

# QDTrack: Quasi-Dense Similarity Learning for Appearance-Only Multiple Object Tracking

Tobias Fischer\*, Thomas E. Huang\*, Jiangmiao Pang\*, Linlu Qiu, Haofeng Chen, Trevor Darrell, Fisher Yu

**Abstract**—Similarity learning has been recognized as a crucial step for object tracking. However, existing multiple object tracking methods only use sparse ground truth matching as the training objective, while ignoring the majority of the informative regions in images. In this paper, we present Quasi-Dense Similarity Learning, which densely samples hundreds of object regions on a pair of images for contrastive learning. We combine this similarity learning with multiple existing object detectors to build Quasi-Dense Tracking (QDTrack), which does not require displacement regression or motion priors. We find that the resulting distinctive feature space admits a simple nearest neighbor search at inference time for object association. In addition, we show that our similarity learning scheme is not limited to video data, but can learn effective instance similarity even from static input, enabling a competitive tracking performance without training on videos or using tracking supervision. We conduct extensive experiments on a wide variety of popular MOT benchmarks. We find that, despite its simplicity, QDTrack rivals the performance of state-of-the-art tracking methods on all benchmarks and sets a new state-of-the-art on the large-scale BDD100K MOT benchmark, while introducing negligible computational overhead to the detector.

**Index Terms**—Multiple Object Tracking, Quasi-Dense Similarity Learning.



## 1 INTRODUCTION

MULTIPLE Object Tracking (MOT) is a fundamental and challenging problem in computer vision, widely used in safety monitoring, autonomous driving, video analytics, and other applications. Contemporary MOT methods [1], [2], [3], [4], [5] mainly follow the tracking-by-detection paradigm [6]. That is, they detect objects on each frame and then associate them according to the estimated similarity between each instance. Recent works [1], [2], [4], [7] show that if the detected objects are accurate, the spatial proximity between objects in consecutive frames, measured by Intersection over Union (IoU) or center distance, is a strong prior to associate the objects. However, this location prior is often violated in more complex scenarios with non-linear object motion, varying video frame rate, or complex camera motion, since the movement of objects on the image plane depends highly on these factors. To remedy this problem, some methods introduce motion estimation [8], [9] or displacement regression [4], [10], [11] to ensure accurate distance estimation. Object appearance similarity usually takes a secondary role [3], [12] to strengthen object association or re-identify vanished objects, because extracted appearance features cannot effectively distinguish different objects. Thus, the search region is constrained to local neighborhoods to avoid distractions.

On the contrary, humans can easily associate identical objects only through appearance. We conjecture this is because the image and object information is not fully utilized for learning object similarity. Previous methods regard instance similarity learning as a post-hoc stage after object detection or only use sparse ground truth

bounding boxes as training samples [3]. These processes ignore the majority of the informative regions on the images. We hypothesize that, because objects in an image are rarely identical to each other, a nearest neighbor search in a learned feature space should associate and distinguish instances without bells and whistles. In addition, we observe that besides the ground truth and detected bounding boxes, which sparsely distribute on the images, many possible object regions can provide valuable training supervision.

In this paper, we propose quasi-dense similarity learning, which densely matches hundreds of informative regions on a pair of images for contrastive learning. The quasi-dense samples cover a wide range of informative regions on the images, providing both more positive examples and hard negatives. Because one sample has more than one positive counterpart on the reference image, we extend the InfoNCE loss [13] commonly used in contrastive learning [14], [15], [16] to multiple positives which makes quasi-dense learning feasible. Each sample is thus trained to discriminate an instance from all possible object regions on the reference image simultaneously. This provides stronger supervision than using only a handful ground truth labels and enhances the instance similarity learning. To extract feature embeddings for each region, we use a lightweight embedding extractor that works with most existing object detectors.

