

Real-time and Intelligent Moving Targets Tracking based on UAV Remote Sensing Video Camera and Brain-like Computing Chips

Kun Hu^{1*}, Yuxuan Wu¹, Mingyang Huang¹, Chen Wang¹, Haoheng Wu¹, Tianyu Cai², Qingle Zhang¹, Shichao Wang³, Bo Li³

¹Institute of Artificial Intelligence, Beihang University, 100191, Beijing, China - (kunhu, buaawyx, 22373155, 22421008, 22421009, qlzhang)@buaa.edu.cn

²The University of Sydney, Camperdown NSW 2050, Sydney, Australia - tcai7097@uni.sydney.edu.au

³Geovis Technology Co.,Ltd, 101399, Beijing, China - (wangsc01, lib)@geovis.com.cn

Keywords: Remote sensing, Moving target tracking, Real-time, Video camera, Brain-inspired computing chips.

Abstract

Moving target tracking technology based on Unmanned Aerial Vehicles (UAV) is widely used in many fields such as automatic inspection and emergency response. The existing moving target tracking methods usually have the problems of large computation and low tracking efficiency. Limited by the computing power of the UAV platform, real-time tracking and analysis of multiple targets based on the video data collected by UAV platform is a difficult task. In this paper, we proposed a novel Target Specific Filtering Tracking with Memory (TSFMTrack) method designed for UAV-based real-time tracking tasks, which involves a Tracklet Filtering Module (TFM) for capturing object appearance features and a Tracklet Matching Module (TMM) for bounding box association in each frame. By experimental comparison with other State-Of-The-Art (SOTA) methods on popular MOT and UAV tracking datasets, the TSFMTrack have shown obvious advantages in accuracy, computational efficiency and reliability. Furthermore, we deployed the TSFMTrack on the brain-inspired chip Lynchip KA200, the experimental results have shown that the TSFMTrack is effective on edge computational platform and suitable for UAV real-time tracking tasks.

1. Introduction

In recent years, there has been rapid development in Unmanned Aerial Vehicle (UAV) remote sensing technology. Videos captured by the UAVs are used for intelligent targets tracking analysis. With UAV platforms featuring unique advantages such as compact size, agile maneuverability, and enhanced safety features, object tracking analysis based on which has found ubiquitous utility in diverse areas such as emergency response, traffic management, factory inspection, and so on. However, there are still several major challenges for real-time and intelligent moving targets tracking algorithms. In complex actual situations, the following obstructive factors hinder the realizing of real-time and accurate moving target tracking.

- **Limited energy and computational resources.** Limitations in power supply and payload significantly constrain the speed at which real-time processing and analysis of UAV remote sensing images can occur. Also, batteries have long been a barrier, but tethered systems can help compensate the weakness, allowing flights of several hours. To achieve real-time, accurate, and robust motion target tracking, algorithms must strike a balance between accuracy and efficiency. Meanwhile, it should be ensured a sufficient lightweight design to conserve energy for other energy-consuming controlling functions in complex environments, thereby enabling the collection of more geographic information during each flight.
- **Influence of camera motion.** UAV-mounted cameras exhibit fast movement and continuous angle changes, resulting in images reflecting the relative motion between the UAV and ground objects. Failure to correct this can lead to significant errors in target trajectory prediction during

motion target tracking. Additionally, UAVs commonly encounter mechanical vibrations during flight, particularly in strong winds, which can result in motion blur, making it difficult to obtain clear information about the appearance and motion characteristics of targets.

- **Viewpoint changes.** During sampling, UAVs often fly around objects, capturing different sides of 3D ground objects, leading to diverse changes in object appearance. If without timely online learning and model updates, trackers may misjudge target trajectories or even lose track of targets.
- **Low image resolution.** The large visual range of UAVs results in background information being insufficient, leading to reduced object resolution in captured images and weakened model representation. This diminished representation can impair tracker discriminative abilities, ultimately resulting in tracking failures.
- **Illumination variations and visual occlusion.** The lighting conditions for UAVs can change rapidly, ranging from bright to dim environments or transitioning between indoor, canopy, shadowed, and sunlit areas. Furthermore, UAVs frequently encounter complex and poorly lit natural environments during flight, such as nighttime, rainy, or foggy conditions, making it challenging for trackers to distinguish objects from the background. Also, partial or complete occlusion can hinder obtaining information about objects, making it easy to lose track of them.

In terms of multiple object tracking (MOT), the Tracking-by-Detection (TbD) paradigm is one of the mainstream approaches. Comprised of detecting phase and tracking phase, this paradigm aims to first determine the locations of various

* Corresponding author

targets and then correlate them between frames, generating estimated tracks of targets.

For TbD trackers, the performance of detection algorithms are crucial, with notable contributions coming from the YOLO series (Glenn, 2022, Glenn, 2024, Ge et al., 2021). These real-time detectors leverage anchor-based convolutional neural networks (CNN) to solve the detection problem through regression, thus achieving remarkable inference speed with relatively high accuracy. Complementing detection, object tracking techniques have witnessed significant advancements, for instance, SORT (Wojke et al., 2017), DeepSORT (Pujara and Bhamare, 2022) and their variants (Cao et al., 2023, Maggolino et al., 2023, Aharon et al., 2022, Zhang et al., 2022). These methods merge Kalman Filters with advanced trajectory matching algorithms along with CNNs to enhance tracking robustness, particularly in scenarios characterized by occlusions and nonlinear motion dynamics.

Moreover, recent advancements in attention mechanisms and correlation filter-based approaches offer promising improvement inspirations for enhancing tracking accuracy and real-time performance. TrackFormer (Meinhardt et al., 2022) and correlation filter-based trackers like MOSSE (Bolme et al., 2010) and ECO (Danelljan et al., 2017), leverage attention mechanisms and Discrete Fourier Transformation respectively to tackle complex tracking scenarios while maintaining low computational complexity.

This study pursues to cope with the challenges of real-time and intelligent moving targets tracking based on UAV remote sensing video cameras, especially the challenges concerning the unique environment of UAV platforms, and further delineates potential directions to guide the progression of research in UAV-based moving target tracking. The primary contributions of this work can be summarized into the following aspects.

In Section 2, we first conducted a comprehensive review of related existing literature in the field. The core idea and methodological researches are explained in Section 3. The integration of the developed tracker into a cohesive code repository facilitates accessibility and reproducibility. Experimental evaluation and onboard testing of our proposed method is included in Section 4, where our TSFMTTrack is deployed on Lynchip KA200 brain-inspired chip for comprehensive inference testing of UAV remote sensing videos. Experiments are undertaken on four authoritative UAV benchmark datasets, namely MOT17 (Sun et al., 2019), MOT20 (Dendorfer et al., 2020) and UAVDT (Du et al., 2018), to comprehensively assess the performance of the TSFMTTrack in complex scenarios. SOTA target tracking algorithms such as ByteTrack (Zhang et al., 2022) and BoT-SORT (Aharon et al., 2022) were deployed on the same chip for experimental comparison.

