

FOLT: Fast Multiple Object Tracking from UAV-captured Videos Based on Optical Flow

Mufeng Yao
School of Computer Science,
Shanghai Key Laboratory of Data
Science, Fudan University
mfyao21@m.fudan.edu.cn

Jiaqi Wang
School of Computer Science,
Shanghai Key Laboratory of Data
Science, Fudan University
21212010033@m.fudan.edu.cn

Jinlong Peng
Tencent Youtu Lab
jeromepeng@tencent.com

Mingmin Chi*
School of Computer Science,
Shanghai Key Laboratory of Data
Science, Fudan University
Zhengzhou Zhongke Institute of
Integrated Circuit and System
Application
mmchi@fudan.edu.cn

Chao Liu
School of Computer Science,
Shanghai Key Laboratory of Data
Science, Fudan University
Zhengzhou Zhongke Institute of
Integrated Circuit and System
Application
chaoliu@fudan.edu.cn

ABSTRACT

Multiple object tracking (MOT) has been successfully investigated in computer vision. However, MOT for the videos captured by unmanned aerial vehicles (UAV) is still challenging due to small object size, blurred object appearance, and very large and/or irregular motion in both ground objects and UAV platforms. In this paper, we propose FOLT to mitigate these problems and reach fast and accurate MOT in UAV view. Aiming at speed-accuracy trade-off, FOLT adopts a modern detector and light-weight optical flow extractor to extract object detection features and motion features at a minimum cost. Given the extracted flow, the flow-guided feature augmentation is designed to augment the object detection feature based on its optical flow, which improves the detection of small objects. Then the flow-guided motion prediction is also proposed to predict the object's position in the next frame, which improves the tracking performance of objects with very large displacements between adjacent frames. Finally, the tracker matches the detected objects and predicted objects using a spatially matching scheme to generate tracks for every object. Experiments on Visdrone and UAVDT datasets show that our proposed model can successfully track small objects with large and irregular motion and outperform existing state-of-the-art methods in UAV-MOT tasks.

CCS CONCEPTS

• **Computing methodologies** → **Tracking**.

*Corresponding author

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MM '23, October 29–November 3, 2023, Ottawa, ON, Canada.

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0108-5/23/10...\$15.00
<https://doi.org/10.1145/3581783.3611868>

KEYWORDS

multiple object tracking, optical flow, motion modeling, feature fusion

ACM Reference Format:

Mufeng Yao, Jiaqi Wang, Jinlong Peng, Mingmin Chi, and Chao Liu. 2023. FOLT: Fast Multiple Object Tracking from UAV-captured Videos Based on Optical Flow. In *Proceedings of the 31st ACM International Conference on Multimedia (MM '23)*, October 29–November 3, 2023, Ottawa, ON, Canada. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/3581783.3611868>

1 INTRODUCTION

Multiple object tracking (MOT) aims at identifying objects at each moment from a given video and is widely used in computer vision [7, 20, 29], such as autonomous driving [14], human-computer interaction [6], and pedestrian tracking [30, 34]. A common approach in Multiple Object Tracking (MOT) involves two primary stages: detection and association [37]. The detection step identifies all objects in each frame, while the association step links objects across consecutive frames to establish complete trajectories for each object [28]. Recently, MOT in unmanned aerial vehicle (UAV) platforms has attracted extensive research interest. Compared with conventional MOT [8, 26], MOT in UAV view faces more challenges.

Firstly, both the ground object and the UAV platform have fast and irregular motion, which makes the tracker difficult to follow the object. Secondly, the fast and irregular motion will reduce the image quality [25] and affect the detection of the object. Thirdly, the object size in aerial view is small, which not only increases the detection difficulty of the detector but also makes the appearance feature of the object unreliable, reducing the accuracy of appearance-matching methods. We define the mean relative acceleration (MRA) of video sequences in Eq (1) to measure the complexity of object motion. MRA of a sequence is calculated as the mean of object center acceleration normalized by object size. As shown in Fig. 1, the MRA of objects in the Visdrone [52] dataset is much higher than in the MOT17/20 dataset, showing that object motion patterns in UAV view are more complex than in the conventional MOT task,

which requires our tracker modeling the object motion more accurately. Fig. 1 also indicate that object size in aerial-view videos is smaller compared with traditional street-view videos, which limits the detection performance of the current tracker.

To address these problems, we propose FOLT (Fast Optical flow Tracker), which utilizes optical flow to model the object motion, augment the object detection feature, and improve tracking performance in UAV view. The FOLT first use a modern detector and flow-estimator to extract object detection feature and estimate a pixel-wise optical flow map. Based on the detection features and flow map, the flow-guided feature augmentation is proposed to augment object features at the current frame using the combination of previous features and current optical flow, which improves the detection of small objects. Then, the flow-guided motion prediction is proposed to model the object's motion and predict its position at a future time, which improves the tracking of objects with large and irregular motion. Finally, the predicted objects are spatially matched with the detected objects and output tracking results at the current time.

By augmenting the detection feature and explicitly modeling the motion of objects, FOLT can not only improve the detection accuracy but also obtain more stable tracking results in fast-moving scenes of UAV view. Experiments indicate that our FOLT performs better than previous appearance-matching-based methods in both accuracy and speed, which support our belief that appearance-matching strategy is not important in UAV-MOT scenes.

The main contributions of this paper are summarized as follows:

- We propose the flow-guided feature augmentation, which combines previous detection features with current features according to optical flow, improves the detection of small objects and mitigates the motion blur problems in UAV-MOT tasks.
- We propose the flow-guided motion prediction, which predicts object position according to optical flow, surpassing the commonly used Kalman Filter in both accuracy and speed.
- The FOLT proposed based on flow-guided feature augmentation and flow-guided motion prediction reached state-of-the-art on two public MOT-in-UAV-view datasets, which promoted the progress of multiple object tracking in UAV scenes.

2 RELATED WORK

2.1 Motion modeling in MOT

Current MOT methods generally follow the tracking-by-detection scheme [1, 11, 29, 32, 35, 36], with the detection step using a deep neural network to output the detection result of each frame, and the association step completes the inter-frame association based on object appearance [24, 27, 33, 39, 46] or motion information [12, 23, 42, 49]. The association step can be categorized as appearance based [43, 44, 48, 50] and motion based [2, 3, 37, 47, 51]. As Fig. 1 and Fig. 9 show, objects in aerial view are very small and blurry, which will result in unreliable appearance features for the association. Therefore, this paper focuses on motion-based association to avoid unreliable appearance feature problems. Motion-based association first predicts the position of each object trajectory in the next frame, then it matches the objects detected in the next frame