Besides similarity, the inference pipeline, which measures the instance similarity and maintains a track history, also plays an important role in the tracking performance, since it needs to consider false positives, missed detections, newly appeared objects, and terminated tracks. To better deal with these cases, we introduce the *bi-directional softmax* similarity metric that enforces bi-directional consistency. In particular, objects that do not have matching targets in the other frame will lack a bi-directional matching and thus have low similarity scores to all other objects. Furthermore, we include unmatched objects in the previous frame, which we call backdrops, for matching to better filter false positives that could otherwise act as distractors in following frames. We compose object detectors, quasi-dense similarity learning, and

- T. Fischer, T. E. Huang, F. Yu are with the Department of Information Technology and Electrical Engineering, ETH Zurich.  
E-mail: tobias.fischer@vision.ee.ethz.ch
  - J. Pang is with the Shanghai AI Laboratory.
  - L. Qiu is with the Department of EECS, Massachusetts Institute of Technology.
  - H. Chen is with the Computer Science Department, Stanford University.
  - T. Darrell is with the Department of EECS, UC Berkeley.
- \* Equal contribution.

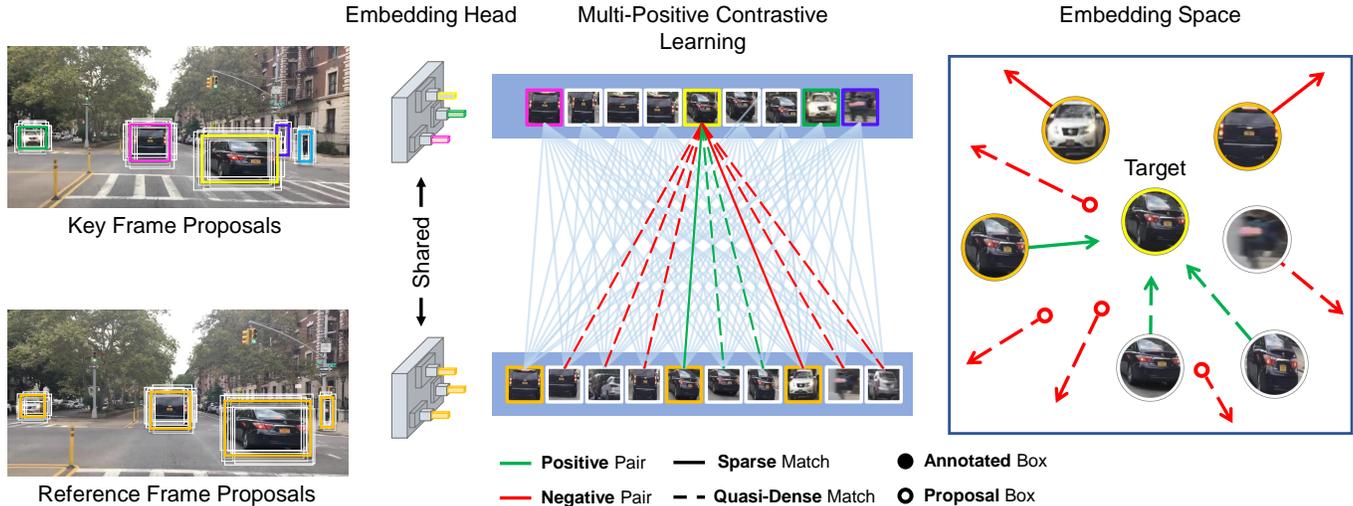


Fig. 1: **Training pipeline.** After we extract feature embeddings for all quasi-dense samples on a pair of key and reference images, we apply dense matching between them and optimize the learned representation with multiple positive contrastive learning. The resulting embedding space effectively discriminates different instances.

our inference pipeline to build *Quasi-Dense Tracking* (QDTrack) models. Since the publication of our initial work [17], QDTrack has been widely adopted for other tracking problems, such as segmentation tracking [18], long-tailed multi-object tracking [19], and 3D object tracking [20].

