# Sampling-Resilient Multi-Object Tracking

## Technical Report

### Abstract

Multi-Object Tracking (MOT) is a cornerstone operator to support video surveillance applications. To process large-scale live video streams in real time, there is critical demand for an efficient and accurate tracking mechanism. In this paper, we study a new scenario called down-sampled multi-object tracking, which performs object detection and association only upon a subset of video frames, to improve efficiency with tolerable accuracy decline.

The problem is challenging for state-of-the-art MOT methods, which exhibit significant performance degradation when the frame reduction ratio is high. In this paper, we devise a sampling-resilient tracker with more accurate motion estimation and robust data association. We adapt Kalman filter with sparse observations, by augmenting its state with more informative representation and devising a dynamic parameter matrix update mechanism to reduce the divergence between observations and internal state estimation. To associate the detected bounding boxes robustly, we propose a comprehensive similarity metric that systematically integrates multiple spatial matching signals.

Experiments on three benchmark datasets show that our proposed tracker is the most suitable for real-time tracking and achieves the best trade-off between efficiency and accuracy. Compared with ByteTrack, we can further reduce the processing time by $2\times$ in MOT17 and $3\times$ in DanceTrack and reach the same level of tracking accuracy.

## 1   Introduction

Multi-object tracking (MOT) aims at detecting and tracking all the moving objects from video clips or live streams, while maintaining a unique identifier for each object. Massive research efforts have been devoted into this domain with fruitful progress. The proposed trackers have witnessed great success in numerous applications, such as smart video surveillance [34], traffic monitoring [27], customer behavior analysis [19] and sports analytics [16].

In this paper, we focus on real-time MOT and study a new scenario called down-sampled multi-object tracking. It performs object detection and association only upon a sampled set of video frames to further reduce inference time, with a tolerable accuracy decline. The task has the potential to achieve an ideal trade-off between tracking efficiency and accuracy. It has not been well explored in the past literature and is useful for scalable video analytic applications. In practice, the computation overhead of deep neural networks remains prohibitive for many surveillance cameras or low-end
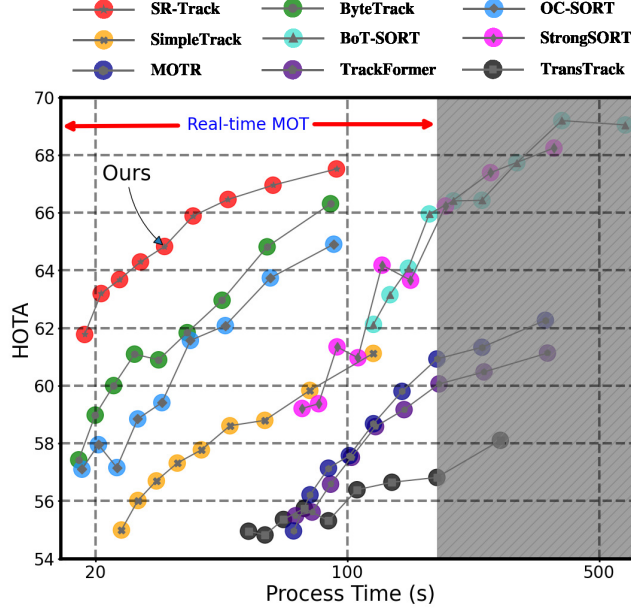
Figure 1: We plot the trade-off between tracking efficiency and accuracy for state-of-the-art MOT methods, by adjusting sampling rate in the test set of MOT17. With fewer video frames sampled, the total process time is reduced, but at the cost of lower tracking accuracy. Since the total length of the test data is 177.2 seconds, the non-shaded area implies real-time tracking which requires online inference to catch up with the speed of video streaming. The results on HOTA show that our proposed SR-Track establishes clear superiority to support real-time tracking — it is the most efficient at the same accuracy level. Note that the process time in the x-axis is in log scale.

IoT devices. Their captured video streams have to be collected and transmitted to an edge server or cloud server for real-time processing [15]. Given a smart surveillance system with thousands of cameras, there is a critical demand for an efficient and accurate tracking mechanism so that the video streams can be processed with affordable computation devices.

The new problem of down-sampled MOT is challenging because the motion pattern becomes more difficult to capture and the position estimation in the next frame becomes less accurate. In addition, the data association strategy such as IoU that works well in dense frames fails in the scenario of sparse frames. Therefore, directly applying state-of-the-art MOT methods on the down-sampled frames would result in significant performance degradation. As shown in Figure 1, we report the trade-off between efficiency and accuracy in terms of MOTA, IDF1 and HOTA, by adjusting different sampling rates. When the video frames are sampled at a very low rate, it's indeed that the processing time can be significantly reduced. However, the tracking accuracy also declines to an impractical level. Detailed performance analysis of these trackers will be presented in Section 4.5. These findings lead to the conclusion that existing MOT

solutions are not sampling resilient.

To devise a more sampling-resilient MOT model, we follow the tracking-by-detection paradigm as used in ByteTrack [39] and propose SR-Track with more accurate motion estimation and robust data association mechanisms. First, to estimate the next position of an object more accurately with sparse observations, we augment Kalman filter (KF) with more informative state representation to capture the evolution of the object motion state. By measuring the divergence between observation and internal state estimation, we can dynamically estimate the noise scale of the motion process under sparse observations. Furthermore, we propose an aligned state update mechanism that enhances the adaptability of the motion model to new observations. Second, to robustly associate detected bounding boxes under widening intervals, we reveal interesting findings from an experimental analysis on existing metrics. We propose a comprehensive similarity metric that integrates multiple spatial matching clues, including overlap, center point distance and aspect ratio of the bounding boxes.

We also notice alternative efforts dedicated to improving KF for more accurate motion estimation. For example, StrongSORT [8] and GiaoTracker [9] adaptively modulate the observation noise scale according to the quality of object detection. In OC-SORT [5], the authors claim that the assumption of consistent velocity direction does not hold due to the non-linear motion of objects and state noise. They also argue that non-linear variants such as Extended KF [11] and Unscented KF [28] are very difficult to implement for online tracking due to the lack of prior knowledge on the complex motion pattern. Therefore, they design an update strategy under occlusion to reduce noise and add the velocity consistency (momentum) term into the cost matrix for better matching between tracklets and observations. *However, these approaches focus on reducing noise with frequent observations. We will experimentally show that they fail to work well in down-sampled MOT with sparse observations.*

Experiments are conducted on three benchmark datasets, among which DanceTrack is the most challenging due to frequent crossover and diverse body gestures. The results show that our proposed tracker outperforms most trackers in terms of both efficiency and accuracy. For the real-time trackers that can achieve similar FPS, our SR-Track exhibits clearly higher accuracy accuracy. Compared with ByteTrack, the state-of-the-art real-time tracker, we can further reduce the processing time by $2\times$ in MOT17 and $3\times$ in DanceTrack, with the same level of tracking accuracy.

## 2 Related Work

In this section, we review state-of-the-art multi-object trackers and divide them into two categories, namely *tracking-by-detection* and *joint-detection-and-tracking*, according to whether its object detection network is a separate module or requires joint training.

### 2.1 Tracking-by-Detection Methods

SORT [4], DeepSORT [32], OC-SORT [5], StrongSORT [8], BoT-SORT [1] and Byte-Track [39] are representative tracking-by-detection methods. They treat MOT as a pipeline of object detection and association, and optimize each module separately.