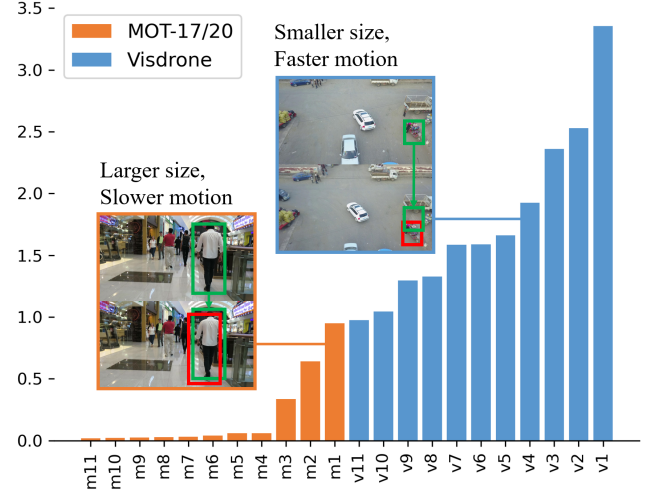


Figure 1: Mean relative accelerations (MRA) of conventional MOT dataset (MOT17/20) and UAV-MOT dataset (Visdrone). All sequences of the MOT17/20 train set (11 in total) and the top 11 sequences of the Visdrone train set are displayed. Green and red rectangle: Visdrone has a smaller object size and faster motion compared with MOT17/20.

according to their position adjacency, to generate trajectory at the next frame. The Kalman Filter [21] is the most commonly used motion modeling method and widely adopted by other motion-based trackers [2, 5, 47]. SORT [2] is the classical tracker that applies the Kalman Filter [21] to predict object position. ByteTrack [47] makes significant progress by replacing SORT with a more powerful detector and a more complex two-stage matching method. OCSORT [5] improves the original Kalman Filter in an object-centric manner and performs better in irregular motion scenes. However, when both objects and the UAV platform are in fast-moving pattern, these Kalman Filter (KF) based methods is hard to follow the large and irregular motion. In addition to KF based method, SiamMOT [37] uses the cross-correlation of the local object regions between adjacent frames to predict the position of the object, CenterTrack [51] directly designs a CNN network to predict the object position, UAVMOT [28] uses the topology information of objects in adjacent frames to associate them. However, these methods introduce too much computation costs and can not reach real-time multiple object tracking. In this paper, we use a modern light-weight optical flow estimator to realize high speed and accurate motion modeling, which effectively improves tracking in UAV view.

2.2 Feature fusion in MOT

Several studies fuse motion information with appearance information in a simple weighted summation scheme[43, 44, 48]. For example, FairMOT[48] compute the cosine distance in feature space and Jaccard distance in image space, then calculate their weighted average value and match the objects according to the averaged distance. However, these methods didn't fuse information between objects in adjacent frames and are limited in improving object detection on small and blurred objects. Others use long-short-term memory (LSTM) networks to fuse the temporal features between

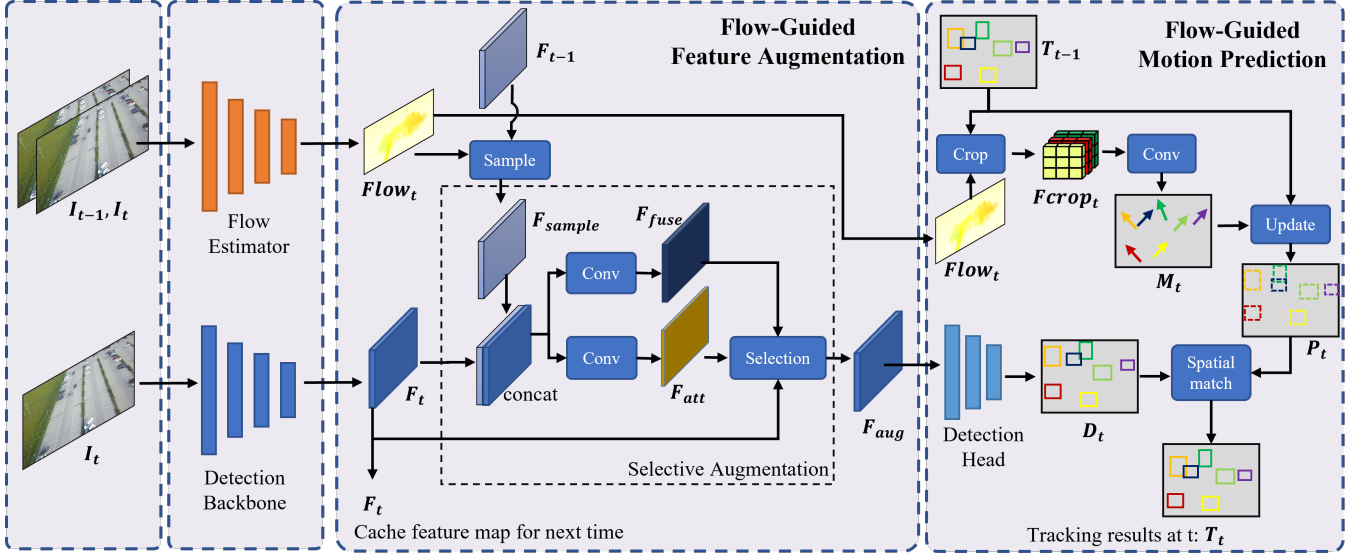


Figure 2: The architecture of our proposed FOLT. I_t, I_{t-1} denotes image frame at time t and time $t - 1$. $Flow_t, F_t$ denotes the optical flow and detection feature at t . F_{t-1} is the cached detection feature of the previous frame. D_t, M_t, P_t denotes the detection results at t , the motion prediction of objects at t , and the predicted objects position at t . D_t and P_t are spatially matched to output T_t : the final track results at t .

adjacent frames[31, 45]. However, these methods didn't consider motion information when fusing adjacent features. Different from previous methods, we use optical flow to sample features in their previous position and selectively fuse it with the current feature in an attentive way, which is helpful to objects with small sizes and blurred appearances.

2.3 Deep flow networks

Optical flow reflects the relative offset of each pixel between adjacent frames and has a wide range of applications in computer vision tasks [19], such as autonomous driving [15] action recognition [38], and pose tracking [4]. It is feasible to estimate the offset of each object between adjacent frames from a given optical flow map and detected object range. Traditional optical flow methods [17, 40] have a large computational overhead and are difficult to apply in real time. Several optical flow neural networks[9, 18, 41] have achieved promising results, providing high accuracy dense optical flow in real-time. However, these networks are still not fast enough to ensemble in a real-time multiple object tracker, since the detection and association have already cost some computation. Recently, the FastFlowNet[22] reached a comparative flow estimation accuracy while only costing 11ms per image in inference. Hence, it is possible to integrate this optical flow network into our tracking algorithm and utilize the extracted flow to enhance our object detection feature and predict the object motion.