In addition to the findings of our initial work [17], we show that quasi-dense instance similarity learning is not limited to video data, but can learn effective instance representations from static images alone. In particular, we show that we can effectively perform tracking even when learning instance similarity without any annotations for association and/or video input. Moreover, in this journal extension we conduct extensive experiments on a wide variety of tracking benchmarks, namely MOT [21], DanceTrack [22], BDD100K [23], Waymo [24], and TAO [25]. In addition, we show the flexibility of our method by combining it with different base models and object detectors. Despite its simplicity, QDTrack rivals the performance of state-of-the-art methods without bells and whistles, and sets a new state-of-the-art on the large-scale BDD100K tracking benchmark. QDTrack allows for joint, end-to-end training of detection and instance similarity, thereby simplifying the training and inference pipelines of MOT frameworks. In addition, our embedding extractor only adds negligible overhead to the inference time of the detector. We hope the simplicity and strengths of QDTrack motivates further research on similarity learning for multiple object tracking.

## 2 RELATED WORK

In MOT, the current leading paradigm is tracking-by-detection [6]. Tracking-by-detection methods detect objects in each individual frame, and subsequently associate the detections over time. They differ in their data association mechanisms and cues that are used in the association process. A variety of approaches have been developed to solve the data association problem, *e.g.*, network flow formulations [26], quadratic pseudo boolean optimization [27], conditional random fields [28], or multi-hypothesis tracking [29]. Many works have focused on finding the best cues to exploit for data association, such as 2D motion [2], [7], [10], [30], [31],

[32], 3D motion [20], [33], [34], [35], [36], or visual appearance similarity [1], [3], [37], [38], [39], [40], [41], [42]. In this work, we focus on learning visual appearance similarity and follow the tracking-by-detection paradigm.

**Location and motion in MOT** Spatial proximity has been proven effective to associate objects in consecutive frames [2], [7]. Some methods use 2D motion, such as predictions of a Kalman Filter [2], [31], [43], [44], optical flow [30], and displacement regression [10], [32], to estimate similarity for object association. However, these methods are brittle when it comes to varying video frame rate and complex camera motion, since the 2D motion of the objects depends highly on these factors. Thus, other methods instead rely on 3D motion cues to associate objects over time, since in 3D camera and object motion can be decomposed. A common paradigm [33], [34] is to track objects with 3D bounding boxes and motion estimates derived from *e.g.* scene flow. In contrast, [35], [36], [45] explored to track and reconstruct objects in 3D by estimating the object’s rigid-body transformation between two frames. Although these methods show promising results, many [3], [43] still rely on an extra appearance similarity model as a complementary component to re-identify vanished objects, complicating the entire framework. Our approach is orthogonal to this line of work, as we rely solely on appearance-based instance similarity and a simple nearest-neighbor search to associate objects.

**Appearance similarity in MOT** In order to strengthen tracking and re-identify vanished objects, some methods exploit appearance similarity extracted from an independent model [1], [3], [37], [38], [39], [40], [41], [42] or add an extra embedding head to the detector for end-to-end training [5], [12], [46], [47]. However, they still learn appearance similarity following the practice in image similarity learning, and measure the instance similarity by cosine distance. An appearance similarity model is trained either as a  $n$ -class classification problem [3], where  $n$  is equal to the number of identities in the whole training set, or using triplet loss [48]. The classification problem is hard to extend to large-scale datasets, while the triplet loss only compares each training sample with two other samples, leading to sub-optimal

instance similarity learning. As a consequence, these methods still rely heavily on motion models and displacement predictions to track objects, and appearance similarity only takes a secondary role. In contrast, QDTrack learns instance similarity from densely-connected contrastive pairs and associates objects with a simple nearest neighbor search in feature space, which allows for a simpler tracking framework without compromising accuracy.

