods GMPHD-HDA and GMPHD-OGM (w/ IOU). GMPHD-HDA employs the GMPHD filtering with hierarchical data association (HDA). GMPHD-OGM is equal to GMPHD-HDA with the proposed occlusion group management (OGM). The IOU and SIOA metrics are used for “Track Merging” in GMPHD-OGM (w/ IOU) and (w/ SIOA), respectively. The optimal values of the merging threshold  $\sigma_m$  are underlined and the best scores are in **bold** in terms of the CLEAR-MOT metrics.

Tracker	$\sigma_m$	MOTA↑	MOTP↑	MT↑	ML↓	FP↓	FN↓	IDS↓	Frag↓	Speed↑
GMPHD-HDA	n/a	34.8 %	72.3 %	14.4 %	47.8 %	4,042	21,646	338	572	212.4 fps
GMPHD-OGM (w/ IOU)	0.2	34.5 %	<b>72.4 %</b>	13.4 %	49.8 %	3,594	22,226	285	550	201.1 fps
	0.3	35.6 %	72.3 %	14.0 %	48.6 %	3,537	21,901	<b>278</b>	548	<b>228.4 fps</b>
	0.4	35.3 %	<b>72.4 %</b>	14.0 %	48.2 %	3,667	21,865	291	562	201.9 fps
	0.5	34.7 %	72.3 %	14.4 %	47.8 %	4,044	21,661	340	577	205.3 fps
	0.3	34.5 %	72.3 %	14.2 %	49.2 %	<b>3,496</b>	22,336	297	<b>559</b>	216.3 fps
GMPHD-OGM (w/ SIOA)	0.4	35.4 %	<b>72.4 %</b>	14.6 %	48.4 %	3,556	21,930	284	<b>540</b>	216.3 fps
	<u>0.5</u>	<b>35.8 %</b>	72.2 %	15.0 %	47.6 %	3,569	21,758	295	545	221.0 fps
	0.6	35.5 %	72.3 %	<b>15.6 %</b>	<b>47.2 %</b>	3,702	21,724	312	556	221.9 fps
	0.7	34.7 %	72.2 %	<b>15.6 %</b>	<b>47.2 %</b>	4,159	<b>21,519</b>	368	567	202.9 fps

TABLE III: Quantitative evaluation results on MOT15 test dataset. The proposed method is compared to state-of-the-art in terms of the CLEAR-MOT metrics. For each mode, i.e, online and offline, the first and the second best scores are highlighted in **red** and **blue** in terms of each metric.

Mode	Tracker	DNN	MOTA↑	MOTP↑	MT↑	ML↓	FP↓	FN↓	IDS↓	Frag↓	Speed↑
Online	CDA_DDAL [14]	O	<b>32.8 %</b>	70.7 %	9.7 %	42.2 %	4,983	35,690	614	1,583	2.3 fps
	HAM_SADF [15]	O	28.6 %	71.1 %	10.0 %	44.0 %	7,485	35,910	<b>460</b>	<b>1,038</b>	18.7 fps
	Proposed*	X	<b>30.7 %</b>	<b>71.6 %</b>	<b>11.5 %</b>	<b>38.1 %</b>	6,502	<b>35,030</b>	1,034	1,351	<b>169.5 fps</b>
	PHD_GSDL [18]	X	30.5 %	71.2 %	7.6 %	41.2 %	6,534	35,284	879	2,208	8.2 fps
	MDP [19]	X	30.3 %	<b>71.3 %</b>	<b>14.0 %</b>	<b>38.4 %</b>	9,717	<b>32,422</b>	680	1,500	1.1 fps
	TBSS [20]	X	29.2 %	<b>71.3 %</b>	6.8 %	43.8 %	<b>6,068</b>	36,779	649	1,508	11.5 fps
	SCEA [21]	X	29.1 %	71.1 %	8.9 %	47.3 %	<b>6,060</b>	36,912	604	1,182	6.8 fps
	EAMTT [22]	X	22.3 %	69.6 %	5.4 %	52.7 %	7,924	38,982	833	1,485	12.2 fps
	RMOT [23]	X	18.6 %	69.6 %	5.3 %	53.3 %	12,473	36,835	684	1,282	7.9 fps
	GMPHD_HDA [16]	X	18.5 %	70.9 %	3.9 %	55.3 %	7,864	41,766	<b>459</b>	1,266	19.8 fps
Offline	GSCR [24]	X	15.8 %	69.4 %	1.8 %	61.0 %	7,597	43,633	514	<b>1,010</b>	<b>28.1 fps</b>
	TC_ODAL [25]	X	15.1 %	70.5 %	3.2 %	55.8 %	12,970	38,538	637	1,716	1.7 fps
	MHT_DAM [28]	O	<b>32.4 %</b>	<b>71.8 %</b>	<b>16.0 %</b>	<b>43.8 %</b>	9,064	<b>32,060</b>	<b>435</b>	826	0.7 fps
	CNNTCM [26]	O	29.6 %	<b>71.8 %</b>	11.2 %	44.0 %	7,786	34,733	712	943	1.7 fps
	SiameseCNN [27]	O	29.0 %	71.2 %	8.5 %	48.4 %	<b>5,160</b>	37,798	639	1,316	<b>52.8 fps</b>
	NOMT [29]	X	<b>33.7 %</b>	<b>71.9 %</b>	<b>12.2 %</b>	44.6 %	7,762	<b>32,547</b>	442	<b>823</b>	11.5 fps
	ELP [30]	X	25.0 %	71.2 %	7.5 %	<b>43.8 %</b>	7,345	37,344	1,396	1,804	5.7 fps
	JPDA_m [31]	X	23.8 %	68.2 %	5.0 %	58.1 %	<b>6,373</b>	70,084	<b>365</b>	869	32.6 fps