2. Related work

2.1 Tracking-by-Detection

There are two main approaches to Multiple Object Tracking (MOT): tracking by detection (TbD) and joint detection and tracking (JDT), with the former being widely used due to its simplicity and modularization. Generally, the TbD method can be divided into two parts: object detection and object tracking.

The YOLO series (Redmon et al., 2016, Glenn, 2022, Glenn, 2024) are commonly used real-time detectors in MOT, outperforming their counterparts in speed and accuracy by modeling

the detection problem as a regression problem and introducing anchor-based CNN to solve it. YOLOX (Ge et al., 2021) removes prior anchors in YOLO and adds other techniques to reduce hyper-parameters and computational cost while still achieving promising performances.

In object tracking, SORT (Wojke et al., 2017) combines Kalman Filter and Hungarian matching algorithm to create a simple yet effective tracker. Building on SORT, DeepSORT (Pujara and Bhamare, 2022) adds a CNN module to learn visual features, improving the tracker's performance in occlusion scenarios. OC-SORT (Cao et al., 2023) employs observation-centric compensation methods to deal with the error accumulation of Kalman filtering in nonlinear motion scenarios. However, OC-SORT's high reliance on image quality results in less effective performance in practical applications. Addressing this issue, Deep-OC-SORT (Maggolino et al., 2023) combines the aforementioned trackers and adds correction terms about objects' appearances to tackle feature degradation. To overcome the limitations of SORT-like trackers, BoT-SORT (Aharon et al., 2022) combines motion and appearance information to optimize bounding box direction. Bytetrack (Zhang et al., 2022) associates almost every detection box to minimize mismatching while maintaining a high running speed.

Although all the aforementioned methods have sound and promising outcomes, most of them rely on high-performance GPUs. In UAV tracking scenarios, computational resources are strictly limited due to the UAV's payload capacity. Additionally, training deep networks suitable for UAV tracking requires a large number of UAV-based datasets, which are currently insufficient to support this need.

2.2 Accurate Tracking with Attention Mechanism

The attention mechanism (Vaswani et al., 2017) has also demonstrated its capability in object tracking. TrackFormer (Meinhardt et al., 2022) models the tracking task as a prediction problem. Using attention in association and encoder-decoder structures to predict tracklets, it outperformed many state-of-the-art traditional trackers in accuracy. TransMOT effectively models relations between a large number of objects mainly through a spatial-temporal graph transformer structure. SMILETrack (Wang et al., 2024) incorporates a Siamese network to capture appearance features. By employing Patch Self-Attention mechanisms, SMILETrack effectively attends to image similarity and enhances performance in the presence of occlusions. Despite attention mechanism still has the problems that 2.1 mentioned, its idea of focusing the major feature is still worth considering in increasing real-time trackers' accuracy.

2.3 Real-time Tracking with Correlation Filter

Correlation Filters (CF) have garnered much attention in UAV-based tracking due to their adaptability, efficiency, and relatively high resilience against background occlusion. One key highlight of correlation filters is that, by applying Discrete Fourier Transformation, they transform cyclic correlation (done by convolution) into element-wise multiplication. This operation significantly reduces computational complexity, allowing CF-based trackers to reach over 24 frames per second (FPS) on a single CPU, meeting real-time requirements for UAV-based tracking.

MOSSE (Bolme et al., 2010) was the first to use CF in object tracking, introducing a minimum squared error regularization method that produced a robust and stable CF tracker. CSK

(Henriques et al., 2012) introduced a circulant matrix to compute cyclic correlation. KCF (Henriques et al., 2015) summarized CSK and reformulated the CF-tracking algorithm to a Kernelized Correlation Filter, with complexity equivalent to linear algorithms. CCOT (Danelljan et al., 2016) introduced implicit interpolation to integrate multi-resolution deep feature maps. However, the key formula in CCOT incurred high computation costs. Additionally, training continuous filters would introduce numerous optimized parameters, leading to overfitting. ECO (Danelljan et al., 2017) was proposed to address these issues. It introduced a factorized convolution operator to build the filter in CCOT, while a compact generative model of the training sample distribution decreased computational costs. MCPF (Zhang et al., 2017) incorporates a particle filter and Multi-task Correlation Filter to handle large-scale variation. CSR-DCF (Lukežič et al., 2018) managed to learn accurate features of irregular objects by introducing a spatial confidence map and channel-wise confidence score. Li et al. (Li et al., 2020b) proposed a discriminative correlation filter (DCF) with a memory queue to preserve keyframes' information, enabling long-term tracking with robustness. Considering the reversibility of motion, BiCF (Lin et al., 2020) adds bidirectional incongruity terms in training to ensure the filter's consistency in forward and backward motion prediction. Also, AutoTrack (Li et al., 2020a) incorporates automatic spatial-temporal regularization by integrating local and global response maps to dynamically regulate spatial and temporal weights, ensuring adaptability across diverse sequences while maintaining computational efficiency.

Correlation filter generates promising outcomes when applied in Single Object Tracking tasks, but its structure impedes its performance on MOT task. By decoupling correlational operation from CF and deploy it in solving MOT task might create a tracker with higher efficiency.

3. Methodology

In this section, a novel detector for real-time UAV MOT task, Target Specific Filtering Track with Memory (TSFMTrack), is presented. The structure of which is illustrated in Fig 1.

To be specific, our TSFMTrack consists of two parts: object detection and tracklet matching. We apply YOLOv8 (Glenn, 2024) for detecting due to its wide range of usage scenarios with better results than previous YOLO detector. The main contribution of our work are mainly in the tracklet matching part, which comprises (a) A Siamese-like Target Filtering Module (TFM) for accurately learns the features and computes similarity score, and (b) A Tracklet Matching Module (TMM) assigns and upgrades the tracklets using Hungarian algorithm.

3.1 Target Filtering Module

To achieve promising tracking quality with high efficiency, a well-designed feature extractor, TFM, is proposed. Though Correlation Filter-based trackers achieves low inference time, they only generates one filter for a frame and thus, not suitable for MOT tasks. Also, since they updates their filters online, their tracking quality are not as good as deep-learning based trackers. Siamese network based trackers integrates template information into the searching region and is suitable in processing multiple similar inputs.

The proposed TFM utilizes the advantages of Siamese Network and Correlation operation to precisely and effectively learn the

discriminative appearance features to for accurate tracking. To extend CF's success to MOT, we decoupled correlation operation from CF and takes the last fragment of a tracklet as a specified filter for this object to conduct correlation operation.

Figure 2 shows the TFM's architecture. It takes the last fragment of a tracklet of the previous frame and bounding boxes of the current frame as inputs, after preprocessed by CNN, the inputs would then be processed by a Correlation Computing Block (CCB). Finally, another CNN is used to calculate the correlation index between tracklets and targets before calculating similarity score between the two inputs.