**Joint detection and tracking** Instead of treating object detection and association as separate modules, Detect & Track [10] is the first work that jointly optimizes object detection and tracking modules. It predicts the displacements of each object in consecutive frames and associates them with the Viterbi algorithm. Tracktor [1] directly adopts a detector as a tracker. CenterTrack [4] and Chained-Tracker [11] predict the object displacements with pair-wise inputs to associate the objects. Other methods focus on learning visual appearance and detection jointly [5], [46], [47], adding an extra embedding head to the detection network. However, these methods do not fully exploit image information for similarity learning. Recent work [49], [50], [51] focuses on leveraging Transformer networks to integrate tracking and detection into a single, query-based architecture. These methods track by propagating queries across timesteps, processing them with a Transformer that outputs the tracking result. In this work, we focus on learning appearance similarity from quasi-dense samples jointly with detection.

**Self-supervised representation learning** The field of self-supervised representation learning has seen significant progress in recent years, fueled by a number of methods relying on contrastive learning [13], [16], [52], [53], [54], [55], [56], [57] that have shown promising performance. The main paradigm of these methods is to learn a representation that is similar for two versions of the same image, where one is distorted with random image augmentations, while enforcing that this representation is dissimilar to other pairs in the current training batch. While this has proven to be very effective, it has not yet drawn much attention when learning the instance similarity in MOT. In this paper, we supervise densely matched quasi-dense samples with multiple positive contrastive learning inspired by [57]. In contrast to image-level contrastive learning, our method allows for multiple positive training, while the methods mentioned above can only handle the case when there is only a single positive target. The promising results of our method show the importance of representation learning for the MOT problem.

**Learning to track from static images** Learning to track objects from static images where no association annotations are available has recently been proposed by multiple methods [4], [47]. CenterTrack [4] proposes to use data augmentation to simulate video input from given a single static image to obtain 2D offsets to learn object motion. FairMOT [47] treats every object in a given detection dataset as a unique class and learns to distinguish between those to learn tracking from static images. In contrast to learning simulated motion or treating every object over a whole dataset as unique, we show that our similarity learning scheme can effectively learn to track objects from static images with comparable accuracy to video input without further modification. We draw inspiration from the success of recent self-supervised representation learning methods and apply our similarity learning scheme between two augmented instances of the same input image.

### 3 METHOD

We propose *quasi-dense similarity learning* to learn a feature embedding space that can associate identical objects and distinguish

different objects for online multiple object tracking. We define *dense matching* to be matching between bounding box candidates at all pixel locations and *sparse matching* to be matching between ground truth box labels as matching candidates. In contrast, *quasi-dense matching* considers potential object candidates specifically at potential object regions. The main ingredients of *Quasi-Dense Tracking* (QDTrack) are object detection, instance similarity learning, and object association.

#### 3.1 Object detection

Our method can be easily coupled with both two-stage and one-stage detectors with end-to-end training. Object detectors contain two components, a feature extractor and a bounding box prediction head. The feature extractor is typically composed of a base model to extract image features and a Feature Pyramid Network (FPN) [58] to obtain multi-scale features. The bounding box prediction head produces dense bounding box candidates, from which we sample quasi-dense samples by filtering with Non-Maximal Suppression (NMS). The resulting samples indicate likely object regions that include multiple overlapped bounding boxes for each object.

#### 3.2 Quasi-dense similarity learning

We use regions that likely contain objects to learn the instance similarity with quasi-dense matching. The full training pipeline is shown in Figure 1. Given a key image  $I_1$  for training, we randomly select a reference image  $I_2$  from its temporal neighborhood. The neighbor distance is constrained by an interval  $k$ . We use the object regions from both images and RoI Align [59] to obtain their corresponding feature maps from the image features. We add an extra lightweight embedding head in parallel with the original bounding box head to extract features for each region. A region is defined as a positive sample to a ground truth object if it has an IoU higher than  $\alpha_1$  or negative if lower than  $\alpha_2$ . A matching of regions on two frames is positive if the two regions are associated with the same ground truth object and negative otherwise.

