

# SIMPLE ONLINE AND REALTIME TRACKING WITH OCCLUSION HANDLING

Mohammad Hossein Nasseri, Hadi Moradi, Reshad Hosseini

Mohammadreza Babaei

School of Electrical and Computer Engineering  
University of Tehran  
Tehran, Iran

Independent researcher

## ABSTRACT

Multiple object tracking is a challenging problem in computer vision due to difficulty in dealing with motion prediction, occlusion handling, and object re-identification. Many recent algorithms use motion and appearance cues to overcome these challenges. But using appearance cues increases the computation cost notably and therefore the speed of the algorithm decreases significantly which makes them inappropriate for online applications. In contrast, there are algorithms that only use motion cues to increase speed, especially for online applications. But these algorithms cannot handle occlusions and re-identify lost objects. In this paper, a novel online multiple object tracking algorithm is presented that only uses geometric cues of objects to tackle the occlusion and re-identification challenges simultaneously. As a result, it decreases the identity switch and fragmentation metrics. Experimental results show that the proposed algorithm could decrease identity switch by 40% and fragmentation by 28% compared to the state of the art online tracking algorithms. The code is also publicly available<sup>1</sup>.

**Index Terms**— Multiple Object Tracking, Occlusion Handling, Target Re-identification, Confidence-Based

## 1. INTRODUCTION

Multiple Object Tracking (MOT) is a challenging problem in computer vision that has a wide variety of applications [1, 2]. In many applications such as autonomous driving, robot navigation, and visual surveillance real-time tracking is highly demanded. So, online MOT methods that only use current and previous frames can be used for these applications. Most recent MOT methods use tracking by detection pipeline which at first targets are detected by a detection algorithm and the results are passed to a data association algorithm to find trajectories [3, 4, 5]. Obviously, the quality of the detection algorithm significantly affects the result of the tracking algorithm. These algorithms utilize motion and appearance information in their data association part. However using the appearance feature imposes high computation cost and hence decreases

the speed of algorithms. As a result, the maximum frame rate of the recent online MOT algorithms in MOT17<sup>2</sup> dataset with a high MOTA metric is below 5 Frame Per Second (FPS) even using high-performance hardware [3, 5, 6]. On the other hand, algorithms that only use motion cues cannot overcome occlusion and object re-identification challenges. For example, the Simple Online and Real-time Tracking (SORT) [7] algorithm that only uses motion cues, achieves high FPS. However the number of ID switch and fragmentation metrics [8] are high that shows its low performance in occlusion handling and target re-identification.

In this paper, a new algorithm is proposed which only uses the location and size of detection bounding boxes to tackle occlusion without imposing high computation cost. Additionally, it can predict occlusions and re-identify lost targets. As a result, it decreases the ID switch and fragmentation metrics significantly.

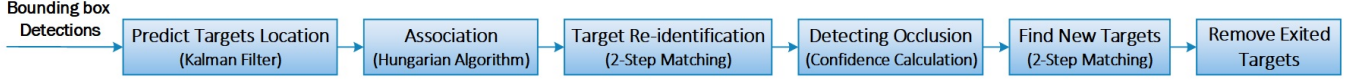
The rest of the paper is as follows: Section 2 provides related work in the area of online multiple object tracking. Section 3 describes the proposed confidence based occlusion handling and target re-identification. In Section 4, the experiments are explained and the results are presented. Finally, in Section 5, we draw our conclusion.

## 2. RELATED WORK

The majority of works in recent years use the tracking by detection paradigm. In this paradigm, first, a detection algorithm locates targets. For example, in MOT17 benchmark, DPM [9], FRCNN [10], and SDP [11] algorithms, which are based on Convolutional Neural Network (CNN), are used for target's bounding box detection. The second step is extracting features from bounding box detections, such as the motion and appearance features that are widely used. For extracting motion features, the trajectory of each target should be predicted from the location and size of its bounding box. In some algorithms, filter-based methods such as Kalman filter [12] with constant velocity assumption is used to predict the bounding box of targets [7, 13]. In some other works, optical flow [14] and LSTM [15] are employed. In [16], the

<sup>1</sup>[https://github.com/mhnasseri/sort\\_oh](https://github.com/mhnasseri/sort_oh)

<sup>2</sup><https://motchallenge.net/>



**Fig. 1:** The block diagram of the proposed algorithm

motion of targets in crowded scenarios are predicted by utilizing more complex methods such as the social force model. Recently reinforcement learning is proposed for the prediction of target location in successive frames [17].

For appearance features, the early works used hand-crafted features such as color histograms [18]. Recent works learn features from CNN networks including Siamese network [19], correlation filters [20] and auto-encoders [21]. However, the computation cost of CNN-based methods is high that gives rise to a significant decrease in the speed of algorithms, even on high-performance hardware. Consequently, there are hardly suitable for online applications.

The third step is calculating the similarity between detections and targets. For calculating affinity from predicted motion, the location and size of bounding boxes are considered and metrics such as Intersection over Union (IoU) is calculated [7, 13]. Also for obtaining similarity between object pairs based on an appearance feature, cosine similarity [3] and other CNN based methods [22] as well as LSTM variants [23] are used.