3.1.1 Correlation Computing Block Use \odot to denote Hadamard multiplication and $*$ to indicate the complex conjugate, then what Correlation Filter does can be generalized as computing the following formula:

$$G = F \odot H^* \quad (1)$$

where F, G, H respectively denotes the 2D Discrete Fourier Transform of the input image, response map and the filter. The output of which is then transformed back to spatial domain by Inverse Discrete Fourier Transform to get the response map. Then the object's current region can be found by searching the maximum response of the map. Here, instead of generating the filter by algorithm, we directly takes the tracking targets as the specified correlation filter for its tracklets and filters the detected targets to get the response map. To take dissimilarity between filters and the unmatched images (i.e. responses generated by the filter and targets not belong to this tracklet) into consideration, the filter itself would be passed to the next procedure after doing the same filtering process. Since the input images are of different sizes, it will first be resized to a fixed size $W \times W$ by CNN before being processed by CCB, where W is set to 127 in our work since its a prime number and is close to 2^7 , which makes it suitable to compute by FFT and would not cause information erosion (Wu et al., 2019) on feature map after correlation computation.

3.1.2 Back Propagation Through CCB Module Back propagation for training TFM requires a differentiable or piecewise differentiable forward function for each part of TFM. In this section, we will prove that the equivalent function of CCB Module is differentiable and calculate the gradient through it. The function of this back propagation is illustrated in Figure 3.

Since the operation CCB Module done are all matrix-wise without cross-channel calculation, the prove can be done under 2D tensor (i.e. matrix) condition and be easily generalized to 3D tensor condition. The bellowing matrix are all $(2n + 1) \times (2n + 1)$ matrix.

For one channel of an image, the process of conducting Fast Fourier Transform and Inverse Fast Fourier Transform is in fact implementing circular convolution. Let Ω represents the kernel matrix for circular convolution, X represents the input matrix and $X(i, j)$ represents the (i, j) -ist element of X . The circular convolution of Ω and X can be written as $F = \Omega \circledast X$, where \circledast stands for circular convolution symbol. Other related values are represented as the equation below.

$$F(i, j) = \sum_{u=-n}^n \sum_{v=-n}^n \Omega(u, v) \cdot X(i + u, j + v) \quad (2)$$

$i, j, u, v, n \in \mathbb{Z}_{2n+1}$

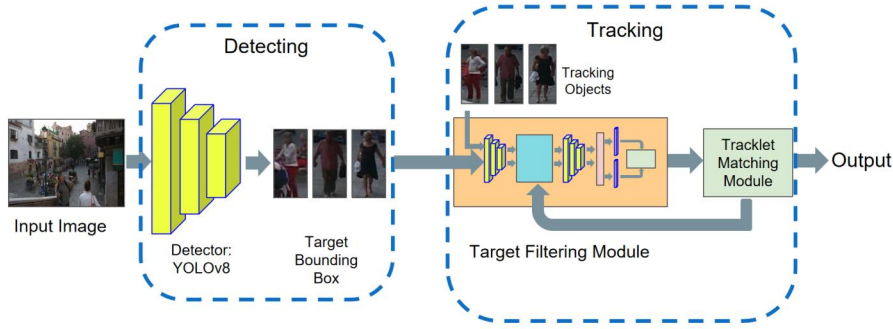


Figure 1. Overall architecture of TSFMTrack

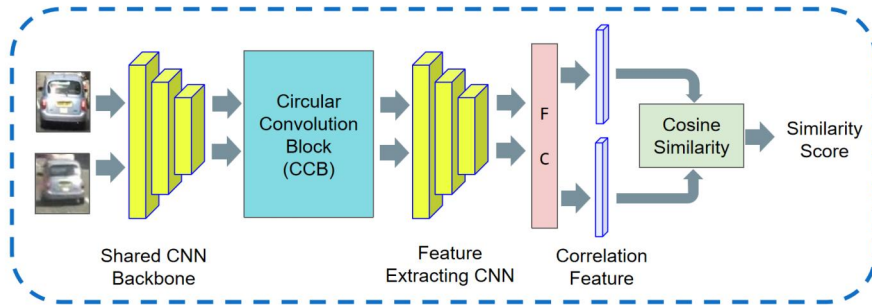


Figure 2. TFM calculates correlational similarity between detected objects at the current frame and tracklets at the previous frame

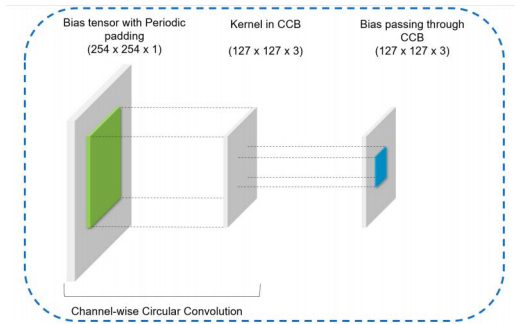


Figure 3. Illustration of back propagation through CCB Module

$$\Omega = \begin{pmatrix} \Omega(-n, -n) & \dots & \Omega(-n, n) \\ \dots & \Omega(0, 0) & \dots \\ \Omega(n, -n) & \dots & \Omega(n, n) \end{pmatrix}$$

$$X = \begin{pmatrix} X(0, 0) & \dots & X(0, 2n) \\ \dots & \dots & \dots \\ X(2n, 0) & \dots & X(2n, 2n) \end{pmatrix}$$

\mathbb{Z}_{2n+1} is the integer ring modulo $2n+1$. Circular convolution is done by element-wise multiplication and addition, so it is easy to prove that $F(i, j)$ is differentiable with respect to $X(k, l)$.

Let X and Y respectively denote the tracklet and the bounding box image, the equivalent function of CCB Module can be written as $F(X, Y) = (\Omega \otimes X, \Omega \otimes Y)$. For clearer representation and narration, we rewrite the above function as $F(X, Y) = (F_1, F_2)$, $F_1 = \Omega \otimes X$, $F_2 = \Omega \otimes Y$.

In formulating the process of back propagation, we denote the bias matrix that feature extracting CNN returned as (\hat{B}_1, \hat{B}_2) . \hat{B}_1, \hat{B}_2 are bias matrices corresponding to the input F_1, F_2 and its element represents the partial derivative of Loss function \mathcal{L} with respect to F , i.e. $\hat{B}(i, j) = \frac{\partial \mathcal{L}}{\partial F(i, j)}$.

Let B_1 denote the bias matrix passed back to the CNN backbone through the CCB Module with respect to \hat{B}_1 , then $B_1(i, j)$ can be calculated as

$$B_1(i, j) = \frac{\partial \mathcal{L}}{\partial X(i, j)}$$

$$= \sum_{\hat{i}=0}^{2n} \sum_{\hat{j}=0}^{2n} \frac{\partial \mathcal{L}}{\partial F_1(\hat{i}, \hat{j})} \cdot \frac{\partial F_1(\hat{i}, \hat{j})}{\partial X(i, j)}$$

$$= \sum_{\hat{i}=0}^{2n} \sum_{\hat{j}=0}^{2n} \hat{B}_1(\hat{i}, \hat{j}) \cdot \frac{\partial F_1(\hat{i}, \hat{j})}{\partial X(i, j)}$$

$i, j, \in \mathbb{Z}$

where \mathbb{Z} is the ring of integer.

Apply homomorphic transformation

$$\mathbb{Z} \rightarrow \mathbb{Z}_{2n+1}$$

$$z \rightarrow z'$$

(4)

to the subscript of Equation 3, and the result of which remains unchanged.