Assume there are  $V$  samples on the key frame as training samples and  $K$  samples on the reference frame as contrastive targets. For each training sample, we can use the non-parametric softmax [13], [16] with cross-entropy to optimize the feature embeddings,

$$\mathcal{L}_{\text{embed}} = -\log \frac{\exp(\mathbf{v} \cdot \mathbf{k}^+)}{\exp(\mathbf{v} \cdot \mathbf{k}^+) + \sum_{\mathbf{k}^-} \exp(\mathbf{v} \cdot \mathbf{k}^-)}, \quad (1)$$

where  $\mathbf{v}$ ,  $\mathbf{k}^+$ ,  $\mathbf{k}^-$  are feature embeddings of the training sample, its positive target, and negative targets in  $K$ . The overall embedding loss is averaged across all training samples, but we only illustrate one training sample for brevity.

We apply dense matching between object regions on the pairs of images. Specifically, each sample in  $I_1$  is matched to all samples in  $I_2$ , in contrast to only using sparse sample crops (mostly ground truth boxes) to learn instance similarity in previous works [48], [60]. Each training sample in the key frame has more than one positive targets in the reference frame, so Eq. (1) can be extended as

$$\mathcal{L}_{\text{embed}} = -\sum_{\mathbf{k}^+} \log \frac{\exp(\mathbf{v} \cdot \mathbf{k}^+)}{\exp(\mathbf{v} \cdot \mathbf{k}^+) + \sum_{\mathbf{k}^-} \exp(\mathbf{v} \cdot \mathbf{k}^-)}. \quad (2)$$

However, this equation does not treat positive and negative targets fairly. Namely, each negative is considered multiple times

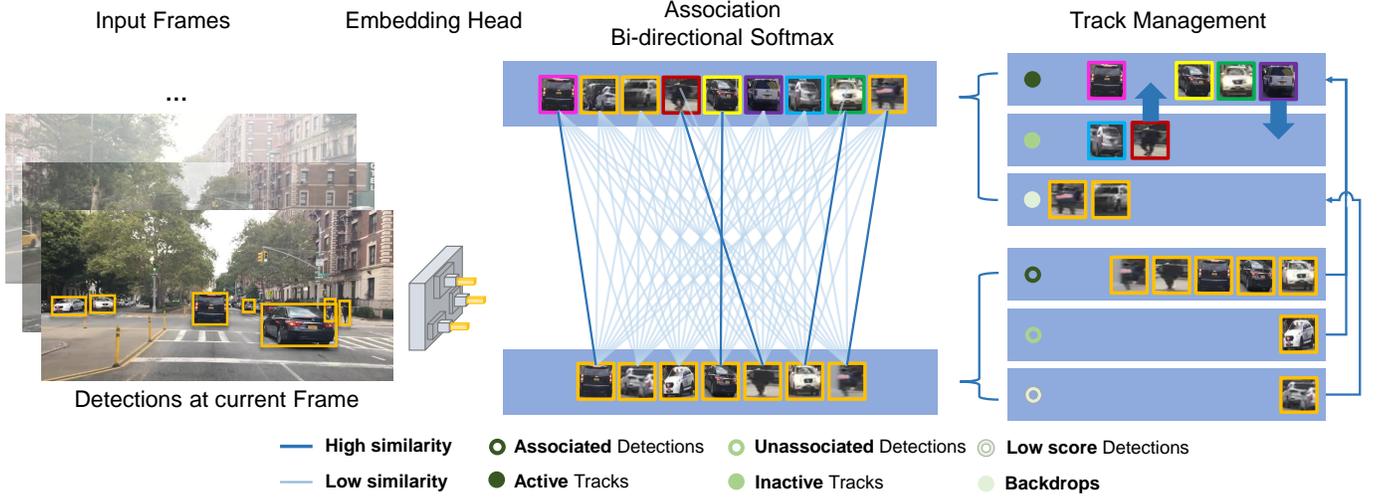


Fig. 2: **Inference pipeline.** First, we extract object detections and their corresponding feature embeddings from the current frame. Next, we use bi-softmax to measure the instance similarity between all detections and matching candidates. Finally, we associate objects with a simple nearest neighbor search in the feature space and update our track history.