Firstly, an existing object detector is adopted to locate objects in each video frame. Early trackers (e.g., SORT and DeepSORT) use Faster RCNN [22] as the default detector, which is replaced by YOLOX [10] in recent trackers. Secondly, an object association mechanism is designed to connect these detected objects into tracklets. Coherence in motion pattern and similarity in visual appearance are two important factors in object association. As to motion pattern, almost all the tracking-by-detection methods adopt Kalman filter for future position estimation. A detected object is assigned to an existing tracklet if its spatial matching distance (e.g., IoU distance) between the two bounding boxes is small. As to visual similarity, DeepSORT [32], StrongSORT and BoT-SORT integrate appearance features into the tracker, which requires additional computation cost to derive visual embedding. The spatial matching score and appearance similarity are combined as the final association metric.

Among these trackers, ByteTrack [39] achieves the best trade-off between efficiency and accuracy. It discards visual similarity and only relies on spatial matching to save computation cost. As a compensation, it introduces a robust association strategy to take into account the detected objects with low confidence.

## 2.2 Joint-Detection-and-Tracking Methods

JDE [30] is a pioneering work that allows object detection and appearance embedding to be learned in a single network. Compared with DeepSORT, its low-level visual features can be shared by the detector and embedding model to avoid re-computation cost. However, the shared network in JDE is biased towards the detector task and unfair to the ReID task. To resolve the competition issue, CSTrack [14] devises a cross-correlation network to learn task-dependent representations. RelationTrack [36] presents global context disentangling (GCD) to decouple th learned features in the two tasks. Fair-MOT [40] adopts another way by implementing two homogeneous branches for the detection and ReID tasks, rather than performing them in a two-stage cascaded style. SimpleTrack [12] is designed to mitigate the issue of object occlusion and presents a new association matrix that combines embedding cosine distance and Giou distance of objects. Note that these works still rely on an online data association strategy based on Kalman filter and appearance similarity to connect the detected boxes.

To push forward the idea of joint training, the following trackers attempt to further incorporate the estimation of inter-frame object motion in the training framework. In other words, Kalman filter is discarded. CenterTrack [42] and TransCenter [35] predict the object offset between adjacent frames to facilitate object tracking. The models are trained to minimize the regression loss of the object offset between adjacent frames. TransCenter [35] proposes a Transformer-based architecture, together with dense but non-overlapping representations for detection, to globally and robustly infer the offset of objects' centers. For GSDT [29] and FUEFT [24], motion and appearance features are fed into a graph neural network (GNN) to predict the association matrix of tracklets and detected bounding boxes. TransTrack [26] utilizes the attention mechanism to model the detection and tracking, and outputs the predicted bounding box of tracked objects. Recently, TrackFormer [18] adopts the concept of track queries and employs the attention mechanism to track the objects in an autoregressive fashion. In the current stage, these trackers are computation expensive to achieve high accuracy and not

4

suitable for real-time tracking.

## 2.3 Low frame rate object tracking

There are some works [13, 31, 38, 6] studying object tracking for cameras inherently with low frame rate. Their algorithm pipelines are focused on robust tracking, and often incur higher computation overhead. For example, [38] adopts a complex matching mechanism based on particle swarm optimization. Since the objective of this paper is to move towards real-time tracking by purposely reducing the number of frames, the above solutions cannot be applied to down-sampled MOT.

# 3 Methodology of SR-Track

Before we present our SR-Track, we first briefly review Kalman filter (KF), which has been widely adopted in object tracking to estimate object location in the subsequent frame. It works as an efficient recursive filter with the stages of estimation and update. KF requires small computational power and provides satisfactory estimation, rendering it well-suited for real-time analysis.

Let $\hat{x}_{k-1}$ be the object state at the $(k-1)^{th}$ frame and $F$ be the state transition matrix. In the estimation step, the state at the $k^{th}$ frame $\hat{x}'_k$ and state estimated covariance matrix $P'_k$ are predicted via the following equations, where $Q_k$ is the process noise covariance matrix. $Q_k$ consists of the errors caused in the motion process and is an important parameter matrix in KF. For example, if the velocity of the detected object changes rapidly, KF can determine an appropriate $Q_k$ matrix to reflect the unreliability of the system at this moment.

$$
\begin{array}{rcl}
\hat{x}'_k & = & F\hat{x}_{k-1} \\
P'_k & = & FP_{k-1}F^\top + Q_k
\end{array}
\tag{1}
$$

$$
\tag{2}
$$

In the update step, KF blends the new observation with the old information from prior state with the Kalman gain matrix $K_k$. The estimation of $K_k$ is shown in Eq. (3), where $H$ is the observation matrix and $R_k$ is the observation noise covariance matrix. In Eq. (4), the actual observation $z_k$ is obtained to generate a posterior state estimate of $\hat{x}'_k$. The residual $z_k - H\hat{x}'_k$ reflects the divergence between the predicted state and the observed state. Finally, in Eq. (5), the estimation state covariance matrix $P'_k$ is also updated according to the Kalman gain $K_k$.

$$
\begin{array}{rcl}
K_k & = & P'_k H^\top \left( H P'_k H^\top + R_k \right)^{-1} \\
\hat{x}_k & = & \hat{x}'_k + K_k \left( z_k - H\hat{x}'_k \right) \\
P_k & = & \left( I - K_k H \right) P'_k
\end{array}
\tag{3}
$$

$$
\tag{4}
$$

$$
\tag{5}
$$

In the scenario of down-sampled MOT, the observations become sparse and each object appears in fewer number of video frames. Consequently, the uncertainty is amplified and it becomes more challenging to capture the model pattern. The traditional KF as well as its improved variants in StrongSORT and OC-SORT fail to address these

unique challenges. Therefore, we are motivated to devise a new variant KF for sparse observations.

## 3.1 Sparse-Observation Kalman Filter

The pipeline of our proposed Sparse-Observation Kalman Filter (SOKF) is illustrated in Figure 2, with the following three key components.
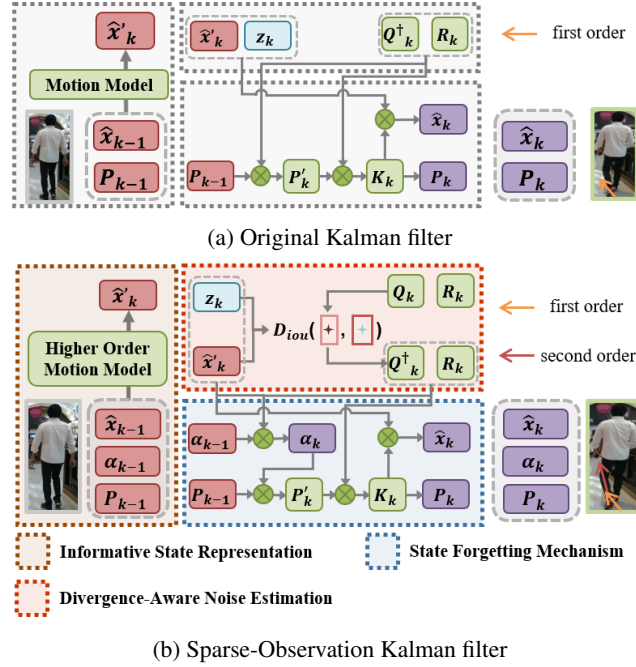


(a) Original Kalman filter



(b) Sparse-Observation Kalman filter

Figure 2: Pipelines of KF and SOKF.