3 PROPOSED METHOD

3.1 Overview

Different from conventional MOT tasks, MOT in the UAV platform faces more challenges such as small object size, similar and blurred object appearance, and large and irregular motion. We define mean

relative acceleration (MRA) to reveal this phenomenon:

$$MRA = \frac{1}{n} \sum_{i=1}^{n-1} (V(t_{i+1}) - V(t_i)), \quad (1)$$

among them, n is the length of video, $V(t_i)$ is the object's relative velocity at time i :

$$V(t_i) = \frac{\sqrt{(cx(t_i) - cx(t_{i-1}))^2 + (cy(t_i) - cy(t_{i-1}))^2}}{\sqrt{w(t_i)^2 + h(t_i)^2}}, \quad (2)$$

where $cx(t_i), cy(t_i)$ denote object center position at time t_i , $w(t_i), h(t_i)$ denote object width and height at time t_i . Capturing from the aerial top-down view, MOT in UAV datasets have higher MRA and smaller object sizes than conventional MOT datasets, which is summarized in Fig 1. The higher MRA means that the object's motion is larger and more irregular, which makes the tracker difficult to follow the object. In addition, the motion blurring and the small object size also increased the detection difficulty in UAV scenes. Aiming at these difficulties, we propose FOLT to improve the tracking accuracy of small and blurred objects in large and irregular motion scenes. As Fig. 2 shows, the proposed FOLT consists of three main components. The feature extraction component consists of a detection backbone and a flow estimator, with the detection backbone extracting the detection feature at current time F_t and the flow estimator estimating a pixel-wise offset map between the current frame and previous frame $Flow_t$. The flow-guided feature augmentation first samples previous detection features F_{t-1} at the current position to output F_{sample} . Then the F_{sample} and F_t are concatenated and fed into two branches of convolution to output fused feature F_{fuse} and attention weights F_{att} . Then the F_{fuse} and F_t are selectively fused in the guide of attention map F_{att} to output the augmented detection feature F_{aug} . Finally, the augmented detection feature F_{aug} is fed into the detection head to

output detection results: D_t . The flow-guided motion prediction takes optical flow map $Flow_t$ and previously tracked object position T_{t-1} as input and uses a convolution layer to predict the object motion M_t , then update the object position prediction P_t using M_t . Finally, the predicted object position P_t is matched with detected objects D_t to formulate T_t tracking results at time t . We use a two-stage spatial matching strategy [47] to match the detection results D_t with predicted position P_t . Among the feature extractor, we select the latest YOLOX [13] as our object detector and tested three pre-trained optical flow networks: FlowNet [9], PWCNet [41], and FastFlowNet [22] as our optical flow extractor. As Table 1 shows, the FastFlowNet achieves the best performance on both tracking accuracy and inference speed. Therefore, we adopt FastFlowNet as our base optical flow extractor in subsequent experiments. Alg 1 shows the tracking process of FOLT.

Algorithm 1 Tracking process of FOLT

Require: An video sequences $\{I_t \in \mathbb{R}^{H \times W \times 3}\}_{t=1}^T$, detection backbone: D_b , flow estimator: $FNet$, detection head: D_h , flow-guided feature augmentation: $FGFA$, flow-guided motion prediction: $FGMP$

Ensure: Tracking results Trk

```

1: Let  $t = 0, Trk = \{\}$ 
2: while  $t < T$  do
3:   if  $t == 0$  then
4:     Input:  $I_t$ 
5:     Extract detection feature  $(F_t^1, F_t^2, F_t^3) = D_b(I_t)$ 
6:     Conduct object detection  $D_t = D_h(F_t^1, F_t^2, F_t^3)$ 
7:   else
8:     Input: two adjacent frames  $I_{t-1}, I_t$ 
9:     Extract detection feature  $(F_t^1, F_t^2, F_t^3) = D_b(I_t)$ 
10:    Extract flow:  $Flow_t = FNet(I_{t-1}, I_t)$ 
11:    Augment detection feature with optical flow:  $F_{aug}^j = FGFA(Flow_t, F_{t-1}^j, F_t^j), j \in \{1, 2, 3\}$ 
12:    Conduct object detection with augmented features  $D_t = D_h(F_{aug}^1, F_{aug}^2, F_{aug}^3)$ 
13:    Predict object motion with  $FGMP: M_t = FGMP(Flow_t, Trk_{t-1})$ 
14:    Update tracks:  $Trk_t = Trk_{t-1} + M_t$ 
15:   end if
16:   Match tracks with current detection:  $Trk_t = Trk_t \cup D_t$ 
17: end while
18: return  $Trk$ 

```

3.2 Feature extraction

We use the YOLOX-S[13] backbone to extract detection features for every frame, which adopts a feature pyramid structure to output a three-stage feature map with different spatial resolutions: F_t^1, F_t^2, F_t^3 . The spatial resolution of the three-stage feature is downsampled by 8, 16, and 32 compared with the original input image. The original flow $Flow_t$ is extracted using FastFlowNet with a spatial resolution (512×384), which is lower than the input resolution of the detection branch (1088×608) for faster inference speed. We rescale the flow value when downsampling it, to keep the right spatial relationship

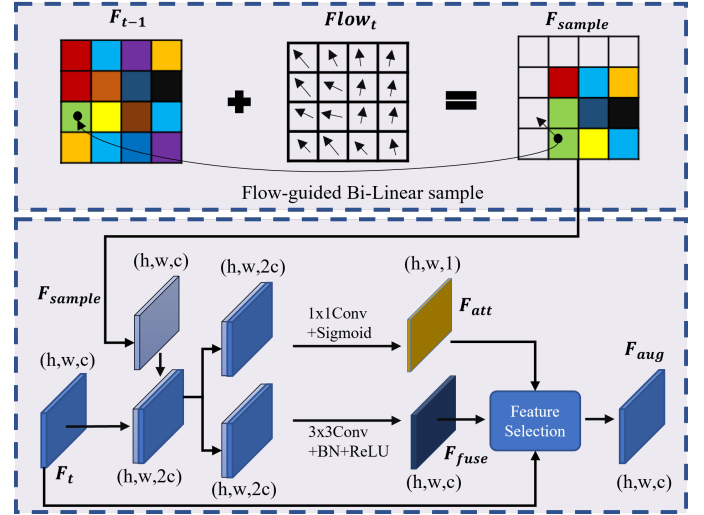


Figure 3: Flow-guided feature augmentation. F_{sample} denotes sampled features from the previous frame, F_t denotes current features, F_{fuse} , F_{att} are the fused feature and attention map, F_{aug} is the final augmented features.

between adjacent feature maps:

$$Flow_t^i(:, :, x) = sxFlow_t^d(:, :, x), \quad (3)$$

$$Flow_t^i(:, :, y) = syFlow_t^d(:, :, y), \quad (4)$$

where sx, sy is the down-sample scale in the x-axis and y-axis, $Flow_t^d(:, :, :)$ is the down-sampled flow map that needs to be rescaled. We only plot one stage of feature augmentation in Fig 2 and Fig 3 for simplicity. The flow-guided motion prediction is conducted in a single stage since the previously tracked objects are already merged in one stage. Therefore we use the original optical flow with resolution 512×384 as the input of the flow-guided motion prediction.