The final step is associating detections to targets. For this task different approaches are used such as Hungarian algorithm [7, 13], dynamic programming [24], and multiple hypothesis tracking [25]. Recently reinforcement learning [26] and graph-based methods [27] are employed for the association task. Alongside the tracking by detection paradigm, some works combine detection and tracking tasks. For example, in [4], a regressor based detector is used for tracking objects. Besides, There are end-to-end methods for object tracking that perform both detection and tracking at the same time [28, 29].

The most referenced algorithms for online object tracking are SORT [7] and Deep SORT[13]. In the Deep SORT algorithm, the appearance cue is added to the SORT algorithm to improve its weakness in target re-identification and imposes higher computation cost. This makes the Deep SORT algorithm inappropriate for real-time applications and decreases the speed of the algorithm dramatically. In this paper, an algorithm based on the SORT algorithm is proposed that deals with occlusion handling and target re-identification efficiently without using appearance cues and losing the speed of the algorithm. In the proposed algorithm, occlusion handling and target re-identification are tackled using only the location and size of bounding boxes and a comparable result to the Deep Sort algorithm is achieved at a very lower computation cost.

### 3. APPROACH

The proposed algorithm is inspired from [7], which is based on tracking-by-detection paradigm. The main steps of the proposed algorithm is shown in Fig. 1. The detected bounding boxes are given as input to the algorithm. These bounding boxes at frame  $t$  are denoted as  $\mathbf{D}_t = \{D_t^1, D_t^2, \dots, D_t^{K_t}\}$  where  $K_t$  is the number of detections at frame  $t$ . Also the targets at frame  $t - 1$  are indicated as  $\mathbf{T}_{t-1} = \{T_{t-1}^1, T_{t-1}^2, \dots, T_{t-1}^{P_{t-1}}\}$  where  $P_{t-1}$  is the number of targets at frame  $t - 1$ . The prediction of these targets at current frame are denoted as  $\hat{\mathbf{T}}_t$  which is another input to the association step of the proposed algorithm. After the association step, most of detections are matched with targets and considered as new observation for correcting estimation and a few of them are not matched. The matched detections at frame  $t$  are indicated by  $\bar{\mathbf{D}}_t$  and unmatched detections are depicted by  $\tilde{\mathbf{D}}_t$ . After the association step, some targets are matched with detections,  $\bar{\mathbf{T}}_t$ , a few of them are marked as occluded, that are  $\mathbf{O}_t$  and a few of them are remained unmatched,  $\tilde{\mathbf{T}}_t$ . At the end, some new targets could be created and added to  $\hat{\mathbf{T}}_t$  or some targets may be discarded and removed from  $\hat{\mathbf{T}}_t$ .

The first step of the algorithm is to model objects and their motions to predict their location in future frames. Like [7], the state of targets are modeled as:

$$x = [u, v, s, r, \dot{u}, \dot{v}, \dot{s}] \quad (1)$$

In the above equation,  $u$  and  $v$  represent the center,  $s$  represents the area, and  $r$  represents the aspect ratio of a bounding box. The  $\dot{u}, \dot{v}$  and  $\dot{s}$  are the rate of change for the corresponding variables. A Kalman filter with a constant velocity model is utilized to model the motion of objects. This filter predicts the location and size of each target in the next frames.

For associating detections to targets, the Intersection over Union (IoU) is calculated. This metric is equal to the intersection of the bounding box of detection with a bounding box of each target divided by the union area of these two bounding boxes:

$$IoU = \frac{I(bb_D, bb_T)}{A(bb_D) + A(bb_T) - I(bb_D, bb_T)} \quad (2)$$

in which  $bb_D$  and  $bb_T$  stand for the bounding box of a detection and the bounding box of a target, respectively. Also,  $A$  represents the area of a bounding box and  $I$  represents the intersection area of two bounding boxes. In the association step, the Hungarian algorithm uses IoU values as a similarity metric to assign detection boxes to the targets.

### 3.1. Target's Confidence

The main idea of the proposed algorithm is introducing a confidence value for each target during tracking. Basically, in online tracking algorithms, the targets matched with detections are marked as live targets and the rest are kept as reserved targets or are discarded. In the proposed algorithm, a confidence measure is defined that plays a key role in determining if a target should be marked as occluded. This confidence measure is a representation of the temporal and spatial characteristics of a target. In the case of temporal characteristics, if a target is observed for several successive frames, the probability of its existence increases. In contrast, if the target is not detected for several consecutive frames, its certainty decreases. In the case of spatial characteristics of a target, if the bounding box of a target is larger, it means that the target covers larger area of the image and the confidence in such a target is higher. It may be noted that for targets with similar sizes, the bigger bounding box is correlated to being closer to the camera. Having discussed all relevant factors, the confidence of  $i^{th}$  target ( $C_i$ ) is defined as:

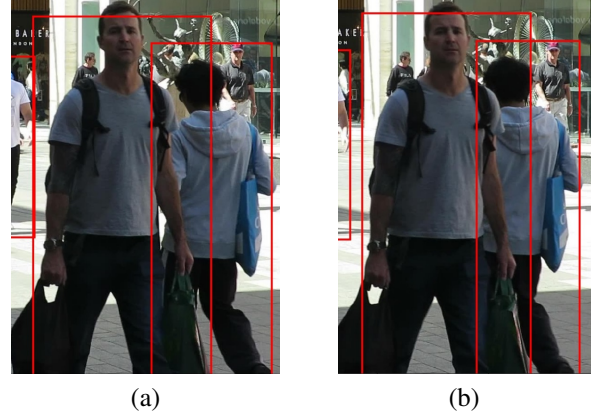
$$C_i = \min(1, \alpha * \frac{Age_i}{t_{so}} * \frac{A_i}{A_{avg}}), \quad (3)$$

in which  $\alpha$  is a hyperparameter.  $Age_i$  is the number of frames a target existed since its appearance. The  $t_{so}$  or time since observed, is the number of successive frames that a target is not matched with any detection.  $A_i$  denotes the area of the target and  $A_{avg}$  is the average area of all targets.  $A_{avg}$  is used to normalize the area of the bounding box of the target, to eliminate the effects of the camera position.

### 3.2. Occlusion Handling

In reality, a target may be partially or fully occluded by objects or other targets. The detection algorithm, only detects targets, in our case persons. So, there are not any corresponding detections for other objects. To detect occlusion by objects, the confidence metric is used. When a target does not match with any detection, if its confidence is greater than the object confidence threshold, which is represented as  $C_O$ , the target is marked as occluded. This can cover the weakness of the detection algorithm when there is no corresponding detection for a target due to an error in the detection algorithm.

The second kind of occlusion is when some other targets lie between the camera and a target. So, there is no corresponding detection to that covered target. In such cases, the intersection of the estimated bounding box of that target with the bounding box of the covering target is high. This property can help to detect such a kind of occlusion. The IoU metric is not appropriate for detecting occlusion. When the area of a target is very larger than another target, the IoU decreases since the intersection is divided into the total area of targets including the bigger one. This situation is common in occlusion, because the target which is closer to the camera and is



**Fig. 2:** Gradual occlusion. The detections are depicted by red rectangles. (a) A target is going to be occluded behind another target. (b) A part of target is occluded which leads to size reduction of its bounding box.

seen bigger, may occlude a far target which is seen smaller. To overcome this weakness, a novel metric, the so-called Coverd Percent (CP) is proposed, that measure how many percent of a target is covered by another target. The CP for the  $i^{th}$  target is defined as:

$$CP_i = \frac{I(bb_i, bb_j)}{A(bb_i)} \quad (4)$$

To detect target-target occlusion, minimum thresholds for CP and confidence metrics are defined. They are represented as  $CP_{min}$  and  $C_T$ , respectively. If the covered percent of an unmatched target with any of other targets is bigger than  $CP_{min}$  and its confidence is larger than  $C_T$  at the same time, it is marked as occluded, too.

For occluded targets, there are no corresponding detections to correct their estimation in the update step of their Kalman filters. For these targets, only the rate of area's change is halved to prevent the size of their estimated bounding boxes decrease unrealistically. This is because occlusion does not happen at once. First, part of the target becomes occluded and gradually becomes fully occluded. So, the size of the corresponding detection is reduced in frames leading to the occlusion, because the bounding box surrounds the observable part of the target.

$$\dot{s}_{t+1} = \frac{\dot{s}_t}{2} \quad (5)$$

The gradual occlusion of targets and the decrease of the size of the corresponding detection is shown in Fig. 2.

### 3.3. Target Re-identification

Target re-identification is done in a cascade matching manner. In the first step, the IoU between detections and all existing





**Fig. 3:** Example of re-identifying an occluded target. Detections are depicted with thin black rectangles, Targets are depicted with colored thick bounding boxes, extended bounding boxes are depicted with dashed rectangles. (a) A target is going to be occluded behind another target. (b) A target is occluded by another target and an extended bounding box is assigned to it. (c) The target is occluded behind another target for several frames and the size of its extended bounding box is increased. (d) The extended bounding box helps to re-identify the occluded target.

targets, including occluded targets are calculated and the association is performed using the Hungarian algorithm. Because the estimated bounding boxes of occluded targets do not correct in the update step of the Kalman filter, there may become a considerable difference between their estimation and real location. As a result, they may not match in the first step.

In the second step, to help the re-identification of unmatched occluded targets, their bounding boxes are extended according to their uncertainty. This uncertainty increases in every next frame where the target is not observed. Then the IoU between extended bounding boxes and remaining unmatched detections should be calculated. Because the sizes of occluded targets do not extend in reality and only their uncertainty increase, the extended IoU is calculated as below:

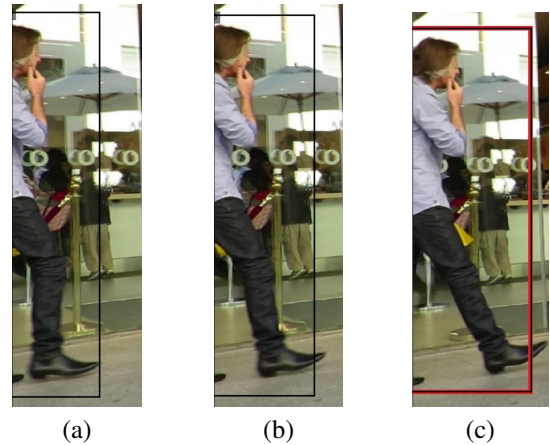
$$IoU_{ext} = \frac{I(bb_D, bb_{ext_T})}{A(bb_D) + A(bb_T) - I(bb_D, bb_{ext_T})} \quad (6)$$