**Informative State Representation.** The object state is normally represented by its position, shape and motion information. As a common practice, the state vector adopted by most MOT trackers is in the form of $(x_c, y_c, s, r, \dot{x}_c, \dot{y}_c, \dot{s})$, where $(x_c, y_c)$ is the center point; $s$ and $r$ capture the area and aspect ratio, respectively; $(\dot{x}_c, \dot{y}_c, \dot{s})$ are the first-order change rate of $(x_c, y_c, s)$ and capture the velocity information. It is worth noting that in some trackers [32, 1], the shape information $s$ and $r$ in the state information are replaced with height $h$ and width $w$. When observations are sparse, even minor changes in an object's position can result in significant and complex alterations over time. The original state representation becomes inadequate and we devise a more informative object state as follows:

$$(x_d, y_d, w, h, \dot{x}_d, \dot{y}_d, \dot{w}, \dot{h}, \ddot{x}_d, \ddot{y}_d, \ddot{w}, \ddot{h})$$

where $(x_d, y_d)$ is the center point of *bottom edge* in the bounding box; $(w, h)$ refer to width and height; $(\dot{x}_d, \dot{y}_d, \dot{w}, \dot{h})$ and $(\ddot{x}_d, \ddot{y}_d, \ddot{w}, \ddot{h})$ store the first-order and second-order change rate of $(x_d, y_d, w, h)$, respectively.

6

Here, we explain the design principles of the state vector. First, we use the center point $(x_d, y_d)$ of the *bottom edge*, rather than the center of the whole rectangle. This is a dataset-oriented trick to leverage the fact that the target persons are often moving on the flat ground. Take DanceTrack as an example, the dancers are moving on the stage and the center point of the *bottom edge* can provide more reliable position information, which is less sensitive to the variance of height. When the status of a dancer changes from stand to squat, the height of the bounding box shrinks, but our definition of center point remains the same.

Secondly, existing solutions utilized a first-order change rate for motion tracking, which we think is not sufficient to capture the rapid position update and complex motion pattern in down-sampled MOT. Our strategy is to incorporate higher order change rates to acquire more informative motion representation. Since more dimensions of state will trigger higher computational and tuning difficulties, we only augment the state with 4-dimensional second-order change rate, in order to strike a balance between performance and cost.

**Divergence-Aware Noise Estimation.** In Kalman filter, the process noise covariance matrix $Q_k$ and the observation noise covariance matrix $R_k$ play very important role because they are used to regulate the impact of prediction and observation on the system state estimation, with the goal of maximizing a posterior estimation. If the noise parameters are incorrect, the accuracy of KF may reduce dramatically.

When the video frames are down-sampled, $R_k$ is not affected because it is mainly determined by the inherent performance of the object detector, which is YOLOX in our implementation. However, $Q_k$ captures motion process noise that represent the estimation noise scale and increases due to the amplified uncertainty in the scenario of sparse observations. Thus, we focus on the update mechanism of $Q_k$ and design an divergence-aware mechanism (DAM) to estimate the process noise covariance matrix $Q_k$. We observe that abrupt change or significant variation on the motion pattern becomes more frequent in the scenario of down-sampled MOT. The traditional Kalman filter cannot immediately update the parameters to keep in line with the dramatic motion change. Our goal is explicitly enlarge the values in matrix $Q_k$ when the uncertainty increases. In other words, when the divergence between estimation and observation becomes large, we use it as a signal to reactively adjust the matrix $Q_k$. Formally, we utilize the IoU of between the estimated bounding box ($H\hat{x}_k$ in Eq. (4)) and the observed bounding box $z_k$ to measure the divergence between estimation and observation.

$$Q_k^{\dagger} = (1 + D_{iou}(H\hat{x}_k, z_k))Q_k \tag{6}$$

$D_{iou}$ is the IoU distance between two bounding boxes and defined as $D_{iou} = 1 - IoU$, where $IoU$ is the overlap area between two bounding boxes divided by their union area.

**State Forgetting Mechanism.** KF captures the noise distribution of the system via updating the state estimation covariance matrix $P_{k-1}$, which determines the influence of the historical state. However, in down-sampled MOT with sparse observations, the historical state estimation becomes less reliable and we should pay more attention to recent observations. Inspired by fading KF [33], we incorporate a state forgetting mechanism into the updating process of $P_{k-1}$. As shown in Eq. (7), we add a fading factor $\alpha_k$ in the transition process of the estimated noise covariance matrix $P'_k$.

7

$$\boldsymbol{P}'_k = \boldsymbol{\alpha}_k \boldsymbol{F}_k \boldsymbol{P}_{k-1} \boldsymbol{F}_k^\top + \boldsymbol{Q}_k^\dagger \tag{7}$$

Initially, we set $\boldsymbol{\alpha} = 1.0$ and update it according to the divergence of estimation and observation $\|\boldsymbol{z}_k - \boldsymbol{H}\hat{\boldsymbol{x}}'_k\|$ in the subsequent iterations. When the divergence between estimation and observation at $(k-1)^{th}$ step is large, the fading factor $\alpha_k$ increases, so as to reduce the influence of the old observation.

## 3.2 Robust Data Association (RDA)

Data association is also a key component in the tracking-by-detection paradigm. The mainstream metrics estimate the spatial matching score according to either IoU (Intersection of Union) [4, 39, 5, 1] or center point distance between two bounding boxes [32, 40, 8]. On the other hand, there also exist certain factors that have been adopted in the loss of object detection (e.g., aspect ratio in CIoU loss [41]), but they are not leveraged by object tracking.

We perform an experimental analysis on these metrics when applied to object tracking across down-sampled video frames. We denote the sample reduction ratio by $RR$, which implies that $\frac{1}{RR}$ frames are sampled. When $RR = 1$, all the frames are preserved. We vary $RR$ from 1 to 9 and for each setting, we randomly collect $10,000$ bounding box association cases that can be successfully solved by at least one of the following metrics, including the overlap, center point distance and aspect ratio of the bounding boxes. We denote them by IoU, DIST, and SCALE, respectively.

Interesting findings can be derived from the results reported in Table 1. The set $S_{metric}$ includes the cases that can be correctly matched by the associated metric. $P_{SCALE}$ represents the cases that can only be solved by SCALE, i.e., IoU and DIST fail in these cases. When there is no down-sampling (with $RR = 1$), it's indeed that IoU or distance-based metric demonstrate very good performance as they are able to correctly identify around $99\%$ of the matching cases. The metric SCALE is inferior to the two metrics as it generates many false negatives. Its complementary effect to IoU and DIST can be negligible because only $0.31\%$ of cases can be uniquely solved by SCALE. This may explain why SCALE is not adopted by the state-of-the-art MOT methods. However, when $RR$ increases, IoU and DIST become less reliable as the sizes of $|S_{IoU}|$ and $|S_{DIST}|$ reduce. It is interesting to find that the factor of SCALE plays a more important role and its size of $P_{SCALE}$ increases with $RR$. This finding motivates us to devise a comprehensive association metric that incorporate all these three factors.

|  | $RR = 1$ | $RR = 3$ | $RR = 5$ | $RR = 7$ | $RR = 9$ |
|---|---|---|---|---|---|
| $|S_{IoU}|$ | 9899 | 9504 | 9169 | 8812 | 8565 |
| $|S_{DIST}|$ | 9891 | 9579 | 9320 | 8999 | 8797 |
| $|S_{SCALE}|$ | 7886 | 6928 | 6444 | 6191 | 6010 |
| $|P_{SCALE}|$ | 31 | 118 | 174 | 234 | 275 |

Table 1: distance metrics analysis on the MOT17 dataset.