3.3 Flow-guided feature augmentation

As Fig 7 shows, objects in UAV-captured videos are typically small in size and their appearance is easily affected by the motion blurring, which increases the difficulty of object detection in UAV-view. Aiming at these challenges, we propose flow-guided feature augmentation (FGFA) to augment object features using previous object features and optical flow between the previous frame and the current frame. As Fig 3 shows, our FGFA first use optical flow $Flow_t$ to sample previous detection feature F_{t-1} using a bi-linear sample to output F_{sample} :

$$F_{sample}(i, j) = BS(F_{t-1}, i, j, dx, dy), \quad (5)$$

where the dx and dy are the optical flow values at position (i, j) : $dx = Flow_t(i, j)(x)$, $dy = Flow_t(i, j)(y)$. The bilinear sample function BS is defined by:

$$BS(F, i, j, dx, dy) = s_x s_y F(x, y) + s_y (1 - s_x) F(x + 1, y) + (1 - s_y) s_x F(x, y + 1) + (1 - s_y) (1 - s_x) F(x + 1, y + 1), \quad (6)$$

where $x = [i + dx]$, $y = [j + dy]$, $[\cdot]$ denote round down operation, $s_x = i + dx - [i + dx]$, $s_y = j + dy - [j + dy]$, F is the feature map. After the flow-guided feature sampling, the F_{sample} and F_t

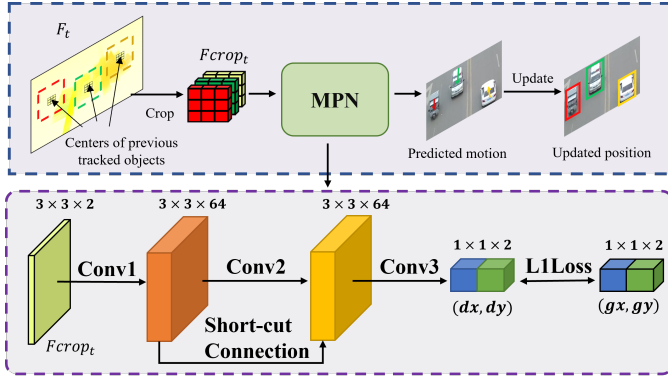


Figure 4: Flow-guided motion prediction. MPN denote motion prediction network, I_t, I_{t-1} denote video frames at time $t, t-1$, F_t denote dense optical flow at time t , F_{crop_t} denote 3×3 pixel of optical flow around the object center, dx, dy denote predicted object center offset at x, y axis, gx, gy denote ground-truth object center offset at x, y axis.

are concatenated in channel dimension and fed into two convolution branches: the attention branch and the fusion branch. The attention branch takes a 1×1 convolution to express $2c$ channels to 1 channel and then uses a sigmoid function to generate attention scores F_{att} for every spatial position. The fusion branch takes a 3×3 convolution following a batch normalization and a ReLU function to generate fused features F_{fuse} . Finally, the fused feature F_{fuse} and original feature F_t are selectively fused according to attention map F_{att} and output the augmented feature F_{aug} :

$$F_{aug} = F_{fuse} * F_{att} + F_t, \quad (7)$$

In accordance with the feature pyramid backbone of YOLOX-S, our flow-guided feature augmentation is conducted in three stages separately:

$$F_{aug}^i = FGFA(F_{t-1}^i, F_t^i, Flow_t^i), \quad (8)$$

We obtain $Flow_t^i$ by down-sample the original flow $Flow_t$ to the same spatial resolution with F_t^i using bilinear sampling.

3.4 Flow-guided Motion Prediction

Previous tracking algorithms generally use the Kalman Filter to estimate the object's future position. However, when both the object and the camera platform have large or irregular motions, the Kalman Filter cannot track the object correctly. We state that the pixel-wise optical flow map estimated by a deep flow network has already given a rough estimate of the position offset for each pixel. Therefore, it is practicable to use the statistics of optical flow within each bounding box to estimate its offset between adjacent frames. Starting from this thought, we propose flow-guided motion prediction to predict every object's motion based on object location and optical flow map. Fig 4 shows the architecture of our proposed flow-guided motion prediction. Given the extracted optical flow map $Flow_t$, we crop the optical flow around the 3×3 neighbor of the predicted object center F_{crop_t} . Then we feed it to the motion prediction network (MPN) to estimate the object motion between adjacent frames. As Fig. 4 shows, the MPN contains 3 layers of convolution, with the second layer having a short-cut connection

Table 1: Effectiveness of different optical flow networks. The best result is marked in bold. The ' \uparrow ' means that the higher result is better. KF denotes the Kalman Filter, FGMP denotes flow-guided motion prediction.

Motion modeling	Flow estimator	MOTA \uparrow	IDF1 \uparrow	FPS \uparrow
KF	-	39.6	50.4	27.0
FGMP	FlowNet-s	39.9	51.3	23.2
FGMP	PWCNet	40.5	54.1	15.6
FGMP	FastFlowNet	40.9	55.3	32.0

Table 2: Ablation studies on Visdrone test-dev set. FGMP indicates flow-guided motion prediction. FGFA indicates flow-guided feature augmentation. The ' \uparrow ' means that the higher result is better. The best result is marked in bold.

Baseline	FGMP	FGFA	MOTA \uparrow	IDF1 \uparrow	FPS \uparrow
\checkmark			39.6	50.4	27.0
\checkmark	\checkmark		40.9	55.3	32.0
\checkmark		\checkmark	40.9	52.3	25.5
\checkmark	\checkmark	\checkmark	42.1	56.9	29.4

inspired by ResNet [16]. Each convolution operation is followed by a batch normalization layer and a Softshrink function. Here, Softshrink is defined by:

$$\text{Softshrink}(x) = \begin{cases} x - \lambda, & x > \lambda \\ 0, & -\lambda \leq x \leq \lambda \\ x + \lambda, & x < -\lambda \end{cases}, \quad (9)$$

where λ is set to 0.5 in the experiments to ignore too lower values. The MPN outputs the predicted offset of the object center $M_t = (dx, dy)$. We use the L1 loss function to supervise the motion prediction:

$$L1(dx, dy, gx, gy) = |dx - gx| + |dy - gy|, \quad (10)$$

where gx, gy is the ground-truth offset of the object in x and y axis. The predicted object offset is then added to the previously tracked object position to output the final position estimate:

$$P_t = T_{t-1} + M_t. \quad (11)$$

Finally, the predicted object position P_t is matched with detected objects D_t using a two-stage spatial matching strategy [47] to generate tracking results at time t .

4 EXPERIMENTS

4.1 Datasets and Metrics

Datasets. To validate the effectiveness of our proposed FOLT, we conduct comparative experiments on the Visdrone [52] and UAVDT [10] datasets. These two datasets are open source multi-class multi-object tracking datasets, and both are collected from the perspective of UAVs. Therefore, it is suitable to study the tracking objects with small sizes and irregular motion problems in these two datasets.

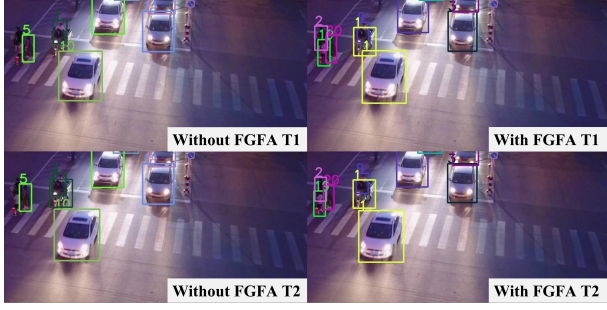


Figure 5: Visualization results on small object scenes. The model without FGFA missed the small-sized pedestrians on top-left corner while the model with FGFA tracked them with ID 2,19,20.