In which  $bb_{ext_T}$  is the extended bounding box of the occluded target. The extended bounding box is only used for computing intersection and the area of the estimated bounding box is used in the denominator. To complete the second step, the calculated  $IoU_{ext}$  is passed to the Hungarian algorithm for the association. In Fig. 3, the extended bounding box of an occluded target is depicted by dashed lines.

### 3.4. Creating New Targets

During tracking, new targets may be entered into the scene. In the first few frames, all unmatched detections, represented by  $\bar{D}_t$  are considered as new targets. But after that, a new target is created after the probability of its existence increases. This is achieved when unmatched detections are presented near a location for three successive frames. For detecting this situa-

tion, the unmatched detection of the current frame, the previous frame, and the two frames before are used. To assign unmatched detections from different frames to each other, a process similar to associating detections to targets is performed. At first, an IoU matrix between  $\bar{D}_t$  and  $\bar{D}_{t-1}$  and another IoU matrix between  $\bar{D}_{t-1}$  and  $\bar{D}_{t-2}$  are calculated. Then each calculated IoU matrix is passed to the Hungarian algorithm. If three unmatched detections from different frames are assigned to each other, a new target is created and its Kalman filter is initialized using corresponding unmatched detections. The process of creating a new target is shown in Fig. 4.



**Fig. 4:** New target creation after three frames. The detections are shown with thin black rectangles and the target is shown with red thick rectangle. (a) A new target is seen for the first time. (b) The new target is seen for the second time. (c) The new target is seen for the third time and a target is associated with it.

		MOTA ↑	MOTP ↑	MT ↑	ML ↓	IDS ↓	FM ↓	FP ↓	FN ↓	FPS ↑
KDNT [30]	BATCH	<b>68.2</b>	79.4	41.0%	19.0%	933	1093	11479	45605	0.7
LMP_p [31]	BATCH	<b>71.0</b>	<b>80.2</b>	46.9%	21.9%	434	<b>587</b>	7880	44564	0.5
MCMOT_HDM [32]	BATCH	62.4	78.3	31.5%	24.2%	1394	1318	9855	57257	<b>35</b>
NOMTwSDO16 [33]	BATCH	62.2	79.6	32.5%	31.1%	<b>406</b>	642	<b>5119</b>	63352	<b>3</b>
EAMTT [34]	ONLINE	52.5	78.8	19%	34.9%	910	<b>1321</b>	<b>4407</b>	81223	12
POI [30]	ONLINE	<b>66.1</b>	79.5	34%	20.8%	805	3093	5061	<b>55914</b>	10
SORT [7]	ONLINE	59.8	<b>79.6</b>	25.4%	22.7%	1423	1835	8698	63245	<b>60</b>
Deep SORT [13]	ONLINE	61.4	79.1	32.8%	18.2%	<b>781</b>	2008	12852	56668	40
<b>Proposed</b>	ONLINE	61.1	79.04	31.62%	21.34%	848	1331	12296	57738	<b>162.7</b>

**Table 1:** The results of our proposed tracking algorithm compared to the results of state-of-the-arts on MOT16. The best are in the bold format. The second best are in the blue format.

	Detection	MOTA ↑	MT ↑	ML ↓	IDS ↓	FM ↓	FP ↓	FN ↓	FPS ↑
Tractor	ALL	53.5	19.5%	36.6%	2072	4611	<b>12201</b>	248047	1.5
GCNNMatch	ALL	<b>57.0</b>	<b>23.3%</b>	<b>34.6%</b>	<b>1957</b>	<b>2798</b>	12283	<b>228242</b>	1.3
<b>Proposed</b>	ALL	44.3	14.4%	45.6%	2191	5243	21796	290065	<b>137.5</b>
Tractor	DPM	52.2	14.9%	37.5%	635	-	<b>2908</b>	86275	-
GCNNMatch	DPM	<b>55.5</b>	<b>21.5%</b>	37.6%	564	<b>782</b>	2937	<b>80242</b>	-
<b>Proposed</b>	DPM	29.7	5.5%	64.7%	<b>530</b>	1448	3048	128696	-
Tractor	FRCNN	52.9	16.2%	34.7%	648	-	<b>3918</b>	83904	-
GCNNMatch	FRCNN	<b>56.1</b>	<b>22.3%</b>	<b>33.9%</b>	<b>647</b>	<b>934</b>	4015	<b>77950</b>	-
<b>Proposed</b>	FRCNN	44.9	14.7%	40.3%	795	1571	9102	93669	-
Tractor	SDP	55.3	18.1%	32.9%	789	-	5375	77868	-
GCNNMatch	SDP	<b>59.5</b>	<b>26.0%</b>	32.4%	<b>746</b>	<b>1082</b>	<b>5331</b>	70050	-
<b>Proposed</b>	SDP	58.4	22.8%	<b>32.0%</b>	866	2224	9646	<b>67700</b>	-

**Table 2:** The results of our proposed tracking algorithm compared to the results of state-of-the-art algorithms on MOT17. The best are in the bold format. The second best are in the blue format.