Then, associate Equation 2 with Equation 3, we get

$$\begin{aligned}
 B_1(i, j) &= \sum_{\hat{i}=0}^{2n} \sum_{\hat{j}=0}^{2n} \hat{B}_1(\hat{i}, \hat{j}) \frac{\partial}{\partial X(i, j)} \cdot \\
 &\quad \left(\sum_{u=-n}^n \sum_{v=-n}^n \Omega(u, v) \cdot X(\hat{i} + u, \hat{j} + v) \right) \\
 &= \sum_{\hat{i}=0}^{2n} \sum_{\hat{j}=0}^{2n} \sum_{u=-n}^n \sum_{v=-n}^n \hat{B}_1(\hat{i}, \hat{j}) \cdot \\
 &\quad \frac{\partial}{\partial X(i, j)} \left(\Omega(u, v) \cdot X(\hat{i} + u, \hat{j} + v) \right) \\
 &\quad i, j, u, v, n \in \mathbb{Z}_{2n-1}
 \end{aligned} \tag{5}$$

Notice that $\frac{\partial}{\partial X(i, j)} \left(\Omega(u, v) \cdot X(\hat{i} + u, \hat{j} + v) \right) \neq 0$ only when $\hat{i} + u \equiv i \pmod{2n+1}$ and $\hat{j} + v \equiv j \pmod{2n+1}$, Equation 5 can be simplified as

$$\begin{cases} B_1(i, j) = \sum_{\hat{i}=0}^{2n} \sum_{\hat{j}=0}^{2n} \hat{B}_1(\hat{i}, \hat{j}) \cdot \Omega(u, v) \\ \hat{i} + u \equiv i \pmod{2n+1} \\ \hat{j} + v \equiv j \pmod{2n+1} \end{cases} \tag{6}$$

that is

$$\begin{aligned}
 B_1(i, j) &= \sum_{u=-n}^n \sum_{v=-n}^n \hat{B}_1(i - u, j - v) \cdot \Omega(u, v), \\
 &\quad i, j, u, v, n \in \mathbb{Z}_{2n-1}
 \end{aligned} \tag{7}$$

Let $u' = -u, v' = -v$, then

$$\begin{aligned}
 B_1(i, j) &= \sum_{u'=-n}^n \sum_{v'=-n}^n \hat{B}_1(i + u', j + v') \cdot \Omega(-u', -v'), \\
 &\quad i, j, u', v', n \in \mathbb{Z}_{2n-1}
 \end{aligned} \tag{8}$$

i.e.

$$B_1 = \hat{B}_1 \otimes \text{Rotate}_\pi(\Omega) \tag{9}$$

in which $\text{Rotate}_\pi(X)$ stands for rotating the matrix X for π rad. This is equal to PXQ , where

$$P = Q = \begin{pmatrix} & & & & 1 \\ & & & & \\ & & & & \\ & & & & \\ 1 & & & & \end{pmatrix} \tag{10}$$

Equation (9) represents the circular convolution operation between the bias matrix \hat{B}_1 and the rotated convolutional kernel Ω , which can be efficiently computed by CCB.

3.2 Tracklet Matching Module

Tracklet matching is also an indispensable step in object tracking. Matching the tracklets properly can have a positive impact on tracking outcomes. Our tracking algorithm is built upon the

recent SoTA tracker SMILETrack, which is the extension of ByteTrack. By associating almost every tracklet using a two-stage matching strategy, these two trackers outperformed their counterparts in overall performances. However, when practically deployed on UAV, tracking algorithms would encounter challenges from unpredictable conditions, e.g. camera motion and occlusion, that affect the tracker's tracking quality. To address these exterior condition, we apply Camera Motion Compensation to accurately track objects in interfered scenarios.

Also, we introduce Key Frame Feature Memory Queue to introduce historical views to the tracker, allowing it to adapt to appearance changing and reduces the interference of transient intense interference. We design our TMM method to address the aforementioned problem and integrates the overall matching pipeline in SMILETrack, achieving a interference-robust accurate tracking strategy.

Let O, T, S denote the set of objects, the set current tracklet list and the matrix of object-tracklet similarity respectively. Elements in set O are sorted in a descending order, any element O_i with a detection score lower than 0.1 would be considered as background noise and be removed before matching. Then we divide O into O^H and O^L , respectively denotes elements with a detection score above/below the median score. The similarity score between the i -th object and the j -th tracklet, i.e. $S(i, j)$ can be calculated by

$$S(i, j) = S_{DIOU}(I, J) + \alpha S_{corr}(i, j) \tag{11}$$

where S_{DIOU} is the DIOU similarity and S_{corr} is the output of TFM with the input of O_i and the tracklet T_i .

Using O, T, S , the mainstream of TMM can be stated as:

- **Stage 1.** Finds the matches between O^H and T . We first predict every tracklets' new position in the using Kalman filter, then the Hungarian algorithm is applied to perform linear assignment using the similarity matrix S^H . The unmatched objects of O^H and the unmatched tracklet of T are then placed in O_{Remain}^H and T_{Remain}^H .
- **Stage 2.** Match the objects in O^L and T_{Remain}^H . The unmatched objects O_{Remain} and tracklets T_{Remain} would pass a gate function before further operation.

In summary, the pseudo-code of TMM can be stated as Algorithm 1.

Moreover, for the training process of the two convolutional networks, we employed a variant of DIOU metric (Zheng et al., 2019) for similarity calculation, and trained the parameters in the cnns via computing the L2 loss.

The measurement of similarity can be represented as:

$$\mathcal{L}_{DIOU} = 1 - IoU + \frac{\rho^2(B_1, B_2)}{c^2} \tag{12}$$

in which $\rho^2(B_1, B_2)$ stands for the distance of the central points of B_1 and B_2 , c is the diagonal length of the smallest enclosing box covering the two boxes, and IoU is calculated with

$$IoU = \frac{|B_1 \cap B_2|}{|B_1 \cup B_2|}$$

Algorithm 1: Framework of TMM

Input: Set of detected objects: O , Set of tracklets of the last frame: T , Set of Memory Queue for each tracklet: Q , Set of posterior state estimate vector and covariance matrix: $\{(\tilde{x}_i, P_i)\}$, Scaled rotation matrix M_t and translation parameter \vec{t} .