\* The final proposed model is GMPHD-OGM (w/ SIOA).

- [20] X. Zhou, P. Jiang, Z. Wei, H. Dong, and F. Wang, “Online Multi-Object Tracking with Structural Invariance Constraint,” in Proc. Brit. Mach. Vis. Conf. (BMVC) pp. 1–13, Sep. 2018.
- [21] J. Yoon, C.-R. Lee, M.-H. Yang, K.-J. Yoon, “Online Multi-object Tracking via Structural Constraint Event Aggregation,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 1392–1400, Jun. 2016.
- [22] R. Sanchez-Matilla, F. Poiesi, and A. Cavallaro, “Multi-target tracking with strong and weak detections,” in Proc. Eur. Conf. Comput. Vis. Workshops (ECCVW), pp. 84–99, Oct. 2016.
- [23] J. Yoon, M.-H. Yang, J. Lim, and K.-J. Yoon, “Bayesian Multi-Object Tracking Using Motion Context from Multiple Objects,” in Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV), pp. 1–13, Jan. 2015.
- [24] L. Fagot-Bouquet, R. Audigier, Y. Dhome, and F. Lerasle, “Online multi-person tracking based on global sparse collaborative representations”, in Proc. IEEE Conf. Image Processing (ICIP), pp. 2414–2418, Sep. 2015.
- [25] S.-H. Bae and K.-J. Yoon, “Robust Online Multi-Object Tracking based on Tracklet Confidence and Online Discriminative Appearance Learning,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 1218–1225, Jun. 2015.
- [26] B. Wang, L. Wang, B. Shuai, Z. Zuo, T. Liu, K. L. Chan, and G. Wang, “Joint Learning of Convolutional Neural Networks and Temporally Constrained Metrics for Tracklet Association,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW), pp. 386–393, Jun. 2016.
- [27] L. Leal-Taixé, C. Canton-Ferrer, and K. Schindler, “Learning by Tracking: Siamese CNN for Robust Target Association,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW), pp. 33–40, Jun. 2016.
- [28] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg, “Multiple Hypothesis Tracking Revisited,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), pp. 4696–4704, Dec. 2015.
- [29] W. Choi, “Near-Online Multi-target Tracking with Aggregated Local Flow Descriptor,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), pp. 3029–3037, Dec. 2015.
- [30] N. McLaughlin, J. M. D. Rincon, and P. Miller, “Enhancing Linear Programming with Motion Modeling for Multi-target Tracking,” in Proc. IEEE Winter Conf. Appl. Comput. Vis. (WACV), pp. 71–77, Jan. 2015.
- [31] S. H. Rezatofighi, A. Milan, Z. Zhang, Q. Shi, A. Dick, and I. Reid, “Joint Probabilistic Data Association Revisited,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), pp. 3047–3055, Dec. 2015.
- [32] L. Leal-Taixé, M. Fenzi, A. Kuznetsova, B. Rosenhahn, and S. Savarese, “Learning an image-based motion context for multiple people tracking,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 3542–

TABLE IV: Quantitative evaluation results on MOT17 test dataset. The proposed method is compared to state-of-the-art in terms of the CLEAR-MOT metrics. For each mode, i.e, online and offline, the first and the second best scores are highlighted in **red** and **blue** in terms of each metric.