### 3.5. Removing Targets

An unmatched target is removed, when its uncertainty increases above a threshold. Like the occlusion handling section, the uncertainty is proportional to the number of frames in which the target is visible. Any unmatched target is retained for at least  $k_{min}$  frames. But more confident targets are retained for more frames proportional to their age. To prevent retaining targets for more than needed, targets are retained for at most  $k_{max}$  frames. An unmatched target is retained when:

$$t_{su} > \min(k_{min} + \frac{Age_T}{c_k}, k_{max}) \quad (7)$$

in which  $t_{su}$ , time since updated, is the number of successive frames the estimated bounding box of a target is not corrected. Because Kalman filters of occluded targets are updated, they do not enter the deletion process.

## 4. EXPERIMENTS

To compare the proposed algorithm with SORT [7] and DEEP SORT [13] algorithms, it is evaluated on MOT16 benchmark using private detections from POI [30] paper. For a fair

comparison, like POI and DEEP SORT papers, detections have been thresholded at a confidence score of 0.3. Also for comparing the proposed algorithm with the most recent algorithms, their results on MOT17 are presented.

### 4.1. Evaluation Metrics

For evaluating multiple object tracking algorithms, the frequently used CLEAR MOT metrics [35] are reported, including Multiple Object Tracking Accuracy (MOTA), and Multiple Object Tracking Precision (MOTP). Also, other popular metrics are used including Mostly Tracked (MT), Mostly Lost (ML), the number of False Negatives (FN), False Positives (FP), ID-Switches (IDS), and track Fragmentation (FM).

### 4.2. Results

The result of running the proposed algorithm on the MOT16 benchmark alongside the result of baseline algorithms such as SORT and DEEP SORT come in the table 1. In the MOTA and MOTP metrics, the results of the proposed algorithm are comparable with other online algorithms and it is slightly lower than the Deep SORT algorithm. For the IDS metric,

---

**Algorithm 1:** SORT with occlusion handling

---

**Data:**  $D_t, T_{t-1}$   
**Result:**  $\bar{D}_t, \bar{T}_t, \bar{O}_t$

- 1  $\hat{T}_t = KF\_Predict(T_{t-1})$
- 2  $A_{avg} = Area\_Average(\hat{T}_t)$
- 3  $\tilde{T}_t, \tilde{D}_t, O_t, \bar{T}_t, \bar{D}_t = Associate(\hat{T}_t, D_t, A_{avg})$
- 4  $\hat{T}_t = KF\_Correction(\tilde{T}_t, \tilde{D}_t)$
- 5  $O_t = Correction(O_t)$
- 6  $T_t \leftarrow \hat{T}_t + \bar{T}_t + O_t$
- 7
- 8 **Function**  $Associate(\hat{T}_t, D_t, A_{avg})$
- 9   **Cascade Matching:**
- 10    $P = IoU(D_t, \hat{T}_t)$
- 11    $\tilde{T}_t, \tilde{D}_t, \bar{T}_t, \bar{D}_t = Hungarian(P)$
- 12   **Target Re-identification:**
- 13    $P_d = IoU(\bar{D}_t, \bar{D}_{t-1})$
- 14    $P_{ext} = IoU_{ext}(\tilde{T}_t, \tilde{D}_t)$
- 15    $\tilde{T}_t^n, \tilde{D}_t^n = 2\_Step\_Matching(P_d, P_{ext})$
- 16    $\hat{T}_t \leftarrow \tilde{T}_t + \tilde{T}_t^n$
- 17    $\bar{T}_t \leftarrow \bar{T}_t - \tilde{T}_t^n$
- 18    $\tilde{D}_t \leftarrow \tilde{D}_t + \tilde{D}_t^n$
- 19    $\bar{D}_t \leftarrow \bar{D}_t - \tilde{D}_t^n$
- 20   **Detecting Occlusion:**
- 21    $O_t \leftarrow \emptyset$
- 22    $P = CP(\hat{T}_t)$
- 23   **for**  $u \in \bar{T}_t$  **do**
- 24      $C_u = \min(1, \alpha * \frac{Age_u}{t_{so_u}} * \frac{A_u}{A_{avg}})$
- 25      $P_u = \max(P)$
- 26     **if**  $C_u > C_O$  **then**
- 27        $O_t \leftarrow O_t + u$
- 28        $\bar{T}_t \leftarrow \bar{T}_t - u$
- 29     **else if**  $C_u > C_T$  **and**  $P_u > min\_coverage$  **then**
- 30        $O_t \leftarrow O_t + u$
- 31        $\bar{T}_t \leftarrow \bar{T}_t - u$

---

in comparison to the SORT algorithm, the value from 1423 is decreased to 848 which is a 40% reduction in this metric and it is only 8% above the Deep SORT algorithm. Also, in fragmentation metric, the score of the proposed algorithm is 1331 which is 28% lower than the SORT score which is 1835, and is even 34% lower than the Deep SORT algorithm with the value of 2008. This is because the Deep SORT algorithm does not detect occlusions and uses only the appearance to re-identify targets after they appear again. All of these comparable or even better results are achieved with the only use of geometric features and with a much lower computation cost that leads to a much higher speed, even with ordinary hardware.