Let $\boldsymbol{D}_{iou}$ denote the overlap distance between two bounding boxes and $\boldsymbol{D}_{dist}$ denote the normalized distance between two center points of the bounding boxes.

$$\boldsymbol{D}_{dist} = \frac{\|(x_d, y_d)_1, (x_d, y_d)_2\|^2}{c^2} \tag{8}$$

where $c$ is is the diagonal length of the smallest enclosing box covering the bounding boxes. For the factor of aspect ratio, we define $\boldsymbol{D}_{scale}$ as

$$\boldsymbol{D}_{scale} = \frac{4}{\pi^2} \left( \arctan \frac{w_1}{h_1} - \arctan \frac{w_2}{h_2} \right)^2 \tag{9}$$

where $w_i$ and $h_i$ are the width and height of the two bounding boxes, respectively. To integrate these three distances, we define $\boldsymbol{D}_{rda}$ as follows. The idea is to first use IoU and DIST if these two metrics can provide confident matching results. This is because as revealed in Table 1, these two factors normally provide better results than SCALE. We use $\frac{\boldsymbol{D}_{dist}+\boldsymbol{D}_{iou}}{2}$ to reversely approximate for the confidence. This is a reasonable estimation because it implies that the estimated bounding box is close to the region of the detected object. If this value is smaller than a threshold $\sigma$, the tracking confidence is high and we directly set $\boldsymbol{D}_{rda} = \frac{\boldsymbol{D}_{dist}+\boldsymbol{D}_{iou}}{2}$. Otherwise, we need to incorporate $\boldsymbol{D}_{scale}$ as a complementary factor and set $\boldsymbol{D}_{rda}$ as a linear combination of the three factors.

$$\boldsymbol{D}_{rda} = \begin{cases} \frac{\boldsymbol{D}_{dist}+\boldsymbol{D}_{iou}}{2} & \frac{\boldsymbol{D}_{dist}+\boldsymbol{D}_{iou}}{2} < \sigma \\ \frac{\boldsymbol{D}_{dist}+\boldsymbol{D}_{iou}+2\boldsymbol{D}_{scale}}{4} & \text{otherwise} \end{cases} \tag{10}$$

## 4 Experiment

In this section, we compare SR-Track with state-of-the-art MOT methods on three benchmark datasets. At the end of the section, we also compare with OTIF [2] in terms of answering selection and aggregation queries in a video database.

### 4.1 Experimental Setup

**Benchmark Datasets.** We use three benchmark datasets for performance evaluation, including MOT17 [20], MOT20 [7] and DanceTrack [25]. MOT17 contains 14 videos (7 for training and 7 for testing) of pedestrians in both indoor and outdoor scenes. MOT20 contains 8 videos (4 for training, 4 for testing) in crowded environments such as train stations, town squares and a sports stadium. DanceTrack is a recent dataset proposed to emphasize the importance of motion analysis. The frequent crossover, uniform appearance and diverse body gestures of dancers bring particular challenges. It provides 100 videos and the split ratio for training, validation and test dataset is $40 : 24 : 35$.

Note that the testing videos of these datasets are not annotated and their performance evaluation requires the submission of the generated tracklets to the official website. Since the focus of this paper is down-sampled MOT, we directly use the annotated

| | RR = 3 | | | RR = 5 | | | RR = 7 | | | RR = 9 | | | FPS |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | HOTA | MOTA | IDF1 | HOTA | MOTA | IDF1 | HOTA | MOTA | IDF1 | HOTA | MOTA | IDF1 | |
| | | | | | | Dataset MOT17 | | | | | | | |
| SimpleTrack [12] | 59.8 | 69.3 | 75.6 | 58.6 | 66.3 | 73.4 | 57.3 | 63.9 | 72.1 | 56.0 | 61.4 | 70.2 | 22.5 |
| OC-SORT [5] | 63.7 | 68.6 | 74.5 | 61.5 | 63.7 | 71.0 | 58.8 | 59.3 | 67.6 | 57.9 | 57.8 | 66.4 | 29.0 |
| ByteTrack [39] | 64.8 | 73.8 | 76.3 | 61.8 | 70.3 | 72.1 | 61.0 | 67.9 | 71.0 | 58.9 | 65.2 | 68.8 | **29.6** |
| SR-Track (ours) | **67.0** | **75.7** | **78.8** | **65.9** | **72.7** | **76.3** | **64.2** | **69.6** | **73.6** | **63.2** | **67.5** | **72.4** | 28.9 |
| | | | | | | Dataset MOT20 | | | | | | | |
| SimpleTrack [12] | 52.2 | 65.3 | 68.0 | 51.4 | 63.7 | 67.7 | 50.0 | 61.6 | 65.8 | 47.8 | 58.8 | 62.5 | 7.0 |
| OC-SORT [5] | 56.3 | 69.6 | 72.1 | 54.6 | 68.0 | 69.8 | 53.2 | 66.4 | 68.1 | 50.7 | 63.3 | 64.4 | **18.7** |
| ByteTrack [39] | 56.0 | 71.3 | 71.1 | 55.5 | 70.1 | 70.8 | 54.2 | 68.4 | 69.7 | 50.7 | 65.7 | 65.2 | 17.5 |
| SR-Track (ours) | **57.3** | **71.8** | **73.6** | **57.6** | **71.3** | **74.1** | **58.1** | **70.6** | **74.8** | **55.8** | **68.9** | **71.2** | 17.0 |
| | | | | | | Dataset DanceTrack | | | | | | | |
| OC-SORT [5] | 49.2 | 84.6 | 48.7 | 40.8 | 77.8 | 41.6 | 36.1 | 69.7 | 37.9 | 33.4 | 63.1 | 34.8 | 29.0 |
| ByteTrack [39] | 40.7 | 82.3 | 46.9 | 35.5 | 74.7 | 39.8 | 32.8 | 68.5 | 37.0 | 32.0 | 63.0 | 35.8 | **29.6** |
| SR-Track (ours) | **53.8** | **88.1** | **53.6** | **46.6** | **84.3** | **45.4** | **42.6** | **79.5** | **40.7** | **38.9** | **74.7** | **37.2** | 29.0 |

Table 2: Comparison with real-time trackers on three benchmark datasets with varying frame reduction ratio $RR$.

videos for performance evaluation. For MOT17 and MOT20, we split the videos into two parts with equal length and use them for training and testing, respectively. For DanceTrack, we use the training set for model training and report the performance on its validation set.

**Performance Metrics.** To evaluate the overall tracking accuracy, we adopt MOTA [3], IDF1 [23] and HOTA [17]. Generally speaking, MOTA score is biased towards measuring the quality of object detection and IDF1 emphasizes the effect of accurate association. HOTA is a recent metric proposed to explicitly balance the effect of detection, association and localization.

As to efficiency, we adopt $FPS$ as a straightforward metric. It refers to the number of video frames that can be processed per second. In addition, we propose a new metric called **Time@HOTA**. The underlying motivation is that we can adjust $RR$ to generate a trade-off curve between processing time and HOTA, as shown in Figure 1. It can be expected that with larger processing time (i.e., smaller $RR$), we can obtain higher HOTA. Time@HOTA measures the processing time required to reach a specified HOTA. For example, $Time@62 = 19$ for our SR-Track at dataset MOT17 implies that it takes 19 seconds for SR-Track to process the testing videos in MOT17 with an accuracy level of $HOTA = 62$.

**Comparison Methods.** We compare SR-Track with representative and open-sourced trackers in both categories of tracking-by-detection and joint-detection-and-tracking. Among these competitors, we consider ByteTrack [39], OC-SORT [5] and Simple-Track [12] as **real-time trackers** because they can achieve as high FPS as our SR-Track. The remaining approaches, including TransTrack [26], TrackFormer [18], MOTR [37], StrongSORT [8] and BoT-SORT [1], are called **expensive trackers** as they exchange processing time for higher tracking accuracy.