Output: Updated \hat{T} , \hat{Q} and $\{(\tilde{x}_i, \hat{P}_i)\}$.

```

while not at end of tracking do
     $(\tilde{x}_i, P_i) \leftarrow (M_t \tilde{x}_i + \vec{t}, M_t P_i M_t^T)$  // CMC
     $(\tilde{x}_i, \hat{P}_i) \leftarrow KF((\tilde{x}_i, P_i))$  // Kalman Filter
    /* Initialize Object list */
     $O \leftarrow O \setminus O_{lowscore}$ ;
    Split  $O$  to  $O^H$  and  $O^L$ ;
    /* Association 1 */
    Assigns  $O^H$  to tracklets;
     $O_{remain}^H \leftarrow$  unmatched objects in  $O^H$ ;
     $T_{remain}^H \leftarrow$  unmatched tracklets in  $T^H$ ;
     $T^H \leftarrow$  matched tracklets and objects;
    Update KFF Memory Queue;
    /* Association 2 */
    Assigns  $O^L$  to remaining tracklets;
    Delete unmatched objects in  $O^L$ ;
     $T_{remain}^L \leftarrow$  unmatched tracklets in  $T^L$ ;
     $T^L \leftarrow$  matched tracklets and objects;
    /* Update tracklet */
     $\hat{T} \leftarrow T^H \cup T^L$ ;
    for  $O_i$  in  $O_{remain}^H$  do
        if  $O_i.score > \tau$  then
             $\hat{T} \leftarrow \hat{T} \cup \{(O_i, T_{new})\}$  // New tracklet
        else
            delete  $O_i$ 
     $T_{matched}^L.keep \leftarrow 0$ ;
    for  $T_i$  in  $T_{remain}^L$  do
         $T_i.keep \leftarrow T_i.keep + 1$  if  $T_i.keep > 30$  then
            delete  $T_i$ 
    
```

Thus the loss function can be computed as:

$$L = \frac{1}{\xi \times \eta} \sum_{s=1}^{\eta} \sum_{t=1}^{\xi} (I_{s,t} - J_{s,t})^2 \quad (13)$$

In this scenario, I represents the bounding box of the input object, while J represents the output of tracklet matching. The similarity we've just calculated serves as an intermediary value approximating K , which isn't a direct output. This introduces a disparity compared to directly matching the tracklets, making it impractical to apply perceptual loss. Additionally, the limited parameters in our proposed method reduce the risk of overfitting, to the extent that regularization terms like L1 Loss may even hinder performance rather than enhance it. Generally speaking, employing L2 loss works most effectively under this circumstance.

3.2.1 Camera Motion Compensation Deep-learning based tracking-by-detection trackers rely heavily on the images' quality, which would be influenced by camera motion in real-life scenarios. In a dynamic camera situation like UAV, this phenomena would be prevalent and could result in increasing ID switches or false negatives.

The principle of CMC is visualized in Figure 4. We apply CMC

to correct the Kalman state following the formula down below:

$$\begin{cases} x'_k = M_k x_k + T_k \\ P'_k = M_k P_k M_k^T \end{cases} \quad (14)$$

Where x_k and x'_k , P_k and P'_k respectively denotes the KF's predicted state vector and covariance matrix before and after CMC operation. It is worth pointing out that the CMC update is done before the Kalman extrapolation step so that the prediction stage is from the CMC-corrected states, which could prevent error accumulation.

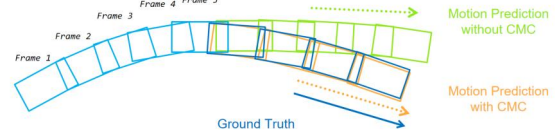


Figure 4. Motion prediction with/without CMC

3.2.2 Key Frame Feature Memory Queue Although the interference from the external environment to the image at a certain moment is random, it is known from the law of large numbers that the sum of multiple interference in a system approximately follows a normal distribution, so the overall interference to the image can be approximated as fluctuating within a small range over a long time span, so if the tracker can somehow gain temporal information while tracking, it should have suppresses and smooths the impact of noises.

To enable an access to such information, we introduce Key Frame Feature Memory Queue (KFF Memory Queue), making our tracker more temporal aware. Basically, a memory queue with the length of N is maintained for every active tracklets, storing fragments belongs to this tracklet. In motion prediction, the contribution of each frame are not equal. That is, it only need a small amount of frames to define and illustrate a given smooth motion, these frames are known as key frame. What the memory queue do is to find out such frames and adding those key frames in memory queue. In practice, key frames is mostly chosen as the starting and ending frame of a transition.

Since motions of tracked objects are often complex, it is necessary to add an additional frame in between starting and ending point of a motion. Also, following the Kalman Filter's prior hypothesis, motion can be seen as linear (or smooth) in a relatively long span of time, e.g. 1 second. Based on the above theory, we set the value of N as 6, for it could contain up to 2 seconds' key frame feature.

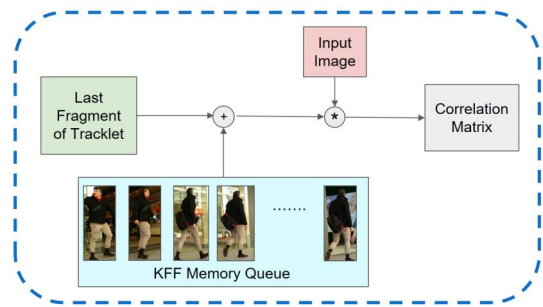


Figure 5. KFF Memory Queue + CCB structure. The enqueued frames are used to update filters

Enqueue and Dequeue. A frame will only enqueue only when it is considered as key frame, i.e. it is the turning point of two different motion or rapid change of appearance, etc. These criteria can be simplified as low similarity score. If a new added frame in a tracklet has a similarity score lower than a threshold τ , i.e. $S(i, j) < \tau$, it would be added into the memory queue of this tracklet. It is also worth pointing out that the enqueueing operation would only performed at stage 1 of TMM. If a queue is not updated for a given period of time t_0 , it would be considered as "inactive" and the queue would be removed.

Key frame in TFM. Once the fragment is enqueued, it is used to update the filter in CCB. Let f_k be the k th fragment in the queue (f_1 is the first enqueued fragment), X be the tracking targets, then a target specific temporal filter can be calculated by

$$\Omega = X \oplus \sum_{i=1}^N \lambda_i f_k \quad (15)$$

where \oplus is element-wise addition. Generally speaking, earlier enqueued frames have a lower similarity to the current appearance than latter ones, so a decaying factor needs to be introduced. For a given interference, its effect on the current frame roughly decays exponentially comparing to its initial intensity. For computational convenience, these two decaying factors are all implemented by multiplying the queue fragments $\{f_i\}$ by an exponentially increasing sequence $\{\lambda_i\}$. Without this Memory Queue, TMM cannot utilize temporal information to tackle exterior interference.

For straightforward understanding, the pseudo-code of KFF Memory Queue can be written as Algorithm 2.

Algorithm 2: Framework of KFF Memory Queue.

Input: Set of matched high-score Object and Tracklets:

$O_{matched}^H \times T_{matched}^H$, Similarity matrix: S ,
 Memory Queue Q .