---

**Algorithm 2:** Target creation and removal

---

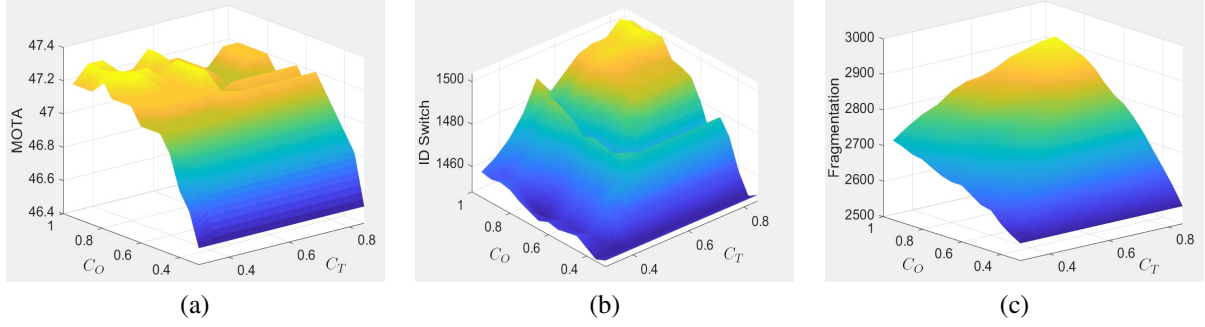
**Data:**  $\bar{D}_t, \bar{D}_{t-1}, \bar{D}_{t-2}, T_t, \bar{T}_t$   
**Result:**  $\bar{D}_t, \bar{D}_{t-1}, T_t$

- 1  $N_t \leftarrow \emptyset$
- 2  $R_t \leftarrow \emptyset$
- 3 **if**  $frame\_num < min\_hits$  **then**
- 4    $N_t \leftarrow \bar{D}_t$
- 5    $\bar{D}_t \leftarrow \emptyset$
- 6 **else**
- 7    $N_t = Find\_New\_Target(\bar{D}_t, \bar{D}_{t-1}, \bar{D}_{t-2})$
- 8    $R_t = Remove\_Target(\bar{T}_t)$
- 9    $T_t \leftarrow T_t + N_t - R_t$
- 10
- 11 **Function**  $Find\_New\_Target(\bar{D}_t, \bar{D}_{t-1}, \bar{D}_{t-2})$
- 12    $N_t \leftarrow \emptyset$
- 13    $P_d = IoU(\bar{D}_t, \bar{D}_{t-1})$
- 14    $P_{db} = IoU(\bar{D}_{t-1}, \bar{D}_{t-2})$
- 15    $N_t, N_{t-1} = 2\_Step\_Matching(P_d, P_{db})$
- 16    $\bar{D}_t \leftarrow \bar{D}_t - N_t$
- 17    $\bar{D}_{t-1} \leftarrow \bar{D}_{t-1} - N_{t-1}$
- 18
- 19 **Function**  $Remove\_Target(\bar{T}_t)$
- 20    $R_t \leftarrow \emptyset$
- 21   **for**  $u \in \bar{T}_t$  **do**
- 22     **if**  $t_{su_u} > \min(k_{min} + \frac{Age_u}{c_k}, k_{max})$  **then**
- 23        $R_t \leftarrow R_t + u$

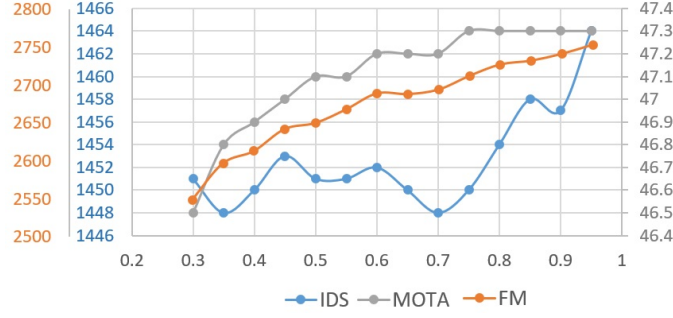
---

The MOT17 benchmark includes three sets of detections for every sequence. The detection algorithms are DPM, FRCNN, and SDP. The results of a tracking algorithm on this benchmark are the aggregation of its results in each detection set. In Table 2 the result of two state-of-the-art tracking algorithms and the proposed algorithm are separated for each detection set.

Between these three detection algorithms, the DPM has the lowest performance. The performance of the FRCNN algorithm is better than DPM, but lower than SDP and the SDP algorithm has the best performance. For the DPM and FRCNN detection sets, there is a significant difference between the result of the proposed algorithm and these two algorithms. But for the SDP detection set, the results of the proposed algorithm are better than results of the algorithm proposed in [4] and is comparable with results of algorithm discussed in [5]. These two tracking algorithms use appearance features and their speed even with high-performance hardware is lower than 2 FPS which is not appropriate for real-time applications. But the proposed algorithm reaches comparable results with the only use of geometric features while it is much faster than them.