## 4.2 Implementation Details

Our SR-Track follows the paradigm of tracking-by-detection. For object detector, we directly adopt the trained YOLOX provided by previous trackers (e.g., ByteTrack for

MOT17 and MOT20, OC-SORT for DanceTrack). The future positions of the detected objects are estimated by the Sparse-Observation Kalman filter (SOKF) proposed in this paper. The objects between different frames are associated via the robust data association (RDA). Therefore, there is no training required for our SR-Track.

As to our proposed Kalman filter, we set $\boldsymbol{\alpha}_0 = 1.0$ for the adaptive fading factor and adopt the previous work [1] to estimate process noise covariance matrix $\boldsymbol{Q}_k$ and observation noise covariance matrix $\boldsymbol{R}_k$ for each frame. All the experiments are conducted using PyTorch and ran on a desktop with 10th Intel(R) Core(TM) i9-10980XE CPU @ 3.00GHz and NVIDIA GeForce RTX 3090Ti GPU.

### 4.3 Comparison with Real-time Trackers

In the first experiment, we compare our SR-Track with the real-time trackers under different reduction ratios (with $RR$ set to 3, 5, 7 and 9, respectively). As shown in Table 2, these trackers demonstrate similar inference speed. OC-SORT, ByteTrack and SR-Track use YOLOX as the object detector and their association ignores visual similarity. Although SimpleTrack adopts appearance similarity for person ReID, it trains the object detector and visual embedding with a single network to avoid recomputation cost. Its FPS is slightly lower than other real-time trackers.

Among these real-time trackers, SR-Track achieves the highest MOTA, IDF1 and HOTA across all the datasets, owing to its Kalman filter designed for the observation-sparse scenario. The performance gap between ByteTrack and our SR-Track is widened when $RR$ increases. In MOT20, the HOTA of SR-Track is higher than ByteTrack by $2.3\%$ when $RR = 3$, which is enlarged to $10\%$ when $RR = 9$.

DanceTrack is a challenging dataset with complex motion pattern and frequent crossover of dancers, which are difficult for existing trackers to perform correct association. Thus, their derived IDF1 and HOTA in DanceTrack are generally lower than those in MOT17 and MOT20. OC-SORT outperforms ByteTrack in this dataset because it is specially designed for DanceTrack and occlusion with excessive nonlinear motion. Nevertheless, the performance of OC-SORT degrades to be close to ByteTrack when RR increases, implying that its strategy is not robust to the observation-sparse scenario. These two models are both significantly inferior to our SR-Track. When $RR = 9$, we boost the HOTA from 33.4 in OC-SORT to 38.9, with $16.5\%$ improvement.

### 4.4 Comparison with Expensive Trackers

In Table 3, we compare SR-Track with the expensive trackers under $RR = 5$ and $RR = 9$. For TransTrack, TrackFormer, MOTR, StrongSORT, their performance is clearly inferior to our SR-Track in terms of both tracking efficiency and accuracy. BoT-SORT is the only method whose accuracy can be slightly better than our SR-Track in MOT17. However, its tracking speed is very slow and the FPS is 6 times lower than SR-Track. Furthermore, similar to previous findings in Table 2, the advantage of SR-Track becomes more obvious when $RR$ increases. In MOT20 with $RR = 9$, our SR-Track can even achieve higher accuracy than BoT-SORT, with 7 times faster

11

|  | $RR = 5$ | | | $RR = 9$ | | | FPS |
| --- | --- | --- | --- | --- | --- | --- | --- |
|  | HOTA | MOTA | IDF1 | HOTA | MOTA | IDF1 | |
| MOT17 | | | | | | | |
| TransTrack [26] | 56.8 | 66.1 | 66.6 | 54.8 | 61.2 | 62.2 | 10.0 |
| TrackFormer [18] | 59.1 | 66.2 | 68.2 | 55.6 | 60.6 | 64.2 | 7.4 |
| MOTR [37] | 59.8 | 65.5 | 68.8 | 56.2 | 61.0 | 65.4 | 7.5 |
| StrongSORT [8] | 63.6 | 61.9 | 70.9 | 59.3 | 53.2 | 62.9 | 7.1 |
| BoT-SORT [1] | **66.4** | **74.3** | **77.9** | **63.1** | **70.1** | **73.1** | 4.5 |
| SR-Track (ours) | 65.9 | 72.7 | 76.3 | **63.1** | 67.5 | 72.4 | **28.9** |
| MOT20 | | | | | | | |
| TransTrack [26] | 31.6 | 47.3 | 44.6 | 30.5 | 44.9 | 42.4 | 7.2 |
| TrackFormer [18] | 47.4 | 70.6 | 56.8 | 43.3 | 65.3 | 51.3 | 4.1 |
| MOTR [37] | 42.8 | 50.6 | 56.1 | 38.0 | 43.0 | 49.7 | 4.2 |
| StrongSORT [8] | 56.5 | 67.4 | 72.8 | 50.7 | 61.2 | 66.6 | 1.4 |
| BoT-SORT [1] | **57.7** | 71.1 | 73.9 | 54.0 | 67.2 | 69.3 | 2.4 |
| SR-Track (ours) | 57.6 | **71.3** | **74.1** | **55.8** | **68.9** | **71.2** | **17.0** |

Table 3: Comparison with expensive trackers on the datasets of MOT17 and MOT20, under different settings of $RR$.

tracking speed. The results on DanceTrack are not available because we lack sufficient hardware resources to re-train these models.

## 4.5 Trade-off Between Efficiency and Accuracy

At the beginning of the paper, we have reported the trade-off between processing time and tracking accuracy for MOT17. The results on DanceTrack in terms of IDF1 and HOTA are presented in Figure 3. ByteTrack is fast and accurate because it does not rely on visual similarity and improves the association mechanism by taking into account detected objects with low confidence. OC-SORT outperforms ByteTrack in dataset DanceTrack because OC-SORT is better at capturing complex motion patterns. Strong-SORT and BoT-SORT utilize visual similarity by extracting appearance features and achieve high accuracy, but at the cost of significantly higher computation overhead. SimpleTrack, the most recent work proposed in the paradigm of joint training of object detection and embedding, achieves modest performance. However, since the joint training is difficult to coordinate, it does not demonstrate superiority in terms of effectiveness. Finally, TransTrack jointly trains object detection, object ReID and motion estimation in the same framework. Its performance is not satisfactory due to insufficient annotated samples for training and its online inference cost is also expensive.

We also study the performance of these trackers in terms of our proposed metric Time@HOTA. As reported in Tables 4 and 5, ByteTrack achieves the best Time@HOTA among the comparison trackers in MOT17. Our SR-Track can further reduce its processing time by half with a given HOTA requirement. For example, it takes SR-Track 19s to generate tracking results in MOT17 with HOTA=62, whereas ByteTrack requires 37.3s. In dataset DanceTrack, the advantage of SR-Track in terms of Time@HOTA is enlarged to $3\times$. If there are a large number of dancing video streams that require real-time tracking, it can estimated that our SR-Track only consumes roughly $1/3$ of GPU devices to achieve the same performance as ByteTrack.
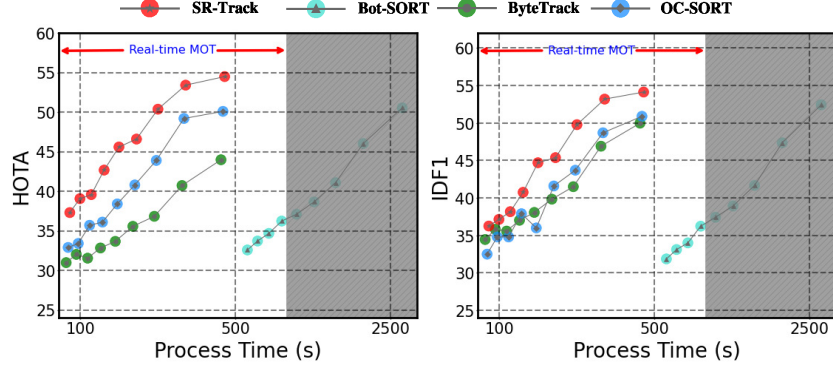
Figure 3: Trade-off analysis in DanceTrack.