Mode	Tracker	DNN	MOTA↑	MOTP↑	MT↑	ML↓	FP↓	FN↓	IDS↓	Frag↓	Speed↑
Online	MOTDT [37]	O	<b>50.9 %</b>	76.6 %	17.5 %	35.7 %	24,069	<b>250,768</b>	2,474	5,317	18.3 fps
	HAM_SADF [15]	O	48.3 %	<b>77.2 %</b>	17.1 %	41.7 %	20,967	269,038	<b>1,871</b>	<b>3,020</b>	5.0 fps
	DMAN [39]	O	48.2 %	75.7 %	<b>19.3 %</b>	38.3 %	26,218	263,608	<b>2,194</b>	5,378	0.3 fps
	Proposed*	X	<b>49.9 %</b>	77.0 %	<b>19.7 %</b>	38.0 %	24,024	255,277	3,125	<b>3,540</b>	<b>30.7 fps</b>
	MTDF [38]	X	49.6 %	75.5 %	18.9 %	<b>33.1 %</b>	37,124	<b>241,768</b>	5,567	9,260	1.2 fps
	AM_ADM [40]	X	48.1 %	76.7 %	13.4 %	37.7 %	25,061	265,495	2,214	5,027	5.7 fps
	PHD_GSDL [18]	X	48.0 %	<b>77.2 %</b>	17.1 %	<b>35.6 %</b>	23,199	265,954	3,998	8,886	6.7 fps
	GMPHD_HDA [16]	X	43.7 %	76.5 %	11.7 %	43.0 %	25,935	287,758	3,838	5,056	9.2 fps
	EAMTT [22]	X	42.6 %	76.0 %	12.7 %	42.7 %	<b>20,711</b>	288,474	4,488	5,720	12.0 fps
	GMPHD_N1Tr [41]	X	42.1 %	<b>77.7 %</b>	11.9 %	42.7 %	<b>18,214</b>	297,646	10,698	10,864	9.9 fps
Offline	GMPHD_KCF [42]	X	39.6 %	74.5 %	8.8 %	43.3 %	50,903	284,228	5,811	7,414	3.3 fps
	GM_PHD [43]	X	36.4 %	76.2 %	4.1 %	57.3 %	23,723	330,767	4,607	11,317	<b>38.4 fps</b>
	MHT_DAM [28]	O	50.7 %	<b>77.5 %</b>	20.8 %	36.9 %	22,875	252,889	2,314	2,865	0.9 fps
	MHT_bLSTM [50]	O	47.5 %	<b>77.5 %</b>	18.2 %	41.7 %	25,981	268,042	2,069	3,124	1.9 fps
	eHAF [44]	X	<b>51.8 %</b>	77.0 %	<b>23.4 %</b>	37.9 %	33,212	<b>236,772</b>	1,834	<b>2,739</b>	0.7 fps
	FWT [45]	X	<b>51.3 %</b>	77.0 %	21.4 %	<b>35.2 %</b>	24,101	247,921	2,648	4,279	0.2 fps
	JCC [46]	X	51.2 %	75.9 %	20.9 %	37.0 %	25,937	247,822	<b>1,802</b>	2,984	1.8 fps
	TLMHT [47]	X	50.6 %	<b>77.6 %</b>	17.6 %	43.4 %	<b>22,213</b>	255,030	<b>1,407</b>	<b>2,079</b>	<b>2.6 fps</b>
	EDMT [48]	X	50.0 %	77.3 %	<b>21.6 %</b>	<b>36.3 %</b>	32,279	<b>247,297</b>	2,264	3,260	0.6 fps
	IOU [49]	X	45.5 %	76.9 %	15.7 %	40.5 %	<b>19,993</b>	281,643	5,988	7,404	<b>1,522.9 fps</b>

\* The final proposed model is GMPHD-OGM (w/ SIOA).