The Visdrone dataset consists of a training set (56 sequences), validation set (7 sequences), test-dev set (17 sequences), and test-challenge set (16 sequences). There are 10 categories in the Visdrone dataset: pedestrian, person, car, van, bus, truck, motor, bicycle, awning-tricycle, and tricycle. Each object within the above categories is annotated by a bounding box, category number, and unique identification number. In experiments of Visdrone, we use the full ten categories in training while only using five categories in testing, i.e., car, bus, truck, pedestrian, and van in evaluation, as the evaluation toolkit offered by Visdrone officials only evaluating in these five categories.

UAVDT dataset is a car-tracking dataset in aerial view, it includes different common scenes, such as squares, arterial streets, and toll stations. UAVDT dataset consists of a training set (30 sequences), and a test set (20 sequences), with three categories: car, truck, and bus. In experiments of UAVDT, all three categories are evaluated using Visdrone’s official evaluation toolkits.

Metrics. To evaluate our tracking efficiency in MOT tasks, we select MOTA and IDF1 as our main evaluation metrics. The MOTA is calculated as:

$$MOTA = 1 - \frac{FP + FN + ID_s}{GT}, \quad (12)$$

where FP denotes numbers of false positives, FN denotes numbers of false negatives, ID_s denote numbers of ID switches, and GT denotes numbers of ground-truth objects. The IDF1 is calculated as:

$$IDF1 = \frac{2IDTP}{2IDTP + IDFP + IDFN}, \quad (13)$$

where IDTP, IDFP, and IDFN denote numbers of true positive, false positive, and false negative that consider ID information.

4.2 Implementation Details

In all experiments, we keep the same train-test split as the official split of the Visdrone and UAVDT datasets. We use the YOLOX-S[13] model as the base object detector on both datasets, the input image size is 1088×608 . We use the stochastic gradient descent method to optimize the detector, the learning rate is set to 0.000005, the batch size is set to 4, each dataset is trained for 10 epochs, and the training and testing are completed on a single 2080TI graphic card. We use L1-Loss to train our motion prediction network and the object localization branch of the detector, and use the cross-entropy

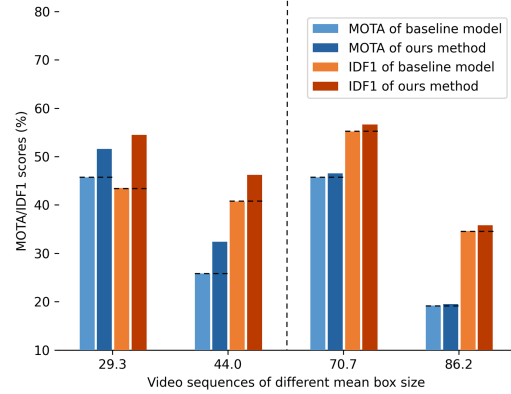


Figure 6: Comparisons of baseline and FOLT in Visdrone test-dev set. The left two groups are the results of two videos with smaller object sizes, while the right two groups are with larger object sizes. The improvement in MOTA and IDF1 compared with the baseline model are explicitly higher in videos with smaller objects than in videos with larger objects.

loss to train the object classification branch of the detector. In the ablation study, we select MOTA and IDF1 as our evaluation metrics, which are commonly used in MOT tasks. In comparison with state-of-the-art methods, we select the accuracy evaluation tool officially provided by the Visdrone dataset to complete the comparison of all metrics.

4.3 Ablation study

Baseline model. The baseline model we compared with is the model that use YOLOX-S as detector, extended Kalman Filter as motion modeling [5], and two-stage spatial matching as association [47].

Different optical-flow extractor. We evaluate the effectiveness of different optical-flow extractors in Table 1. As Table 1 shows, the modern optical flow network FastFlowNet performs better than previous old methods (FlowNet and PWCNet) in both tracking accuracy (40.9 in MOTA and 55.3 in IDF1) and inference speed (32.0 in FPS). Therefore, we adopt FastFlowNet as our base optical flow estimator in subsequent experiments.

Flow-guided feature augmentation. In this section, we evaluate the effectiveness of flow-guided feature augmentation (FGFA). As shown in Table 4, our FGFA performs better than the single frame feature and simple feature warp strategy, with MOTA 1.3 higher than the second best (40.9 to 39.6) and IDF1 1.6 higher than the second best (52.3 to 50.4), while only slower the inference speed a little (27.0 to 25.5 FPS). Fig 5 shows that after using FGFA, our model successfully tracks small pedestrians (with ID 2, 10, and 20), while the model without FGFA fails to track all of them. Fig 6 shows that after using FGFA, the tracking improvements (shown in black dash lines) of small objects are larger than normal objects. In summary, the qualitative and quantitative results above confirm the effectiveness of our proposed FGFA in improving object tracking with small size and motion blur.

Flow-guided motion prediction. As Table 5 shows, compared with the Kalman Filter, using optical flow as a motion estimator not

Table 3: Comparisons of the proposed FOLT with state-of-the-art methods on Visdrone and UAVDT test sets. The best result is marked in bold. The '↑'('↓') means that the higher (lower) result is better. Our proposed FOLT surpasses state-of-the-arts in both tracking accuracy and inference speed.

Dataset	Method	Pub & Year	Speed(FPS)↑	MOTA↑	MOTP↑	IDF1↑	FP↓	FN↓	IDs↓
VisDrone	GOG [36]	CVPR2011	2.0	28.7	76.1	36.4	17706	144657	1387
	SORT [2]	ICIP2016	23.5	14.0	73.2	38.0	80845	112954	3629
	IOUT [3]	AVSS2017	27.3	28.1	74.7	38.9	36158	126549	2393
	SiamMOT [37]	CVPR2021	11.2	31.9	73.5	48.3	24123	142303	862
	ByteTrack [47]	ECCV2022	27.0	35.7	76.8	37.0	21434	124042	2168
	UAVMOT [28]	CVPR2022	12.0	36.1	74.2	51.0	27983	115925	2775
	OCSORT [5]	CVPR2023	26.4	39.6	73.3	50.4	14631	123513	986
	FOLT (Ours)	Ours	29.4	42.1	77.6	56.9	24105	107630	800
UAVDT	GOG [36]	CVPR2011	2.0	35.7	72	0.3	62929	153336	3104
	SORT [2]	ICIP2016	23.5	39.0	74.3	43.7	33037	172628	3350
	IOUT [3]	AVSS2017	27.3	36.6	72.1	23.7	42245	163881	9938
	SiamMOT [37]	CVPR2021	11.2	39.4	76.2	61.4	46903	176164	190
	ByteTrack [47]	ECCV2022	27.0	41.6	79.2	59.1	28819	189197	296
	UAVMOT [28]	CVPR2022	12.0	46.4	72.7	67.3	66352	115940	456
	OCSORT [5]	CVPR2023	26.4	47.5	74.8	64.9	47681	148378	288
	FOLT (Ours)	Ours	29.4	48.5	80.1	68.3	36429	155696	338

Table 4: Effectiveness of our flow-guided feature augmentation module. The best result is marked in bold. The '↑' means that the higher result is better.