**Fig. 5:** The changes of three evaluation metrics on confidence thresholds are presented; (a) MOTA; (b) IDS; (c) FM



**Fig. 6:** MOTA, IDS and FM versus  $C_O$  for a fixed value of  $C_T$

#### 4.3. Sensitivity to parameters

The proposed algorithm relies on several parameters from which  $C_O$  and  $C_T$ , i.e. the confidence thresholds for detecting an occluded target, are the most important ones. Thus, the sensitivity of the algorithm has been evaluated with respect to these two parameters. The sensitivity evaluation is done on the train data only since the test data was not accessible for exhaustive analysis. The evaluation is done on both MOT16 and MOT17 datasets, separately. In Fig. 5, the changes in MOTA, IDS, and FM metrics based on the changes in  $C_O$  and  $C_T$  are plotted for MOT17. Evidently, increasing  $C_O$  improves MOTA, and increasing  $C_T$  slightly decreases it. Increasing both  $C_O$  and  $C_T$  increases IDS in general but there are intermediate values in which IDS is near to its minimum. Also, increasing both  $C_O$  and  $C_T$  increases the FM which is not desired. To optimize the performance of the proposed algorithm, a cost function proportional to the IDS and proportional to the inverse of the MOTA is defined. After minimizing this cost function, the best results for MOT17 are achieved with  $C_O = 0.75$  and  $C_T = 0.35$ . The general behavior of the algorithm is the same on MOT16 too and  $C_O = 0.9$  and  $C_T = 0.55$  are chosen for it. Tuning these parameters can indeed improve the performance of the algorithm, but the important point is that the sensitivity of the algorithm to these parameters is low. For example in MOT17, the minimum and maximum of MOTA for the whole range of  $C_O$  and  $C_T$  is 46.5 and 47.3, respectively which is less than 1 percent. Also

the minimum and maximum of IDS are 1448 and 1504 which is less than 4 percent. In addition, for a fixed value of  $C_T$ , MOTA, IDS and FM are plotted versus  $C_O$  in Fig. 6. Needless to say, by increasing  $C_O$  MOTA and FM grow and IDS changes non-monotonically. To balance between these metrics, 0.75 is chosen for  $C_O$  at which MOTA is maximum, IDS is near its minimum and FM is not too high.

## 5. CONCLUSION

In this paper, a novel algorithm for tracking multiple objects is proposed to handle occlusions and re-identification of lost targets efficiently. This algorithm only uses geometric cues including the location and size of the bounding box of detections. As a result, it is very fast and is appropriate for real-time applications. The performance of this algorithm is comparable to other algorithms that use appearance features and could reach comparable results in terms of MOTA, MOTP, and IDS for the MOT16 dataset in comparison to the Deep SORT algorithm. For the FM metric, the results are better than the Deep SORT algorithm. Moreover, in the MOT17 benchmark, its result is comparable with the state-of-the-art algorithms for the SDP detection set. The results of the proposed algorithm on the MOT17 dataset show that if the detection algorithm has good performance, there is no need to use appearance features to achieve desired performance and only geometric cues can be used to achieve a much higher speed.