- 3549, Jun. 2014.
- [33] A. Milan, L. Leal-Taixé, K. Schindler, and I. Reid, “Joint Tracking and Segmentation of Multiple Targets,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 5397–5406, Jun. 2015.
- [34] A. Milan, S. Roth, and K. Schindler, “Continuous Energy Minimization for Multitarget Tracking,” IEEE Trans. Pattern Anal. Mach. Intell., vol. 36, no. 1, pp. 58–72, Aug. 2014.
- [35] C. Dicle, O. I. Camps, and M. Sznaier, “The Way They Move: Tracking Targets with Similar Appearance,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), pp. 2304–2311, Dec. 2013.
- [36] H. Pirsiavash, D. Ramanan, and C. C. Fowlkes, “Globally-Optimal Greedy Algorithms for Tracking a Variable Number of Objects,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 1201–1208, Jun. 2011.
- [37] C. Long, A. Haizhou, Z. Zijie and S. Chong, “Real-time Multiple People Tracking with Deeply Learned Candidate Selection and Person Re-identification,” in Proc. IEEE Conf. Multi. Expo (ICME), pp. 1–6, Oct. 2018.
- [38] Z. Fu, F. Angelini, J. Chambers, and S. M. Naqvi, “Multi-Level Cooperative Fusion of GM-PHD Filters for Online Multiple Human Tracking,” IEEE trans. Multimedia. (Early Access), Mar. 2019.
- [39] J. Zhu, H. Yang, N. Liu, M. Kim, W. Zhang, and M.-H. Yang, “Online Multi-Object Tracking with Dual Matching Attention Networks,” in Proc. Eur. Conf. Comput. Vis. (ECCV), pp. 366–382, Feb. 2019.
- [40] S.-H. Lee, M.-Y. Kim, and S.-H. Bae, “Learning Discriminative Appearance Models for Online Multi-Object Tracking with Appearance Discriminability Measures,” IEEE Access, vol. 6, pp. 67316–67328, Nov. 2018.
- [41] N. L. Baisa and A. Wallace, “Development of a N-type GM-PHD filter for multiple target, multiple type visual tracking.” Journal of Visual Communication and Image Representation, vol. 59, pp. 257–271, 2019.
- [42] T. Kutschbach, E. Bochinski, V. Eiselein, and T. Sikora, “Sequential Sensor Fusion Combining Probability Hypothesis Density and Kernelized Correlation Filters for Multi-Object Tracking in Video Data,” in Proc. IEEE Int. Workshop Traffic Street Survill. Safety Secur. (AVSS), pp. 1–6, Sep. 2017.
- [43] V. Eiselein, D. Arp, M. Pätzold, and T. Sikora, “Real-time Multi-Human Tracking using a Probability Hypothesis Density Filter and multiple detectors,” in Proc. IEEE Int. Conf. Adv. Video Signal Based Survill. (AVSS), pp. 325–330, Sep. 2012.
- [44] H. Sheng, Y. Zhang, J. Chen, Z. Xiong, and J. Zhang, “Heterogeneous Association Graph Fusion for Target Association in Multiple Object Tracking,” IEEE Trans. Circuits Syst. Video Technol., (Early Access), Nov. 2018.
- [45] R. Henschel, L. Leal-Taixé, D. Cremers, and B. Rosenhahn, “Fusion of Head and Full-Body Detectors for Multi-Object Tracking,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW), pp. 1428–1437, Jun. 2018.
- [46] M. Keuper, S. Tang, B. Andres, T. Brox, and B. Schiele, “Motion Segmentation & Multiple Object Tracking by Correlation Co-Clustering,” IEEE Trans. Pattern Anal. Mach. Intell., (Early Access), Oct. 2018.
- [47] H. Sheng, J. Chen, Y. Zhang, W. Ke, Z. Xiong, and J. Yu, “Iterative Multiple Hypothesis Tracking with Tracklet-level Association,” IEEE Trans. Circuits Syst. Video Technol., (Early Access), Nov. 2018.
- [48] J. Chen, H. Sheng, Y. Zhang, and Z. Xiong, “Enhancing Detection Model for Multiple Hypothesis Tracking,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. Workshops (CVPRW), pp. 18–27, Jul. 2017.
- [49] E. Bochinski, V. Eiselein, and T. Sikora, “High-Speed Tracking-by-Detection Without Using Image Information,” in Proc. IEEE Int. Workshop Traffic Street Survill. Safety Secur. (AVSS), pp. 1–6, Sep. 2017.
- [50] C. Kim, F. Li, and J. Rehg, “Multi-object Tracking with Neural Gating Using Bilinear LSTM,” in Proc. Eur. Conf. Comput. Vis. (ECCV), pp. 200–215, Sep. 2018.
- [51] S. Chopra, R. Hadsell, Y. LeCun, “Learning a Similarity Metric Discriminatively,” in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), pp. 539–546, Jun. 2005.
- [52] L. Zhao, X. Li, J. Wang, and Y. Zhuang, “Deeply-learned part-aligned representations for person re-identification,” in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), pp. 3239–3248, Oct. 2017.
- [53] J. F. Henriques, R. Caseiro, P. Martins, and J. Batista, “High-Speed Tracking with Kernelized Correlation Filters,” IEEE Trans. Pattern Anal. Mach. Intell., vol. 37, no. 3, pp. 583–596, Mar. 2015.
- [54] K. Bernardin and R. Stiefelhagen, “Evaluating Multiple Object Tracking Performance: The CLEAR MOT Metrics,” EURASIP Journal on Image and Video Processing, vol. 2008 no. 1, pp. 1–10, May 2008.
- [55] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman, “The Pascal Visual Object Classes Challenge: A Retrospective,” International Journal of Computer Vision, vol. 111 no. 1, pp. 98–136, Jan. 2015.
- [56] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, “ImageNet Large Scale Visual Recognition Challenge,” International Journal of Computer Vision, vol. 115 no. 3, pp. 211–252, Dec. 2015.