|  | Time@66 | Time@65 | Time@64 | Time@63 | Time@62 |
|---|---|---|---|---|---|
| TrackFormer [18] | 718.3 | 718.3 | 718.3 | 718.3 | 718.3 |
| MOTR [37] | 708.8 | 708.8 | 708.8 | 708.8 | 317.1 |
| TransTrack [26] | 531.6 | 531.6 | 531.6 | 531.6 | 531.6 |
| SimpleTrack [12] | 235.9 | 235.9 | 235.9 | 235.9 | 235.9 |
| StrongSORT [8] | 193.5 | 167.1 | 166.2 | 120.4 | 112.3 |
| BoT-SORT [1] | 195.1 | 157.3 | 156.6 | 136.8 | 117.2 |
| OC-SORT [5] | 183.3 | 94.6 | 92.6 | 58.2 | 45.1 |
| ByteTrack [39] | 91.7 | 62.3 | 61.5 | 48.3 | 37.3 |
| SR-Track | **45.1** | **31.9** | **31.6** | **22.2** | **19.0** |

Table 4: Time@HOTA in MOT17 (in seconds).

|  | Time@53 | Time@50 | Time@47 | Time@44 | Time@41 |
|---|---|---|---|---|---|
| BoT-SORT [1] | 3989.8 | 2890.2 | 2036.1 | 1662.6 | 1405.1 |
| OC-SORT [5] | 879.6 | 577.5 | 345.5 | 249.2 | 196.7 |
| MOTR [37] | 621.8 | 384.9 | 295.1 | 248.4 | 214.4 |
| ByteTrack [39] | 861.8 | 861.8 | 861.8 | 432.1 | 295.6 |
| SR-Track | **270.6** | **229.5** | **182.8** | **136.8** | **118.7** |

Table 5: Time@HOTA in DanceTrack (in seconds).

## 4.6  Ablation Study

We evaluate the advantage brought by the Sparse-Observation Kalman filter (SOKF) and robust data association (RDA) in Table 6. ByteTrack can be viewed as a variant without these two components. It is not surprising to find that when RDA is removed, the performance on the matching-related metrics, such as AssA and IDF1, drop significantly. In contrast, SOKF is more important for the remaining metrics. We also conduct a break-down analysis on the components of SOKF and examine the effect of

13

our proposed Informative State Representation (ISR), Divergence-Aware Mechanism (DAM) and State Forgetting Mechanism (SFM). We can see that they all contribute to the improvement of tracking accuracy and play important roles on different metrics.

| | HOTA↑ | DetA↑ | AssA↑ | MOTA↑ | IDF1↑ |
|---|---|---|---|---|---|
| SR-Track | **63.2** | **61.2** | **65.8** | **67.5** | **72.4** |
| SR-Track w/o SOKF | 60.1 | 57.8 | 63.7 | 65.2 | 71.5 |
| SR-Track w/o RDA | 61.9 | 61.0 | 63.6 | 67.5 | 69.6 |
| ByteTrack [39] | 58.9 | 57.6 | 61.6 | 65.2 | 68.8 |
| Break-down analysis on SOKF | | | | | |
| SOKF w/o ISR | 62.2 | 59.3 | 66.3 | 67.5 | 71.8 |
| SOKF w/o DAM | 61.9 | 60.0 | 64.7 | 66.6 | 71.5 |
| SOKF w/o SFM | 62.2 | 60.8 | 64.3 | 66.2 | 70.9 |

Table 6: Ablation study of SR-Track on the MOT17 dataset with $RR = 9$.

## 4.7 Experiments Without Down-Sampling

We are also curious to examine the performance of our SR-Track on the original dataset without down-sampling. Table 7 shows the results returned by the leadboard of Dance-Track. SR-Track is the best performer and improves the metrics of HOTA, IDF1 and AssA by 7.3%, 9.4% and 12.9%, respectively.

| | HOTA↑ | DetA↑ | AssA↑ | MOTA↑ | IDF1↑ |
|---|---|---|---|---|---|
| FairMOT [40] | 39.7 | 66.7 | 23.8 | 82.2 | 40.8 |
| QDTrack [21] | 45.7 | 72.1 | 29.2 | 83.0 | 44.8 |
| TransTrack [26] | 45.5 | 75.9 | 27.5 | 88.4 | 45.2 |
| MOTR [37] | 48.4 | 71.8 | 32.7 | 79.2 | 46.1 |
| ByteTrack [39] | 47.3 | 71.6 | 31.4 | 89.5 | 52.5 |
| OC-SORT [5] | 55.1 | 80.3 | 38.0 | 89.4 | 54.2 |
| SR-Track | **59.1** | **81.5** | **42.9** | **92.4** | **59.3** |

Table 7: Performance on DanceTrack test dataset.

## 4.8 Case Analysis

Finally, we perform a case analysis by comparing SR-Track and ByteTrack on MOT17 with $RR = 9$. As shown in Figure 4, we highlight the incorrect association generated by ByteTrack. From frame 4 to frame 5, its Kalman filter makes wrong estimation of the next bounding box, whereas our Sparse-Observation Kalman filter delivers accurate estimation. From frame 16 to frame 17, ByteTrack incurs ID switching caused by occlusion, but our SR-Track, with more robust association metric, is able to resolve the issue.
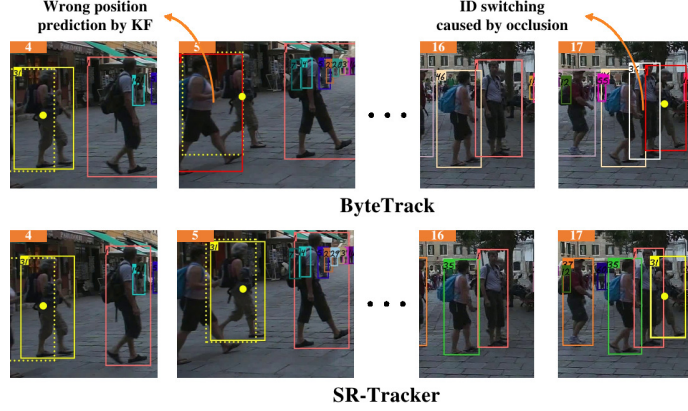
Figure 4: A case study for ByteTrack and SR-Track in MOT17.

## 4.9 Supporting Database Queries

OTIF [2] is a state-of-the-art video database to answer selection and aggregation queries. Its advantages essentially come from its multi-object tracking (MOT) model to extract the tracklets of all the moving objects and index them to support selection and aggregation queries. For instance, to retrieve video frames containing ambulance and firetruck, we can simply perform an intersection between the inverted lists of these two labels. As another example, to estimate the traffic flow with a time period, we can leverage the spatial-temporal index to identify the relevant video frames and aggregate the number of distinct objects within these frames.