Motion modeling	MOTA↑	IDF1↑	FPS↑
Single frame	39.6	50.4	27
Naive feature warp	39.5	50.4	26.2
Flow-guided feature augmentation	40.9	52.3	25.5

Table 5: Effectiveness of our flow-guided motion prediction module. The best result is marked in bold. The '↑' means that the higher result is better.

Motion modeling	MOTA↑	IDF1↑	FPS↑
Kalman Filter	39.6	50.4	27
Mean of flow	40.4	53.1	32.3
Flow-guided Motion Prediction	40.9	55.3	32.0

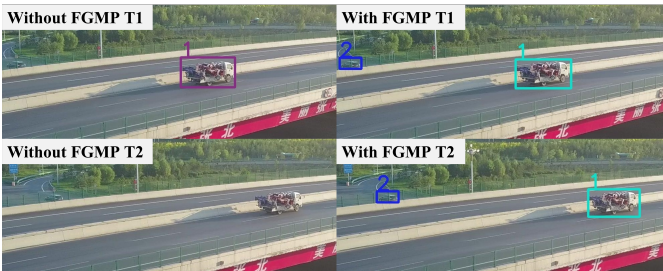


Figure 7: Visualization results on large and irregular motion scenes. The model without FGMP missed the fast-moving trucks on the highway at time T while the model with FGMP tracked it with ID 1.

only improves tracking accuracy (39.6 to 40.4 in MOTA and 50.4 to 53.1 in IDF1) but also increases inference speed (27.0 to 32.3 in FPS). Second, using flow-guided motion prediction (FGMP) further

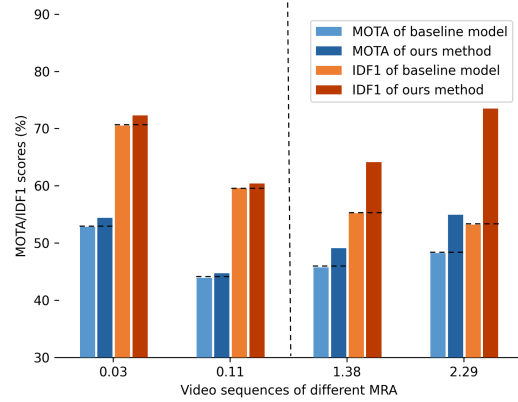


Figure 8: Comparisons of baseline and FOLT in Visdrone test-dev set. The left two groups are the results of two videos with low MRA, while the right two groups are with high MRA. The improvement in MOTA and IDF1 compared with the baseline model are explicitly higher in videos with high MRA than in videos with low MRA.

improves tracking accuracy (40.4 to 40.9 in MOTA and 53.1 to 55.3 in IDF1) but only introduces a little speed cost (32.3 to 32.0 in FPS).

As Fig 8 shows, after using FGMP in tracking, the improvements on objects with large and irregular motion (higher MRA) are larger than objects with slower motion (lower MRA). As Fig 7 shows, after using our FGMP, the tracker successfully tracks the fast-moving truck on the highway (see the object with ID1), while the model without FGMP fails to track the fast-moving truck at time t . In summary, the qualitative and quantitative results above confirm the effectiveness of our proposed FGMP in improving object tracking with large and irregular motion.

Combination. We also evaluate the effectiveness of our FGFA and FGMP which are used in combination. Table 2 shows that both

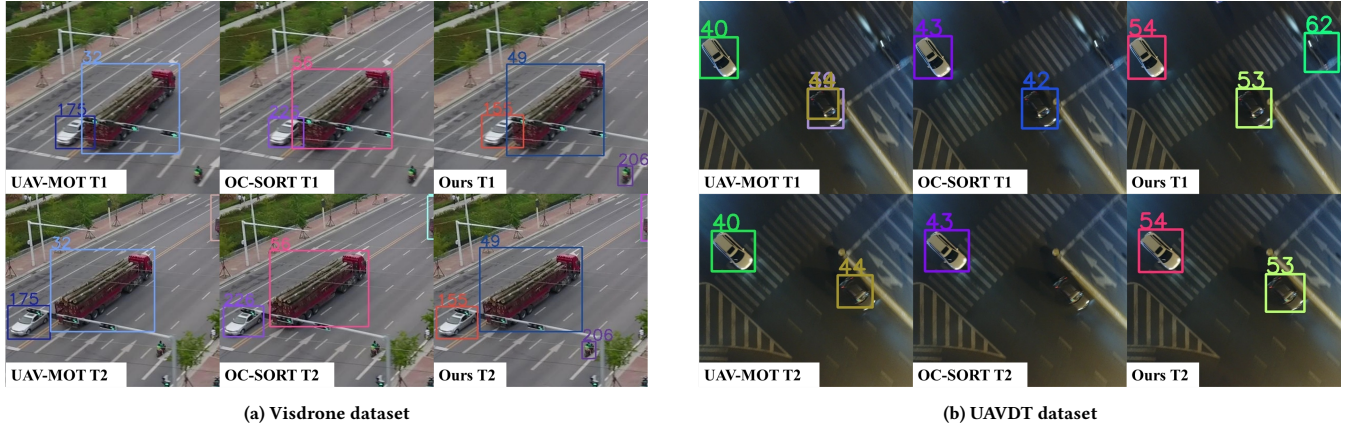


Figure 9: Visual comparisons of our method and state-of-the-arts on Visdrone and UAVDT datasets. The same number denotes the tracked same object in different frames. (a) Visual comparisons on the Visdrone dataset. (b) Visual comparisons on the UAVDT dataset.

FGFA and FGMP can effectively improve the tracking performance (39.6 to 40.9 and 39.6 to 40.9 in MOTA, 50.4 to 52.3 and 50.4 to 55.3 in IDF1). With the FGFA and FGMP commonly used, the best tracking performance is obtained (42.1 in MOTA and 56.9 in IDF1), and the tracking speed is faster than the baseline model (27.0 to 29.4 in FPS).

4.4 Comparison with State-of-the-arts

We conducted comparative experiments on the Visdrone test-dev set and the UAVDT test set. Table 3 shows the experimental results on these two datasets. As shown in Table 3, the FOLT proposed in this paper exceeds the current state-of-the-art (SOTA) methods in several important metrics such as MOTA, IDF1, IDs, and MOTP. For example, the MOTA of FOLT is 2.5% higher than the previous highest method OCSORT (42.1 compared with 39.6), the IDF1 is 5.9% higher than the previous highest method UAVMOT (56.9 compared with 51.0), while the inference speed is also faster than the previous fastest method ByteTrack (29.4 compared with 27.0). FOLT also achieves the best MOTA and IDF1 in the UAVDT dataset, with MOTA 1.0% higher than the previous best method OCSORT (48.5 compared with 47.5), and IDF1 1.0% higher than the previous best method UAVMOT, (68.3 compared with 67.3) showing the strongest ability in tracking objects from a UAV view. In addition, our method that only uses location information in object matching also performs better than those that use both location and appearance information such as UAVMOT and SiamMOT, validating our belief that appearance feature is not reliable in the small object with a blurry appearance.