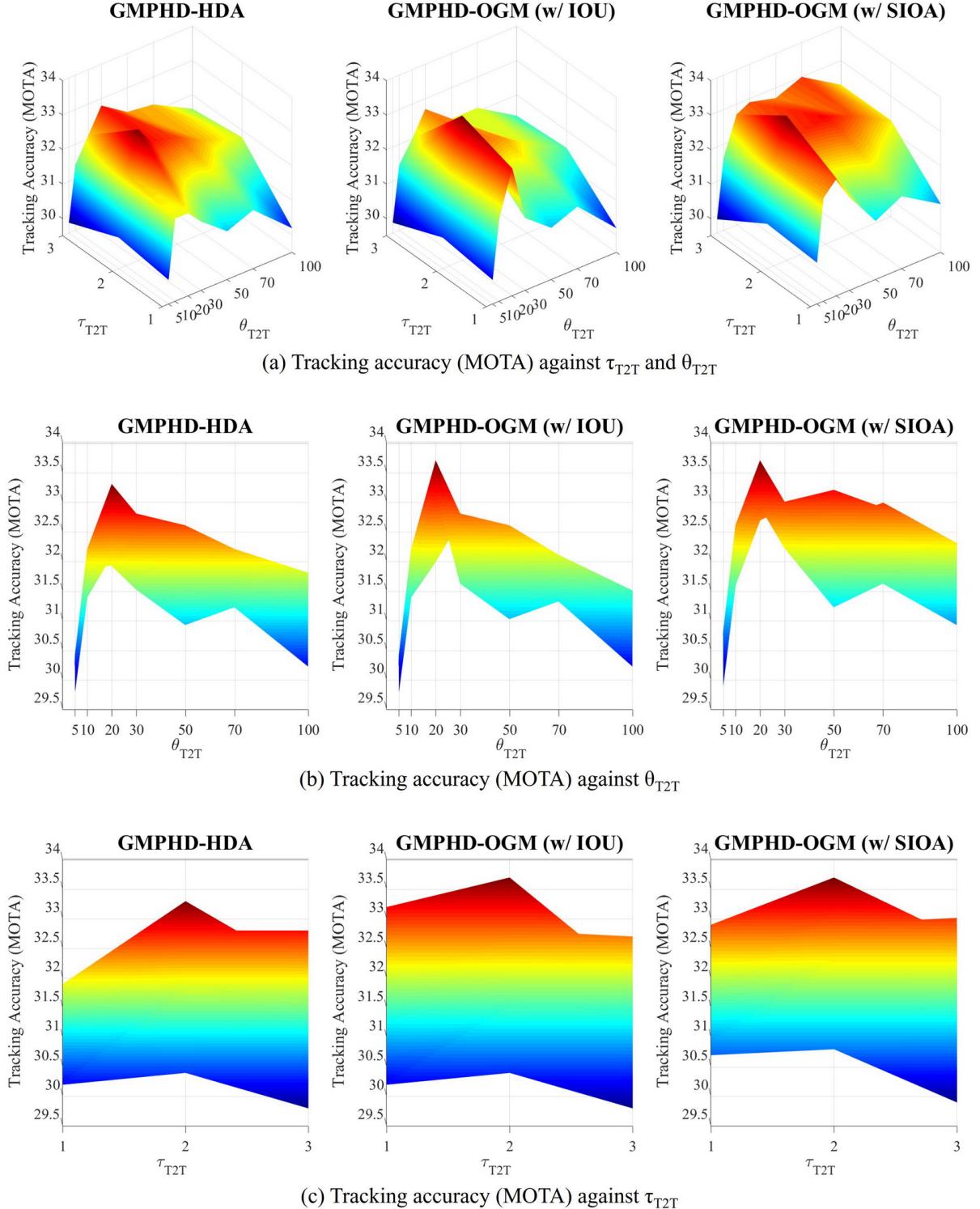


Fig. 7: Ablation study with the two baseline methods, i.e., GMPHD-HDA and GMPHD-OGM (with IOU). The final proposed method is GMPHD-OGM (with SIOA). Three graphs indicates the MOTA scores' distributions against (a) the minimum track length for T2TA ( $\tau_{T2T}$ ) and the maximum frame interval for T2TA ( $\theta_{T2T}$ ), (b)  $\tau_{T2T}$ , and (c)  $\theta_{T2T}$ . GMPHD-OGM (with SIOA) shows overall improvement in upper and lower bound of MOTA score.