## 6. REFERENCES

- [1] Mohammadreza Babaei, Yue You, and Gerhard Rigoll, “Combined segmentation, reconstruction, and tracking of multiple targets in multi-view video sequences,” *Computer Vision and Image Understanding*, vol. 154, pp. 166–181, 2017.
- [2] Mohammadreza Babaei, Yue You, and Gerhard Rigoll, “Pixel level tracking of multiple targets in crowded environments,” in *European Conference on Computer Vision*. Springer, 2016, pp. 692–708.
- [3] Xinshuo Weng, Yongxin Wang, Yunze Man, and Kris Kitani, “Gnn3dmot: Graph neural network for 3d multi-object tracking with multi-feature learning,” *arXiv preprint arXiv:2006.07327*, 2020.
- [4] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe, “Tracking without bells and whistles,” in *Proceedings of the IEEE international conference on computer vision*, 2019, pp. 941–951.
- [5] Ioannis Papakis, Abhijit Sarkar, and Anuj Karpatne, “Gcnmatch: Graph convolutional neural networks for multi-object tracking via sinkhorn normalization,” *arXiv preprint arXiv:2010.00067*, 2020.
- [6] Yang Zhang, Hao Sheng, Yubin Wu, Shuai Wang, Wei Ke, and Zhang Xiong, “Multiplex labeling graph for near-online tracking in crowded scenes,” *IEEE Internet of Things Journal*, vol. 7, no. 9, pp. 7892–7902, 2020.
- [7] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft, “Simple online and realtime tracking,” in *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2016, pp. 3464–3468.
- [8] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler, “Mot16: A benchmark for multi-object tracking,” *arXiv preprint arXiv:1603.00831*, 2016.
- [9] Pedro Felzenszwalb, David McAllester, and Deva Ramanan, “A discriminatively trained, multiscale, deformable part model,” in *2008 IEEE conference on computer vision and pattern recognition*. IEEE, 2008, pp. 1–8.
- [10] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [11] Fan Yang, Wongun Choi, and Yuanqing Lin, “Exploit all the layers: Fast and accurate cnn object detector with scale dependent pooling and cascaded rejection classifiers,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 2129–2137.
- [12] Rudolph Emil Kalman, “A new approach to linear filtering and prediction problems,” 1960.
- [13] Nicolai Wojke, Alex Bewley, and Dietrich Paulus, “Simple online and realtime tracking with a deep association metric,” in *2017 IEEE international conference on image processing (ICIP)*. IEEE, 2017, pp. 3645–3649.
- [14] Lu Wang, Lisheng Xu, Min Young Kim, Luca Rigazico, and Ming-Hsuan Yang, “Online multiple object tracking via flow and convolutional features,” in *2017 IEEE International Conference on Image Processing (ICIP)*. IEEE, 2017, pp. 3630–3634.
- [15] Anton Milan, Seyed Hamid Rezatofighi, Anthony Dick, Ian Reid, and Konrad Schindler, “Online multi-target tracking using recurrent neural networks,” *arXiv preprint arXiv:1604.03635*, 2016.
- [16] Laura Leal-Taixé, Gerard Pons-Moll, and Bodo Rosenhahn, “Everybody needs somebody: Modeling social and grouping behavior on a linear programming multiple people tracker,” in *2011 IEEE international conference on computer vision workshops (ICCV workshops)*. IEEE, 2011, pp. 120–127.
- [17] Boyu Chen, Dong Wang, Peixia Li, Shuang Wang, and Huchuan Lu, “Real-time actor-critic tracking,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 318–334.
- [18] Li Zhang, Yuan Li, and Ramakant Nevatia, “Global data association for multi-object tracking using network flows,” in *2008 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2008, pp. 1–8.
- [19] Minyoung Kim, Stefano Alletto, and Luca Rigazio, “Similarity mapping with enhanced siamese network for multi-object tracking,” *arXiv preprint arXiv:1609.09156*, 2016.
- [20] Kwang-Yong Kim, Jun-Seok Kwon, and Kee-Seong Cho, “Multi-object tracker using kemelized correlation filter based on appearance and motion model,” in *2017 19th International Conference on Advanced Communication Technology (ICACT)*. IEEE, 2017, pp. 761–764.
- [21] Haoyang Feng, Xiaofeng Li, Peixin Liu, and Ning Zhou, “Using stacked auto-encoder to get feature with continuity and distinguishability in multi-object tracking,” in *International Conference on Image and Graphics*. Springer, 2017, pp. 351–361.



- [22] Maryam Babaei, Zimu Li, and Gerhard Rigoll, "A dual cnn-rnn for multiple people tracking," *Neurocomputing*, vol. 368, pp. 69–83, 2019.
- [23] Chanhon Kim, Fuxin Li, and James M Rehg, "Multi-object tracking with neural gating using bilinear lstm," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 200–215.
- [24] Mohib Ullah and Faouzi Alaya Cheikh, "A directed sparse graphical model for multi-target tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1816–1823.
- [25] Chanhon Kim, Fuxin Li, Arridhana Ciptadi, and James M Rehg, "Multiple hypothesis tracking revisited," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4696–4704.
- [26] Liangliang Ren, Jiwen Lu, Zifeng Wang, Qi Tian, and Jie Zhou, "Collaborative deep reinforcement learning for multi-object tracking," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 586–602.
- [27] Qiankun Liu, Qi Chu, Bin Liu, and Nenghai Yu, "Gsm: Graph similarity model for multi-object tracking," .
- [28] Yongxin Wang, Kris Kitani, and Xinshuo Weng, "Joint object detection and multi-object tracking with graph neural networks," .
- [29] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl, "Tracking objects as points," *arXiv preprint arXiv:2004.01177*, 2020.
- [30] Fengwei Yu, Wenbo Li, Quanquan Li, Yu Liu, Xiaohua Shi, and Junjie Yan, "Poi: Multiple object tracking with high performance detection and appearance feature," in *European Conference on Computer Vision*. Springer, 2016, pp. 36–42.
- [31] Margret Keuper, Siyu Tang, Yu Zhongjie, Bjoern Andres, Thomas Brox, and Bernt Schiele, "A multi-cut formulation for joint segmentation and tracking of multiple objects," *arXiv preprint arXiv:1607.06317*, 2016.
- [32] Byungjae Lee, Enkhbayar Erdene, Songguo Jin, Mi Young Nam, Young Giu Jung, and Phill Kyu Rhee, "Multi-class multi-object tracking using changing point detection," in *European Conference on Computer Vision*. Springer, 2016, pp. 68–83.
- [33] Wongun Choi, "Near-online multi-target tracking with aggregated local flow descriptor," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 3029–3037.
- [34] Ricardo Sanchez-Matilla, Fabio Poiesi, and Andrea Cavallaro, "Online multi-target tracking with strong and weak detections," in *European Conference on Computer Vision*. Springer, 2016, pp. 84–99.
- [35] 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. 246309, 2008.