Output: Updated \hat{Q} .

```

for ( $O_i, T_j$ ) in  $O_{matched}^H \times T_{matched}^H$  do
    if  $S(i, j) < s_0$  then
        if if not exist  $Q_j$  then
             $\perp$  create and initialize  $Q_j$ ;
         $\hat{Q}_j \leftarrow Q_j.enqueue(O_i)$ ;
        /* decay value of stored frames */
         $\hat{Q}_j.elements \leftarrow \lambda \times Q_j.elements$ ;
         $\hat{Q} \leftarrow (Q \setminus \{Q_j\}) \cup \{\hat{Q}_j\}$ ;
    
```

4. Experiments

4.1 Implementation Details

We implemented all the experiments using PyTorch and deployed our TSFMTTrack on the Lynchip KA200 brain-inspired computing chip to conduct comprehensive inference testing of UAV remote sensing videos. In terms of datasets, our experiments were conducted on MOT17, MOT20 and UAVDT benchmarks. Metrics such as MOTA (Bernardin and Stiefelwagen, 2008a), IDF1 (Ristani et al., 2016) and HOTA (Luiten et al., 2020) are employed during our experiments, highlighting identity matching. Our detector was initialized on MOT datasets

and fine-tuned on UAVDT datasets. In order to optimize the performance, data augmentation and an SGD optimizer with cosine annealing were applied. Our TMM module is designed to manage tracklets, and we assessed its key parameters in an ablation study.

Table 1. Benchmark evaluation experiment results of TSFMTTrack

Dataset	MOTA \uparrow	IDF1 \uparrow	HOTA \uparrow	FPS \uparrow
MOT20	73.3	75.0	60.1	20.2
MOT17	68.5	59.7	-	30.5
UAVDT	58.4	74.8	84.4	25.6

4.2 Datasets and Metrics

4.2.1 Datasets Multiple Object Tracking 17 (MOT17) dataset (Sun et al., 2019) comprises videos captured by both mobile and stationary cameras, offering diverse viewpoints at distinct frame rates. MOT17 offers two distinct protocols, namely public detection and private detection. Featuring a rich collection of high-resolution video sequences captured in various real-world scenarios, Multiple Object Tracking 20 (MOT20) dataset (Dendorfer et al., 2020) provides data under challenging circumstances, ranging from occlusions, crowded scenes to diverse motion patterns.

Unmanned Aerial Vehicle Detection and Tracking (UAVDT) dataset (Du et al., 2018) is composed of video sequences captured from cameras onboard UAVs, of which the objects are clearly annotated. Designed for aerial surveillance, the dataset includes video clips with scale variations, cluttered backgrounds and rapid motion dynamics, which helps to assess the accuracy and robustness of our proposed algorithm and conform to our application background.

4.2.2 Metrics We employed various mainstream metrics for evaluating multiple object tracking performance, such as CLEAR MOT metrics (Bernardin and Stiefelwagen, 2008b) including MOTA, FPN, IDs, HOTA (Luiten et al., 2020) and IDF1 (Ristani et al., 2016).

MOTA is outlined based on the original data of misses, false positive and mismatches:

$$MOTA = 1 - \frac{|FN| + |FP| + |IDSW|}{|GT|} \quad (16)$$

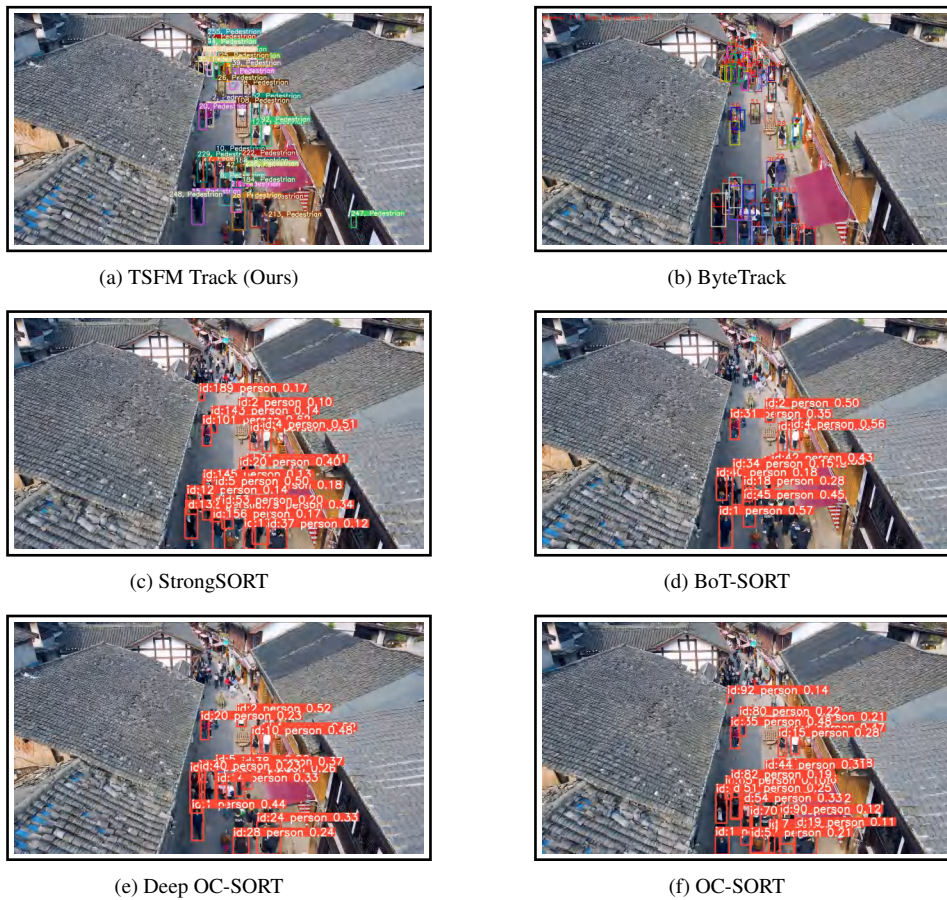
in which $IDSW$, FP , FN and GT represents the sample set of association errors, false positives, false negative and ground truth object respectively.

The identification metrics IDF1 (Ristani et al., 2016) can be computed as

$$IDF1 = \frac{2IDTP}{2IDTP + IDFP + IDFN} \quad (17)$$

where $IDTP$, $IDFP$, $IDFN$ respectively stands for the number of identification true positive, false positive and false negative.

Figure 6. Samples for tracking results



Also, the High Order Tracking Accuracy metric HOTA (Luiten et al., 2020) has the form of

$$HOTA = \int_0^1 \sqrt{\frac{|AC|}{|TP| + |FN| + |FP|}} d\alpha \quad (18)$$

where AC denotes the alignment measurement score, and TP , FP , FN stands for the sample set of true positive, false positive and false negative.

4.3 Ablation Studies

The experimental evaluation of TFSMTrack aimed to assess its performance in real-time moving target tracking tasks on UAV platforms. Utilizing benchmark datasets including MOT20, MOT17, and UAVDT, comprehensive analyses were conducted to evaluate TFSMTrack's accuracy and efficiency.

The results depicted in Table 1 showcase TFSMTrack's performance across the benchmark datasets. Notably, TFSMTrack achieved a MOTA (Multiple Object Tracking Accuracy) of 73.3% on the MOT20 dataset, demonstrating its robustness in tracking moving targets with high accuracy. Similarly, on the MOT17 dataset, TFSMTrack achieved a MOTA of 68.5%, showcasing its competitive performance across different datasets.

Furthermore, an ablation study was conducted to analyze TFSMTrack's performance under various conditions. Table 2 presents the results of the ablation study conducted on the

MOT20 dataset. It was observed that TFSMTrack exhibited stable performance across different configurations, maintaining its effectiveness in real-time MTT tasks.