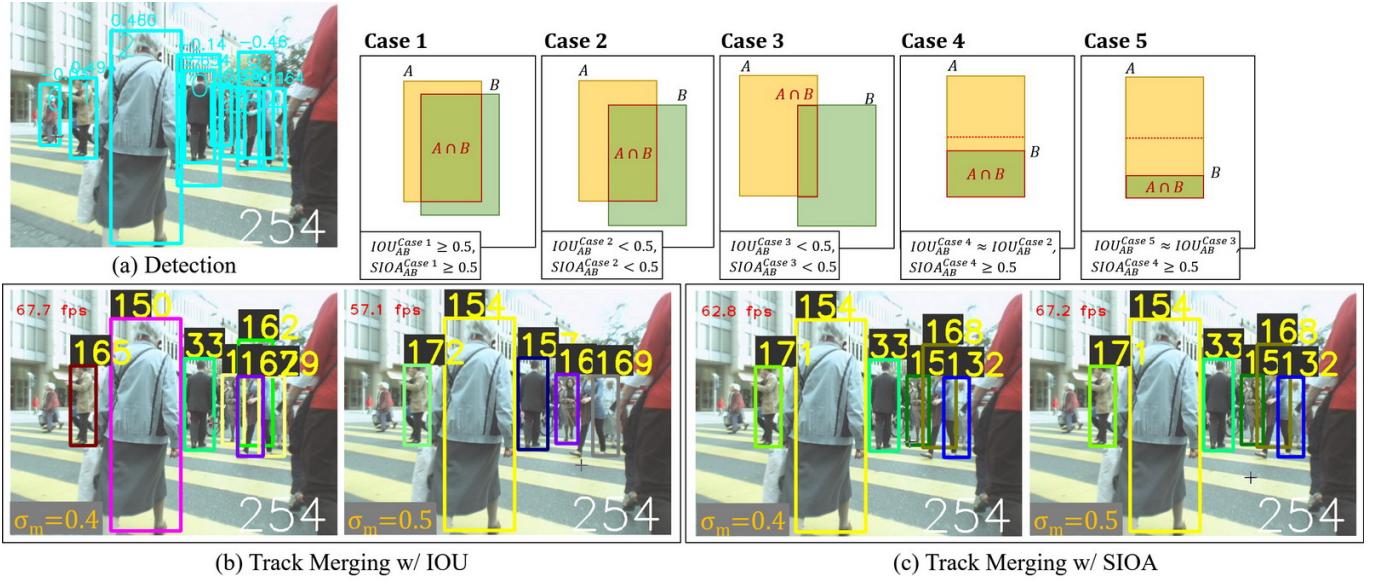


Fig. 8: Case study about “Track Merging” with the qualitative results on MOT17-05-DPM training sequence at frame 254. For the same detection results, the overlapping ratios between the occluded objects are measured with  $IOU \gtrsim 0.4$  and  $SIOA \gtrsim 0.6$ . Under the different merging threshold  $\sigma_m$  values 0.4 and 0.5, the IOU metric is more sensitive than the SIOA metric. The SIOA metric is more robust to merge size variant false positive bounding boxes than IOU metric.