We also visualize the comparison results of the latest best methods (UAVMOT and OCSORT) and our FOLT in the Visdrone dataset and the UAVDT dataset. As Fig 9a shows, our FOLT successfully tracks the small and fast-moving electronic bicycle rider (with object ID 206) in both T1 and T2 frames, while the UAV-MOT and the OC-SORT tracker all failed to track it in both frames. Therefore, our FOLT is better than UAVMOT and OCSORT in tracking objects with small sizes and fast motion. As Fig 9b shows, the camera and the object are both in motion, introducing large and irregular motion

and causing the object feature to blur. In this difficult situation, our FOLT successfully tracks three cars in frame T1 with object ID 54, 53, and 62 respectively, while the UAVMOT has duplicate tracking on the middle car and the OCSORT missed the right-top car that is severely blurred due to large and irregular motion. Our FOLT also successfully tracks two cars in frame T2 with object size 54 and 53, while the OC-SORT missed the fast-moving cars with ID 42. In summary, our FOLT performs better than the state-of-the-arts in objects with small sizes, large and irregular motion, and blurred appearance.

5 CONCLUSION

Multiple objects tracking in UAV view face the difficulty of large and irregular motion of both ground objects and UAV platforms. The small object size and blurred object appearance also hinder the tracking process. In this paper, we propose FOLT to address these challenging problems. The FOLT uses a modern light-weight detector and optical flow estimator to extract object detection features and motion information with high efficiency. Given these extracted features, FOLT introduces the flow-guided feature augmentation (FGFA) to augment detection features based on previous features and optical flow, which improves the detection of small objects and blurred objects. We then propose flow-guided motion prediction (FGMP) which utilizes optical flow and a convolution layer to track objects with large and irregular motion more precisely. Experiments show that both FGFA and FGMP can improve the accuracy of multiple object tracking in UAV view, and the combination of the two methods further improves the accuracy and achieves the best results. Comparison with state-of-the-arts shows that FOLT achieves the best results on the two public UAV-MOT datasets in both tracking accuracy and inference speed.

ACKNOWLEDGEMENT

This work was supported in part by Natural Science Foundation of China under contract 62171139, and in part by Zhongshan science and technology development project under contract 2020AG016.