**Algorithm 1** Inference pipeline of QDTrack for associating objects across a video sequence.

**Input:** frame index  $t$ , detections  $\mathbf{b}_i$ , scores  $s_i$ , detection embeddings  $\mathbf{n}_i$  for  $i = 1 \dots N$ , and track embeddings  $\mathbf{m}_j$  for  $j = 1 \dots M$

- 1: DuplicateRemoval( $\mathbf{b}_i$ )
- 2: **for**  $i = 1 \dots N, j = 1 \dots M$     # compute matching scores
- 3:     $\mathbf{f}(i, j) = \text{bisoftmax}(\mathbf{n}_i, \mathbf{m}_j)$
- 4: **end for**
- 5: **for**  $i = 1 \dots N$     # track management
- 6:     $c = \max(\mathbf{f}(i))$     # match confidence
- 7:     $j_{\text{match}} = \text{argmax}(\mathbf{f}(i))$     # matched track ID
- 8:    **if**  $c > \beta_{\text{match}}$  **and**  $s_i > \beta_{\text{obj}}$     # object match found
- 9:     updateTrack( $j_{\text{match}}, \mathbf{b}_i, \mathbf{n}_i, t$ )    # update track
- 10:    **else if**  $s_i > \beta_{\text{new}}$
- 11:     createTrack( $\mathbf{b}_i, \mathbf{n}_i, t$ )    # create new track
- 12:    **else**
- 13:     addBackdrop( $\mathbf{b}_i, \mathbf{n}_i, t$ )    # add new backdrop
- 14:    **end if**
- 15: **end for**

while each positive is considered only once. Alternatively, we can first reformulate Eq. (1) as

$$\mathcal{L}_{\text{embed}} = \log \left[ 1 + \sum_{\mathbf{k}^-} \exp(\mathbf{v} \cdot \mathbf{k}^- - \mathbf{v} \cdot \mathbf{k}^+) \right]. \quad (3)$$

Then in the multi-positive scenario, it can be extended by accumulating the positive term as

$$\mathcal{L}_{\text{embed}} = \log \left[ 1 + \sum_{\mathbf{k}^+} \sum_{\mathbf{k}^-} \exp(\mathbf{v} \cdot \mathbf{k}^- - \mathbf{v} \cdot \mathbf{k}^+) \right]. \quad (4)$$

We further adopt L2 loss as an auxiliary loss

$$\mathcal{L}_{\text{aux}} = \left( \frac{\mathbf{v} \cdot \mathbf{k}}{\|\mathbf{v}\| \|\mathbf{k}\|} - c \right)^2, \quad (5)$$

where  $c$  is 1 if the match of two samples is positive and 0 otherwise. Note the auxiliary loss aims to constrain the magnitude of the logits and cosine similarity instead of improving the performance. We sample all positive pairs and three times more negative pairs to calculate the auxiliary loss and use hard negative mining.

The entire network is jointly optimized under

$$\mathcal{L} = \mathcal{L}_{\text{det}} + \gamma_1 \mathcal{L}_{\text{embed}} + \gamma_2 \mathcal{L}_{\text{aux}}, \quad (6)$$

where  $\gamma_1$  and  $\gamma_2$  are set to 0.25 and 1.0 by default in this paper.

### 3.3 Object association and track management

Tracking objects across frames purely based on object feature embeddings introduces many challenges. False positives, ID switches, newly appeared objects, and terminated tracks all increase the matching difficulty. We here introduce our inference pipeline that utilizes instance similarity for object association and a track management scheme to address these problems. The entire pipeline is shown in Figure 2 and described in Algorithm 1.

**Duplicate removal** Most object detectors only use intra-class NMS to remove duplicate detections within each class, which results in some detections that are in the same location but with different categories. For object tracking, this is undesirable as it will create duplicate object embeddings. We instead use inter-class NMS to avoid this issue.