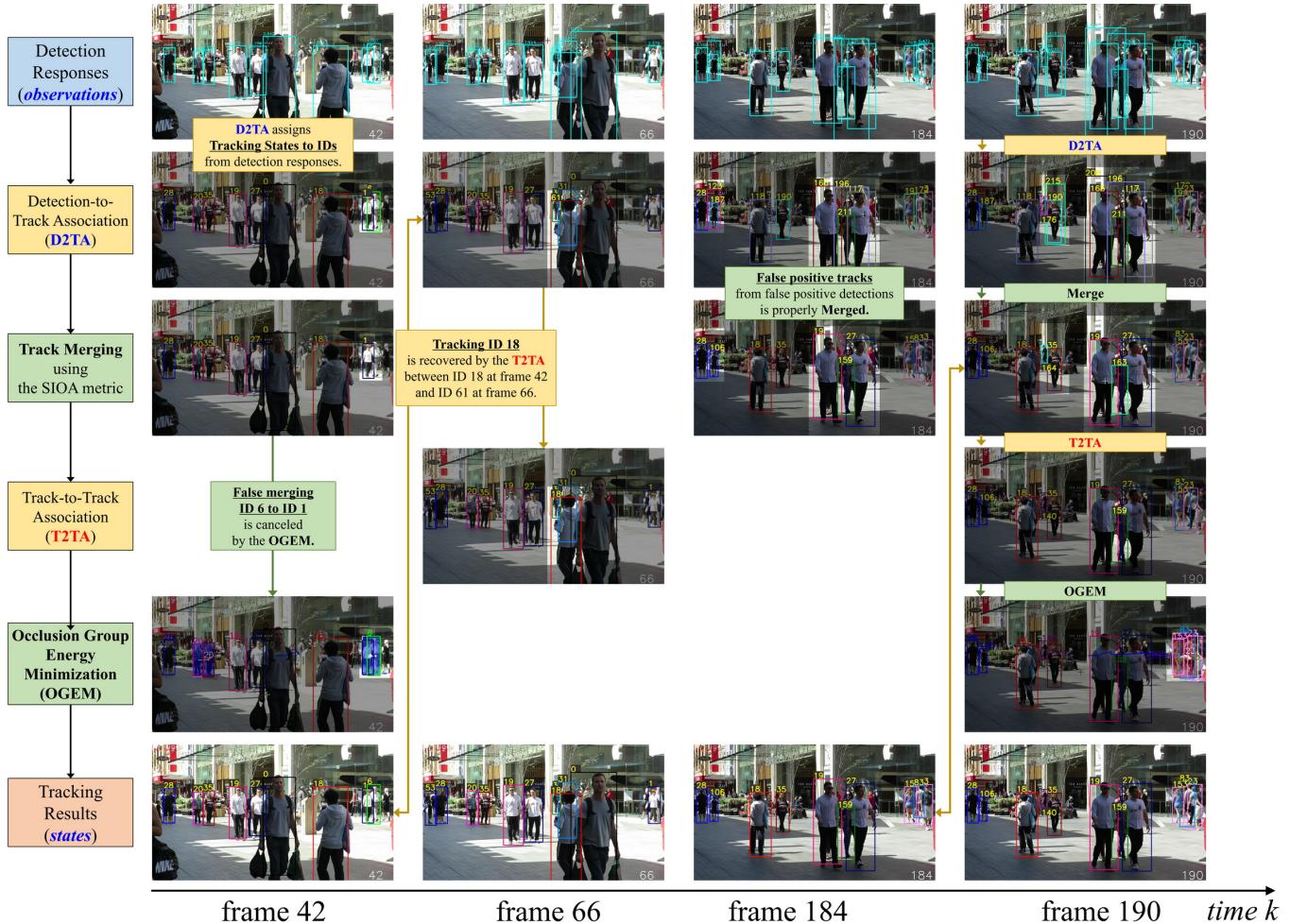


Fig. 9: Illustration of the proposed multi-object tracking process with the qualitative results on MOT17-08-DPM test sequence. The whole process consists of four components which are D2TA, Merge, T2TA, and OGEM. Qualitative tracking results at frame 42, 66, 184, and 190 present that all components are complementary to each other with handling tracking problems.