In this paper, we devise a sampling-resilient MOT model that works well with very low sampling rate and integrate it in our video database called DoveDB. It is then an interesting question to compare the performance of our DoveDB and OTIF. As to experimental setup, we use the original implementation of OTIF using Go language. As to the dataset, we use Jackson Town used in the paper of OTIF. It is collected from a traffic surveillance camera and the video length is around 1 hour. In the following, we evaluate the performance of OTIF with DoveDB in terms of both multi-object tracking and video query processing.

In Table 8, we vary the sampling ratio and report MOTA and IDF1, which are two popular performance metrics to measure tracking accuracy. When we use lower sampling ratio, the inference time can be significantly reduced but the accuracy of tracking also drops. Under the same ratio, DoveDB is much superior over OTIF in terms of both efficiency and accuracy. The reason for the speedup is that our DoveDB relies on Kalman filter to connect the bounding boxes in different frames. Instead, OTIF trains an RNN network to perform data association. As to the tracking accuracy, DoveDB is designed with multiple techniques to track objects with sparse observations. It's not surprising that its tracking accuracy is significantly better than OTIF at low sampling rate.

As mentioned, we can leverage the results from MOT models to answer selection

| | Sampling Ratio=1/8 | | | Sampling Ratio=1/16 | | | Sampling Ratio=1/32 | | |
|---|---|---|---|---|---|---|---|---|---|
| | MOTA | IDF1 | Time | MOTA | IDF1 | Time | MOTA | IDF1 | Time |
| OTIF | 55.5 | 69.6 | 1113.5s | 51.5 | 62.3 | 556.7s | 42.3 | 53.1 | 278.3s |
| DoveDB | **74.6** | **74.7** | **752.**1s | **70.8** | **67.4** | **353.7**s | **65.2** | **59.9** | **181.2**s |

Table 8: Tracking performance between OTIF and DoveDB.

and aggregation queries. In Figure 5, we compare DoveDB with OTIF in terms of the trade-off between efficiency and query accuracy, by varying sampling ratios. The selection query $Q_1$ retrieves all the video frames containing at least one vehicle. In the evaluation, as long as there exists a class label "bus", "car", "truck" or "suv", we consider it a true positive. The aggregation query $Q_2$ estimates the total number of vehicles that appear in the video clip. We use F1-score to measure the accuracy of selection query and Mean Average Error to measure estimation error for the aggregation query. With the same video ingestion time, our DoveDB achieves significantly more accurate retrieval performance than OTIF. This owes to the superiority of DoveDB to handle video frames with low sampling rate. It can achieve the same query accuracy level with OTIF by ingesting much fewer number of video frames.
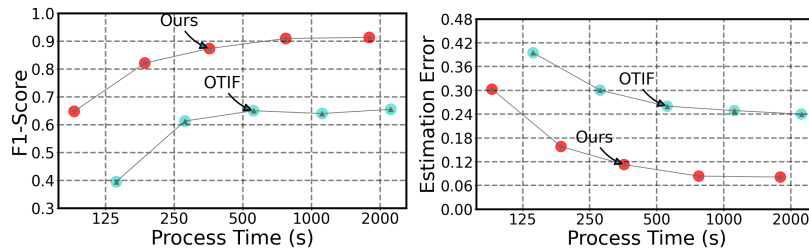


Figure 5: Comparison between DoveDB and OTIF for database queries.

# 5 Conclusion

In this paper, we studied a new scenario of multi-object tracking on down-sampled video frames and devised a sampling-resilient tracker. In particular, we proposed an Sparse-Observation Kalman filter for accurate motion estimation and a comprehensive data association metric for robust inter-frame matching. Experiments on three benchmark datasets show that our proposed SR-Track establishes new SOTA performance for real-time tracking.

# References

[1] Nir Aharon, Roy Orfaig, and Ben-Zion Bobrovsky. Bot-sort: Robust associations multi-pedestrian tracking. *CoRR*, abs/2206.14651, 2022. 3, 6, 8, 10, 11, 12, 13

[2] Favyen Bastani and Samuel Madden. OTIF: efficient tracker pre-processing over large video datasets. In Zachary Ives, Angela Bonifati, and Amr El Abbadi, editors, *SIGMOD '22: International Conference on Management of Data, Philadelphia, PA, USA, June 12 - 17, 2022*, pages 2091–2104. ACM, 2022. 9, 15

[3] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: The CLEAR MOT metrics. *EURASIP J. Image Video Process.*, 2008, 2008. 10

[4] Alex Bewley, ZongYuan Ge, Lionel Ott, Fabio Tozeto Ramos, and Ben Upcroft. Simple online and realtime tracking. In *2016 IEEE International Conference on Image Processing, ICIP 2016, Phoenix, AZ, USA, September 25-28, 2016*, pages 3464–3468. IEEE, 2016. 3, 8

[5] Jinkun Cao, Xinshuo Weng, Rawal Khirodkar, Jiangmiao Pang, and Kris Kitani. Observation-centric SORT: rethinking SORT for robust multi-object tracking. *CoRR*, abs/2203.14360, 2022. 3, 8, 10, 13, 14

[6] Meng-Che Chuang, Jenq-Neng Hwang, Kresimir Williams, and Richard Towler. Tracking live fish from low-contrast and low-frame-rate stereo videos. *IEEE Trans. Circuits Syst. Video Technol.*, 25(1):167–179, 2015. 5

[7] Patrick Dendorfer, Hamid Rezatofighi, Anton Milan, Javen Shi, Daniel Cremers, Ian Reid, Stefan Roth, Konrad Schindler, and Laura Leal-Taixé. Mot20: A benchmark for multi object tracking in crowded scenes. *arXiv: Computer Vision and Pattern Recognition*, 2020. 9

[8] Yunhao Du, Yang Song, Bo Yang, and Yanyun Zhao. Strongsort: Make deepsort great again. *CoRR*, abs/2202.13514, 2022. 3, 8, 10, 12, 13

[9] Yunhao Du, Junfeng Wan, Yanyun Zhao, Binyu Zhang, Zhihang Tong, and Junhao Dong. Giaotracker: A comprehensive framework for MCMOT with global information and optimizing strategies in visdrone 2021. In *IEEE/CVF International Conference on Computer Vision Workshops, ICCVW 2021, Montreal, BC, Canada, October 11-17, 2021*, pages 2809–2819. IEEE, 2021. 3

[10] Zheng Ge, Songtao Liu, Feng Wang, Zeming Li, and Jian Sun. Yolox: Exceeding yolo series in 2021. *arXiv: Computer Vision and Pattern Recognition*, 2021. 4

[11] Simon J. Julier and Jeffrey K. Uhlmann. New extension of the Kalman filter to nonlinear systems. In Ivan Kadar, editor, *Signal Processing, Sensor Fusion, and Target Recognition VI*, volume 3068, pages 182 – 193. International Society for Optics and Photonics, SPIE, 1997. 3

[12] Jiaxin Li, Yan Ding, Hua-Liang Wei, Yutong Zhang, and Wenxiang Lin. Simpletrack: Rethinking and improving the JDE approach for multi-object tracking. *Sensors*, 22(15):5863, 2022. 4, 10, 13

[13] Yuan Li, Haizhou Ai, Takayoshi Yamashita, Shihong Lao, and Masato Kawade. Tracking in low frame rate video: A cascade particle filter with discriminative observers of different life spans. *IEEE Trans. Pattern Anal. Mach. Intell.*, 30(10):1728–1740, 2008. 5

[14] Chao Liang, Zhipeng Zhang, Yi Lu, Xue Zhou, Bing Li, Xiyong Ye, and Jianxiao Zou. Rethinking the competition between detection and reid in multi-object tracking. *CoRR*, abs/2010.12138, 2020. 4