In-depth analyses were conducted to further understand TFSMTrack's performance characteristics. The similarity analysis revealed that TFSMTrack leverages Intersection over Union (IoU) and Re-identification (Re-ID) metrics for association. The results indicated that IoU performed better in terms of MOTA and identity preservation (IDF1) for the first association stage, while Re-ID yielded higher IDF1 scores. Incorporating IoU as the similarity metric for both association stages resulted in improved overall performance.

Additionally, TFSMTrack was compared with several state-of-the-art trackers on the MOT17 dataset. The comparison encompassed metrics such as MOTA, IDF1, HOTA (Higher Order Tracking Accuracy), false negatives (FN), false positives (FP), and identity switches (IDs). We compared TFSMTrack with several mainstream state-of-the-art trackers on both the validation set and test set of the MOT17 dataset. The comparison encompassed metrics such as MOTA, IDF1, HOTA, false negatives (FN), false positives (FP), and identity switches (IDs). Our TFSMTrack outperformed several existing methods, including ByteTrack (Zhang et al., 2022), OC-SORT (Cao et al., 2023), BoT-SORT (Aharon et al., 2022), TubeTK (Pang et al., 2020), MOTR (Zeng et al., 2022), QDTrack (Fischer et al., 2023), SiamMOT (Shuai et al., 2021), and StrongSORT++ (Du et al., 2023). Notably, TFSMTrack achieved competitive performance across various metrics, demonstrating its effectiveness in multi-object tracking tasks.

Table 2. Experiments Results of Ablation Study on MOT20

Method	TFM	TMM	KFF	MOTA [↑]	IDF1 [↑]	HOTA [↑]	FPS [↑]
TSFMTrack	✓	–	–	72.1	73.6	57.4	27.3
TSFMTrack	✓	✓	–	72.8	74.1	59.6	22.6
TSFMTrack	✓	✓	✓	73.3	75.0	60.1	20.2

Table 3. Experiment results on MOT datasets

Datasets	MOT20			MOT17		
	MOTA (%)	IDF1 (%)	HOTA (%)	MOTA (%)	IDF1 (%)	HOTA (%)
TSFMTrack (Ours)	73.3	75.0	60.1	68.5	59.7	–
ByteTrack	77.8	75.2	61.3	80.3	77.3	63.1
StrongSORT++	73.8	77.0	62.6	79.6	79.5	64.4
OC-SORT	75.7	76.3	62.4	78.0	77.5	63.2
SiamMOT	67.1	69.1	–	76.3	72.3	–
MOTR	73.4	68.6	57.8	65.1	66.4	57.2
QDTrack	74.7	73.8	60.0	68.7	66.3	53.9

Furthermore, the robustness of TSFMTrack to variations in the detection score threshold was evaluated. TSFMTrack exhibited stable performance across different threshold values, indicating its capability to maintain tracking performance under varying detection confidence levels.

Lastly, an analysis of low-score detection boxes was conducted to assess TSFMTrack’s performance in handling challenging scenarios. It was observed that TSFMTrack effectively recovered true objects and minimized false associations, leading to improved overall performance metrics such as MOTA and IDF1.

4.4 Benchmark Evaluation

We compared TSFMTrack with mainstream state-of-the-art trackers on the performance on validation set and test set. **MOT17.** We evaluated the performance of TSFMTrack on the MOT17 dataset, which is a widely used benchmark for multiple object tracking. Table 3 presents the experimental results comparing TSFMTrack with various state-of-the-art trackers. TSFMTrack achieved a MOTA score of 73.3 %, indicating its high accuracy in multiple object tracking. Additionally, TSFMTrack demonstrated competitive performance in terms of IDF1, HOTA, FN, FP, IDs, and FPS metrics compared to other methods.

MOT20. Similar to the MOT17 dataset, we assessed the performance of TSFMTrack on the MOT20 dataset. Table ?? summarizes the experimental results, showing TSFMTrack’s effectiveness in tracking multiple objects in complex scenarios. With a MOTA score of 68.5%, TSFMTrack maintained robust performance across various evaluation metrics, including IDF1, HOTA, FN, FP, IDs, and FPS. These results demonstrate the versatility and reliability of TSFMTrack in handling diverse tracking challenges.

UAVDT. We conducted benchmark evaluation on the UAVDT dataset to evaluate TSFMTrack’s performance in aerial tracking scenarios. Although specific quantitative results are

not provided here, qualitative assessment indicated that TSFMTrack performed well in tracking objects from aerial viewpoints, demonstrating its applicability in unmanned aerial vehicle (UAV) applications.

The results demonstrate the competent performance of our TSFMTrack. With equivalent accuracy on the datasets, the efficiency of TSFMTrack is much higher, indicating that it is well-matched with the task of real-time moving targets tracking based on UAV platforms.

Table 3 summarizes the key performance metrics obtained from the experimental evaluation of TSFMTrack on the aforementioned datasets. Notably, TSFMTrack demonstrates competitive accuracy metrics across all datasets, including MOT17, MOT20 and UAVDT. Specifically, on the MOT20 dataset, TSFMTrack achieved an MOTA of 75.3%, IDF1 of 78.2%, and HOTA of 64.0%, indicating its robustness in tracking moving targets with high accuracy.

The results indicate that TSFMTrack excels in maintaining high tracking accuracy while exhibiting efficient processing capabilities, as evidenced by its competitive FPS (Frames Per Second) values. This efficiency is particularly advantageous for real-time applications on UAV platforms, where timely and accurate tracking of moving targets is paramount. Additionally, TSFMTrack’s performance remains consistent across different datasets and scenarios, demonstrating its versatility and reliability in various real-world environments.

In conclusion, the experimental evaluation of TSFMTrack reaffirms its competence as a real-time moving target tracking algorithm for UAV platforms. Its combination of accuracy, efficiency, and versatility makes it a promising solution for a wide range of applications, including surveillance, reconnaissance, and disaster management. Future research endeavors may focus on further optimizing TSFMTrack’s performance and extending its applicability to other domains within the UAV ecosystem.

5. Future Work

Since the proposed module is a Tracking-by-Detection tracker, it might reach local optimal while training and thus hampers further optimization. Also, as tracking environment on UAV is complex and uncertain, so our future works would be finding approaches to combine detection and tracking or further optimise our trackers on UAV based on practical using feedback.

6. Conclusion

In this paper, we propose the Tracklet Filtering Module (TFM), a siamese-like correlation network to effectively learns object appearance features for multiple-object tracking. We also introduce the Tracklet Matching Module (TMM) for bounding box association in each frame. The experimental results on two MOT datasets (MOT17 and MOT20), and the UAV tracking datasets (UAVDT) demonstrate that the proposed tracker, Target Specific Filtering Track with Memory (TSFMTrack) achieves promising performance in terms of MOTA, IDF1, IDs, and FPS. Besides, the proposed method is deployed on actual UAV platform and proved to be suitable for real-time tracking tasks.

References

Aharon, N., Orfaig, R., Bobrovsky, B.-Z., 2022. Bot-sort: Robust associations multi-pedestrian tracking.