REFERENCES

- [1] Seung-Hwan Bae and Kuk-Jin Yoon. 2017. Confidence-based data association and discriminative deep appearance learning for robust online multi-object tracking. *IEEE transactions on pattern analysis and machine intelligence* 40, 3 (2017), 595–610.
- [2] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft. 2016. Simple online and realtime tracking. In *ICIP*. IEEE, 3464–3468.
- [3] Erik Bochinski, Volker Eiselein, and Thomas Sikora. 2017. High-speed tracking-by-detection without using image information. In *AVSS2017*. IEEE, 1–6.
- [4] Thomas Brox, Bodo Rosenhahn, Daniel Cremers, and Hans-Peter Seidel. 2006. High accuracy optical flow serves 3-D pose tracking: exploiting contour and flow based constraints. In *ECCV*. Springer, 98–111.
- [5] Jinkun Cao, Xinhao Weng, Rawal Khrodar, Jiangmiao Pang, and Kris Kitani. 2022. Observation-Centric SORT: Rethinking SORT for Robust Multi-Object Tracking. *arXiv preprint arXiv:2203.14360* (2022).
- [6] Sushil Chandra, Greeshma Sharma, Saloni Malhotra, Devendra Jha, and Alok Prakash Mittal. 2015. Eye tracking based human computer interaction: Applications and their uses. In *2015 International Conference on Man and Machine Interfacing (MAMI)*. IEEE, 1–5.
- [7] S. Chen, Y. Xu, X. Zhou, and F. Li. 2019. Deep Learning for Multiple Object Tracking: A Survey. *IET Computer Vision* 13, 4 (2019), 355–368.
- [8] Patrick Dendorfer, Hamid Rezatofighi, Anton Milan, Javen Shi, Daniel Cremers, Ian Reid, Stefan Roth, Konrad Schindler, and Laura Leal-Taixé. 2020. Mot20: A benchmark for multi object tracking in crowded scenes. *arXiv preprint arXiv:2003.09003* (2020).
- [9] Alexey Dosovitskiy, Philipp Fischer, Eddy Ilg, Philip Hausser, Caner Hazirbas, Vladimir Golkov, Patrick Van Der Smagt, Daniel Cremers, and Thomas Brox. 2015. FlowNet: Learning optical flow with convolutional networks. In *Proceedings of the IEEE international conference on computer vision*. 2758–2766.
- [10] Dawei Du, Yuanqi Qi, Hongyang Yu, Yifan Yang, Kaiwen Duan, Guorong Li, Weigang Zhang, Qingming Huang, and Qi Tian. 2018. The unmanned aerial vehicle benchmark: Object detection and tracking. In *ECCV*. 370–386.
- [11] Loïc Fagot-Bouquet, Romaric Audigier, Yoann Dhome, and Frédéric Lerasle. 2016. Improving multi-frame data association with sparse representations for robust near-online multi-object tracking. In *European Conference on Computer Vision*. Springer, 774–790.
- [12] David A Forsyth, Okan Arikan, Leslie Ikemoto, James O'Brien, Deva Ramanan, et al. 2006. Computational studies of human motion: Part 1, tracking and motion synthesis. *Foundations and Trends® in Computer Graphics and Vision* 1, 2–3 (2006), 77–254.
- [13] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. 2021. YoloX: Exceeding yolo series in 2021. *arXiv preprint arXiv:2107.08430* (2021).
- [14] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. 2013. Vision meets robotics: The kitti dataset. *The International Journal of Robotics Research* 32, 11 (2013), 1231–1237.
- [15] Andreas Geiger, Philip Lenz, and Raquel Urtasun. 2012. Are we ready for autonomous driving? the kitti vision benchmark suite. In *2012 IEEE conference on computer vision and pattern recognition*. IEEE, 3354–3361.
- [16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [17] Berthold KP Horn and Brian G Schunck. 1981. Determining optical flow. *Artificial intelligence* 17, 1–3 (1981), 185–203.
- [18] Junhwa Hur and Stefan Roth. 2019. Iterative residual refinement for joint optical flow and occlusion estimation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5754–5763.
- [19] Joel Janai, Fatma Güney, Aseem Behl, Andreas Geiger, et al. 2020. Computer vision for autonomous vehicles: Problems, datasets and state of the art. *Foundations and Trends® in Computer Graphics and Vision* 12, 1–3 (2020), 1–308.
- [20] L. Kalake, W. Wan, and L. Hou. 2021. Analysis Based on Recent Deep Learning Approaches Applied in Real-Time Multi-Object Tracking: A Review. *IEEE Access* PP, 99 (2021), 1–1.
- [21] Rudolph Emil Kalman. 1960. A new approach to linear filtering and prediction problems. (1960).
- [22] Lingtong Kong, Chunhua Shen, and Jie Yang. 2021. FastflowNet: A lightweight network for fast optical flow estimation. In *ICRA*. IEEE, 10310–10316.
- [23] Louis Kratz and Ko Nishino. 2010. Tracking with local spatio-temporal motion patterns in extremely crowded scenes. In *2010 IEEE computer society conference on computer vision and pattern recognition*. IEEE, 693–700.
- [24] Cheng-Hao Kuo, Chang Huang, and Ramakant Nevatia. 2010. Multi-target tracking by on-line learned discriminative appearance models. In *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, 685–692.
- [25] Eero Kurimo, Leena Lepistö, Jarno Nikkanen, Juuso Grén, Iivari Kuutti, and Jorma Laaksonen. 2009. The effect of motion blur and signal noise on image quality in low light imaging. In *Image Analysis: 16th Scandinavian Conference, SCIA 2009, Oslo, Norway, June 15–18, 2009. Proceedings* 16. Springer, 81–90.
- [26] Laura Leal-Taixé, Anton Milan, Ian Reid, Stefan Roth, and Konrad Schindler. 2015. Motchallenge 2015: Towards a benchmark for multi-target tracking. *arXiv preprint arXiv:1504.01942* (2015).
- [27] Xi Li, Weiming Hu, Chunhua Shen, Zhongfei Zhang, Anthony Dick, and Anton Van Den Hengel. 2013. A survey of appearance models in visual object tracking. *ACM transactions on Intelligent Systems and Technology (TIST)* 4, 4 (2013), 1–48.
- [28] Shuai Liu, Xin Li, Huchuan Lu, and You He. 2022. Multi-Object Tracking Meets Moving UAV. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8876–8885.
- [29] Wenhan Luo, Junliang Xing, Anton Milan, Xiaoqin Zhang, Wei Liu, and Tae-Kyun Kim. 2021. Multiple object tracking: A literature review. *Artificial Intelligence* 293 (2021), 103448.
- [30] T. Meinhardt, A. Kirillov, L. Leal-Taixe, and C. Feichtenhofer. 2021. TrackFormer: Multi-Object Tracking with Transformers. (2021).
- [31] Anton Milan, S Hamid Rezatofighi, Anthony Dick, Ian Reid, and Konrad Schindler. 2017. Online multi-target tracking using recurrent neural networks. In *AAAI*, Vol. 31.
- [32] Jinlong Peng, Yueyang Gu, Yabiao Wang, Chengjie Wang, Jilin Li, and Feiyue Huang. 2020. Dense scene multiple object tracking with box-plane matching. In *ACM International Conference on Multimedia*. 4615–4619.
- [33] Jinlong Peng, Fan Qiu, John See, Qi Guo, Shaoshuai Huang, Ling-Yu Duan, and Wei Yao Lin. 2018. Tracklet siamese network with constrained clustering for multiple object tracking. In *VCIP*. IEEE, 1–4.
- [34] Jinlong Peng, Changan Wang, Fangbin Wan, Yang Wu, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. 2020. Chained-tracker: Chaining paired attentive regression results for end-to-end joint multiple-object detection and tracking. In *ECCV*. Springer, 145–161.
- [35] Jinlong Peng, Tao Wang, Wei Yao Lin, Jian Wang, John See, Shilei Wen, and Erui Ding. 2020. TPM: Multiple object tracking with tracklet-plane matching. *Pattern Recognition* 107 (2020), 107480.
- [36] Hamed Pirsiavash, Deva Ramanan, and Charless C Fowlkes. 2011. Globally-optimal greedy algorithms for tracking a variable number of objects. In *CVPR 2011*. IEEE, 1201–1208.
- [37] Bing Shuai, Andrew Berneshawi, Xinyu Li, Davide Modolo, and Joseph Tighe. 2021. SiamMOT: Siamese Multi-Object Tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 12372–12382.
- [38] Karen Simonyan and Andrew Zisserman. 2014. Two-stream convolutional networks for action recognition in videos. *Advances in neural information processing systems* 27 (2014).
- [39] Daisuke Sugimura, Kris M Kitani, Takahiro Okabe, Yoichi Sato, and Akihiro Sugimoto. 2009. Using individuality to track individuals: Clustering individual trajectories in crowds using local appearance and frequency trait. In *ICCV*. IEEE, 1467–1474.
- [40] Deqing Sun, Stefan Roth, and Michael J Black. 2014. A quantitative analysis of current practices in optical flow estimation and the principles behind them. *International Journal of Computer Vision* 106, 2 (2014), 115–137.
- [41] Deqing Sun, Xiaodong Yang, Ming-Yu Liu, and Jan Kautz. 2018. Pwc-net: Cnns for optical flow using pyramid, warping, and cost volume. In *CVPR*. 8934–8943.
- [42] Valtteri Takala and Matti Pietikainen. 2007. Multi-object tracking using color, texture and motion. In *2007 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 1–7.
- [43] Zhongdao Wang, Liang Zheng, Yixuan Liu, Yali Li, and Shengjin Wang. 2020. Towards real-time multi-object tracking. In *European Conference on Computer Vision*. Springer, 107–122.
- [44] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. 2017. Simple online and realtime tracking with a deep association metric. In *ICIP*. IEEE, 3645–3649.
- [45] Yihong Xu, Aljosa Osep, Yutong Ban, Radu Horaud, Laura Leal-Taixé, and Xavier Alameda-Pineda. 2020. How to train your deep multi-object tracker. In *CVPR*. 6787–6796.
- [46] Bo Yang and Ram Nevatia. 2012. Multi-target tracking by online learning of non-linear motion patterns and robust appearance models. In *CVPR*. IEEE, 1918–1925.
- [47] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. 2022. ByteTrack: Multi-object tracking by associating every detection box. In *ECCV*. Springer, 1–21.
- [48] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. 2021. Fairmot: On the fairness of detection and re-identification in multiple object tracking. *International Journal of Computer Vision* 129, 11 (2021), 3069–3087.
- [49] Xuemei Zhao, Dian Gong, and Gérard Medioni. 2012. Tracking using motion patterns for very crowded scenes. In *ECCV*. Springer, 315–328.
- [50] Liang Zheng, Zhi Bie, Yifan Sun, Jingdong Wang, Chi Su, Shengjin Wang, and Qi Tian. 2016. Mars: A video benchmark for large-scale person re-identification. In *European conference on computer vision*. Springer, 868–884.
- [51] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. 2020. Tracking objects as points. In *European Conference on Computer Vision*. Springer, 474–490.
- [52] Pengfei Zhu, Longyin Wen, Dawei Du, Xiao Bian, Heng Fan, Qinghua Hu, and Haibin Ling. 2020. Detection and tracking meet drones challenge. *arXiv preprint arXiv:2001.06303* (2020).