[15] Shengzhong Liu, Tianshi Wang, Jinyang Li, Dachun Sun, Mani B. Srivastava, and Tarek F. Abdelzaher. Adamask: Enabling machine-centric video streaming with adaptive frame masking for DNN inference offloading. In João Magalhães, Alberto Del Bimbo, Shin'ichi Satoh, Nicu Sebe, Xavier Alameda-Pineda, Qin Jin, Vincent Oria, and Laura Toni, editors, *ACM MM*, pages 3035–3044. ACM, 2022. 2

[16] Wei-Lwun Lu, Jo-Anne Ting, James J. Little, and Kevin P. Murphy. Learning to track and identify players from broadcast sports videos. *IEEE Trans. Pattern Anal. Mach. Intell.*, 35(7):1704–1716, 2013. 1

[17] Jonathon Luiten, Aljosa Osep, Patrick Dendorfer, Philip H. S. Torr, Andreas Geiger, Laura Leal-Taixé, and Bastian Leibe. HOTA: A higher order metric for evaluating multi-object tracking. *Int. J. Comput. Vis.*, 129(2):548–578, 2021. 10

[18] Tim Meinhardt, Alexander Kirillov, Laura Leal-Taixé, and Christoph Feichtenhofer. Trackformer: Multi-object tracking with transformers. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 8834–8844. IEEE, 2022. 4, 10, 12, 13

[19] Djamal Merad, Kheir-Eddine Aziz, Rabah Iguernaissi, Bernard Fertil, and Pierre Drap. Tracking multiple persons under partial and global occlusions: Application to customers' behavior analysis. *Pattern Recognit. Lett.*, 81:11–20, 2016. 1

[20] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler. Mot16: A benchmark for multi-object tracking. *arXiv: Computer Vision and Pattern Recognition*, 2016. 9

[21] Jiangmiao Pang, Linlu Qiu, Xia Li, Haofeng Chen, Qi Li, Trevor Darrell, and Fisher Yu. Quasi-dense similarity learning for multiple object tracking. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25, 2021*, pages 164–173. Computer Vision Foundation / IEEE, 2021. 14

[22] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. *Advances in neural information processing systems*, 28, 2015. 4

[23] Ergys Ristani, Francesco Solera, Roger S. Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In Gang Hua and Hervé Jégou, editors, *Computer Vision - ECCV 2016 Workshops - Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part II*, volume 9914 of *Lecture Notes in Computer Science*, pages 17–35, 2016. 10

[24] Chaobing Shan, Chunbo Wei, Bing Deng, Jianqiang Huang, Xian-Sheng Hua, Xiaoliang Cheng, and Kewei Liang. Tracklets predicting based adaptive graph tracking. *arXiv preprint arXiv:2010.09015*, 2020. 4

[25] Peize Sun, Jinkun Cao, Yi Jiang, Zehuan Yuan, Song Bai, Kris Kitani, and Ping Luo. Dancetrack: Multi-object tracking in uniform appearance and diverse motion. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition, CVPR 2022, New Orleans, LA, USA, June 18-24, 2022*, pages 20961–20970. IEEE, 2022. 9

[26] Peize Sun, Yi Jiang, Rufeng Zhang, Enze Xie, Jinkun Cao, Xinting Hu, Tao Kong, Zehuan Yuan, Changhu Wang, and Ping Luo. Transtrack: Multiple-object tracking with transformer. *CoRR*, abs/2012.15460, 2020. 4, 10, 12, 13, 14

[27] Wei Tian, Martin Lauer, and Long Chen. Online multi-object tracking using joint domain information in traffic scenarios. *IEEE Trans. Intell. Transp. Syst.*, 21(1):374–384, 2020. 1

[28] E.A. Wan and R. Van Der Merwe. The unscented kalman filter for nonlinear estimation. In *Proceedings of the IEEE 2000 Adaptive Systems for Signal Processing, Communications, and Control Symposium (Cat. No.00EX373)*, pages 153–158, 2000. 3

[29] Yongxin Wang, Kris Kitani, and Xinshuo Weng. Joint object detection and multi-object tracking with graph neural networks. In *IEEE International Conference on Robotics and Automation, ICRA 2021, Xi'an, China, May 30 - June 5, 2021*, pages 13708–13715. IEEE, 2021. 4

[30] Zhongdao Wang, Liang Zheng, Yixuan Liu, Yali Li, and Shengjin Wang. Towards real-time multi-object tracking. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XI*, volume 12356 of *Lecture Notes in Computer Science*, pages 107–122. Springer, 2020. 4

[31] Peter Wei, Haocong Shi, Jiaying Yang, Jingyi Qian, Yinan Ji, and Xiaofan Jiang. City-scale vehicle tracking and traffic flow estimation using low frame-rate traffic cameras. In Robert Harle, Katayoun Farrahi, and Nicholas D. Lane, editors, *Proceedings of the 2019 ACM International Joint Conference on Pervasive and Ubiquitous Computing and Proceedings of the 2019 ACM International Symposium on Wearable Computers, UbiComp/ISWC 2019 Adjunct, London, UK, September 9-13, 2019*, pages 602–610. ACM, 2019. 5

[32] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *2017 IEEE International Conference on Image Processing, ICIP 2017, Beijing, China, September 17-20, 2017*, pages 3645–3649. IEEE, 2017. 3, 4, 6, 8

[33] Qijun Xia, Ming Rao, Yiqun Ying, and Xuemin Shen. Adaptive fading kalman filter with an application. *Automatica*, 30(8):1333–1338, 1994. 7

[34] Ronghua Xu, Seyed Yahya Nikouei, Yu Chen, Aleksey Polunchenko, Sejun Song, Chengbin Deng, and Timothy R. Faughnan. Real-time human objects tracking for smart surveillance at the edge. In *2018 IEEE International Conference on Communications, ICC 2018, Kansas City, MO, USA, May 20-24, 2018*, pages 1–6. IEEE, 2018. 1

[35] Yihong Xu, Yutong Ban, Guillaume Delorme, Chuang Gan, Daniela Rus, and Xavier Alameda-Pineda. Transcenter: Transformers with dense queries for multiple-object tracking. *CoRR*, abs/2103.15145, 2021. 4

[36] En Yu, Zhuoling Li, Shoudong Han, and Hongwei Wang. Relationtrack: Relation-aware multiple object tracking with decoupled representation. *CoRR*, abs/2105.04322, 2021. 4

[37] Fangao Zeng, Bin Dong, Tiancai Wang, Cheng Chen, Xiangyu Zhang, and Yichen Wei. MOTR: end-to-end multiple-object tracking with transformer. *CoRR*, abs/2105.03247, 2021. 10, 12, 13, 14

[38] Xiaoqin Zhang, Weiming Hu, Nianhua Xie, Hujun Bao, and Stephen J. Maybank. A robust tracking system for low frame rate video. *Int. J. Comput. Vis.*, 115(3):279–304, 2015. 5

[39] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. Bytetrack: Multi-object tracking by associating every detection box. *CoRR*, abs/2110.06864, 2021. 3, 4, 8, 10, 13, 14

[40] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. Fairmot: On the fairness of detection and re-identification in multiple object tracking. *Int. J. Comput. Vis.*, 129(11):3069–3087, 2021. 4, 8, 14

[41] Zhaohui Zheng, Ping Wang, Wei Liu, Jinze Li, Rongguang Ye, and Dongwei Ren. Distance-iou loss: Faster and better learning for bounding box regression. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 12993–13000. AAAI Press, 2020. 8

[42] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Tracking objects as points. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part IV*, volume 12349 of *Lecture Notes in Computer Science*, pages 474–490. Springer, 2020. 4