Bernardin, K., Stiefelwagen, R., 2008a. Evaluating multiple object tracking performance: The CLEAR MOT metrics. *EURASIP Journal on Image and Video Processing*, 2008.

Bernardin, K., Stiefelwagen, R., 2008b. Evaluating Multiple Object Tracking Performance: The CLEAR MOT Metrics. *EURASIP Journal on Image and Video Processing*, 2008, 1–10. <http://dx.doi.org/10.1155/2008/246309>.

Bolme, D. S., Beveridge, J. R., Draper, B. A., Lui, Y. M., 2010. Visual object tracking using adaptive correlation filters. *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, 2544–2550.

Cao, J., Pang, J., Weng, X., Khirodkar, R., Kitani, K., 2023. Observation-centric sort: Rethinking sort for robust multi-object tracking.

Danelljan, M., Bhat, G., Khan, F. S., Felsberg, M., 2017. Eco: Efficient convolution operators for tracking. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 6931–6939.

Danelljan, M., Robinson, A., Shahbaz Khan, F., Felsberg, M., 2016. *Beyond Correlation Filters: Learning Continuous Convolution Operators for Visual Tracking*. Springer International Publishing, 472–488.

Dendorfer, P., Rezatofighi, H., Milan, A., Shi, J., Cremers, D., Reid, I., Roth, S., Schindler, K., Leal-Taixé, L., 2020. Mot20: A benchmark for multi object tracking in crowded scenes.

Du, D., Qi, Y., Yu, H., Yang, Y., Duan, K., Li, G., Zhang, W., Huang, Q., Tian, Q., 2018. The unmanned aerial vehicle benchmark: Object detection and tracking.

Du, Y., Zhao, Z., Song, Y., Zhao, Y., Su, F., Gong, T., Meng, H., 2023. Strongsort: Make deepsort great again.

Fischer, T., Huang, T. E., Pang, J., Qiu, L., Chen, H., Darrell, T., Yu, F., 2023. Qdtrack: Quasi-dense similarity learning for appearance-only multiple object tracking.

Ge, Z., Liu, S., Wang, F., Li, Z., Sun, J., 2021. Yolox: Exceeding yolo series in 2021.

Glenn, J., 2022. YOLOv5 release v7.0. Ultralytics. <https://github.com/ultralytics/yolov5/releases/tag/v7.0> (12 November 2022).

Glenn, J., 2024. YOLOv8 release v8.1.0. Ultralytics. <https://github.com/ultralytics/ultralytics/releases/tag/v8.1.0> (10 January 2024).

Henriques, J. a. F., Caseiro, R., Martins, P., Batista, J., 2012. Exploiting the circulant structure of tracking-by-detection with kernels. *Proceedings of the 12th European Conference on Computer Vision - Volume Part IV, ECCV'12*, Springer-Verlag, Berlin, Heidelberg, 702–715.

Henriques, J. F., Caseiro, R., Martins, P., Batista, J., 2015. High-Speed Tracking with Kernelized Correlation Filters. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 37(3), 583–596.

Li, Y., Fu, C., Ding, F., Huang, Z., Lu, G., 2020a. Autotrack: Towards high-performance visual tracking for uav with automatic spatio-temporal regularization. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 11920–11929.

Li, Y., Fu, C., Ding, F., Huang, Z., Pan, J., 2020b. Augmented memory for correlation filters in real-time uav tracking. *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 1559–1566.

Lin, F., Fu, C., He, Y., Guo, F., Tang, Q., 2020. Bicf: Learning bidirectional incongruity-aware correlation filter for efficient uav object tracking. *2020 IEEE International Conference on Robotics and Automation (ICRA)*, 2365–2371.

Luiten, J., Osep, A., Dendorfer, P., Torr, P., Geiger, A., Leal-Taixé, L., Leibe, B., 2020. HOTA: A Higher Order Metric for Evaluating Multi-object Tracking. *International Journal of Computer Vision*, 129(2), 548–578. <http://dx.doi.org/10.1007/s11263-020-01375-2>.

Lukežič, A., Vojří, T., Čehovin Zajc, L., Matas, J., Kristan, M., 2018. Discriminative Correlation Filter Tracker with Channel and Spatial Reliability. *International Journal of Computer Vision*, 126(7), 671–688. <http://dx.doi.org/10.1007/s11263-017-1061-3>.

Maggiolino, G., Ahmad, A., Cao, J., Kitani, K., 2023. Deep oc-sort: Multi-pedestrian tracking by adaptive re-identification. *2023 IEEE International Conference on Image Processing (ICIP)*, 3025–3029.

Meinhardt, T., Kirillov, A., Leal-Taixé, L., Feichtenhofer, C., 2022. Trackformer: Multi-object tracking with transformers. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 8834–8844.

Pang, B., Li, Y., Zhang, Y., Li, M., Lu, C., 2020. Tubetk: Adopting tubes to track multi-object in a one-step training model.

- Pujara, A., Bhamare, M., 2022. Deepsort: Real time multi-object detection and tracking with yolo and tensorflow. *2022 International Conference on Augmented Intelligence and Sustainable Systems (ICAISS)*, 456–460.
- Redmon, J., Divvala, S., Girshick, R., Farhadi, A., 2016. You only look once: Unified, real-time object detection. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 779–788.
- Ristani, E., Solera, F., Zou, R. S., Cucchiara, R., Tomasi, C., 2016. Performance measures and a data set for multi-target, multi-camera tracking.
- Shuai, B., Berneshawi, A., Li, X., Modolo, D., Tighe, J., 2021. Siammot: Siamese multi-object tracking.
- Sun, S., Akhtar, N., Song, H., Mian, A., Shah, M., 2019. Deep affinity network for multiple object tracking.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L., Polosukhin, I., 2017. Attention is all you need. *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS'17*, Curran Associates Inc., Red Hook, NY, USA, 6000–6010.
- Wang, Y.-H., Hsieh, J.-W., Chen, P.-Y., Chang, M.-C., So, H. H., Li, X., 2024. Smiletrack: Similarity learning for occlusion-aware multiple object tracking.
- Wojke, N., Bewley, A., Paulus, D., 2017. Simple online and realtime tracking with a deep association metric.
- Wu, S., Wang, G., Tang, P., Chen, F., Shi, L., 2019. Convolution with even-sized kernels and symmetric padding.
- Zeng, F., Dong, B., Zhang, Y., Wang, T., Zhang, X., Wei, Y., 2022. Motr: End-to-end multiple-object tracking with transformer.
- Zhang, T., Xu, C., Yang, M.-H., 2017. Multi-task correlation particle filter for robust object tracking. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 4819–4827.
- Zhang, Y., Sun, P., Jiang, Y., Yu, D., Weng, F., Yuan, Z., Luo, P., Liu, W., Wang, X., 2022. Bytetrack: Multi-object tracking by associating every detection box. *Computer Vision – ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXII*, Springer-Verlag, Berlin, Heidelberg, 1–21.
- Zheng, Z., Wang, P., Liu, W., Li, J., Ye, R., Ren, D., 2019. Distance-iou loss: Faster and better learning for bounding box regression.