**Bi-directional softmax** Our main inference strategy is bi-directional matching in the feature embedding space. Assume there are  $N$  detected objects in frame  $t$  with feature embeddings  $\mathbf{n}$  and  $M$  matching candidates with feature embeddings  $\mathbf{m}$  from the past  $x$  frames. The instance similarity  $\mathbf{f}$  between objects and their matching candidates is obtained by a bi-directional softmax (bi-softmax):

$$\mathbf{f}(i, j) = \frac{1}{2} \left[ \frac{\exp(\mathbf{n}_i \cdot \mathbf{m}_j)}{\sum_{k=0}^{M-1} \exp(\mathbf{n}_i \cdot \mathbf{m}_k)} + \frac{\exp(\mathbf{n}_i \cdot \mathbf{m}_j)}{\sum_{k=0}^{N-1} \exp(\mathbf{n}_k \cdot \mathbf{m}_j)} \right]. \quad (7)$$

A high score under bi-softmax indicates that the two matched objects are each other's nearest neighbor in the feature space,

thus satisfying bi-directional consistency.  $\mathbf{f}$  can be used to directly associate objects with a simple nearest neighbor search.

**Track management** We use a track management scheme to keep track of inactive and currently active tracks and to handle the matching of objects. Active tracks are tracks that have a matching detection in the previous frame, otherwise they become inactive. Tracks that are inactivate for  $K$  frames will be removed and not be considered for matching. Detections are only considered for matching to existing tracks if the detection confidence is above a threshold  $\beta_{obj}$ . A match is determined if the matching score is higher than a threshold  $\beta_{match}$ .

Objects without a target in the feature space should not be matched to any candidates. Newly appeared objects, vanished tracks, and some false positives fall into this category. Bi-softmax can handle such objects, as it is difficult to achieve high matching scores in both directions due to the uncertainty in matching. Thus, these objects will likely obtain a low bi-softmax score and will not be matched to any existing tracks. For such objects that have a detection confidence higher than a threshold  $\beta_{new}$ , we initialize a new track instead.

Most detections with low confidence that do not match any existing tracks are false positives that introduce uncertainty to the matching process. Previous methods often directly drop them and do not consider them again. We argue that these false positives appear frequently in the following frames, which hurts tracking performance. To remedy this, we keep the unmatched objects as *backdrops* for  $L$  frames and use them as matching candidates. Detections that are matched to backdrops will thus not be matched to existing tracks. Our experiments show that backdrops can reduce the number of false positives.

## 4 EXPERIMENTS

We conduct experiments on a variety of MOT benchmarks including MOT17 [21] and MOT20 [61], DanceTrack [22], BDD100K [23], Waymo [24], and TAO [25], and compare our method extensively to the state-of-the-art. In addition, we show that we can effectively perform tracking even without tracking supervision or video data. We demonstrate the flexibility of our method by combining it with different detection methods and feature-extraction base models and conduct extensive ablation studies on all aspects of our method. Finally, we also present a straightforward extension of our method to segmentation tracking and give insights on the limitations of our method. More detailed oracle and failure case analyses are presented in the appendix.

### 4.1 Datasets

**MOT Challenge** We perform experiments on two of the MOT Challenge benchmarks, namely MOT17 [21] and MOT20 [61]. The MOT Challenge videos contain high-density public spaces such as street scenes and malls with many pedestrians, creating challenging tracking conditions with heavy occlusions. Only pedestrians are evaluated in this benchmark. Since these datasets do not provide official validation sets, we split each training video into two halves: the first half for training and the second half for validation following [4], [5], [47].

The MOT17 dataset contains 7 videos (5,316 images) for training and 7 videos (5,919 images) for testing. The video frame rate is 25 - 30 FPS. The MOT20 dataset includes heavily crowded scenes and contains 4 videos (8,931 images) for training and 4 videos (4,479 images) for testing. The video frame rate is 25 FPS.

**DanceTrack** The DanceTrack [22] benchmark is a large-scale dataset for multi-human tracking consisting mostly of group dancing videos. The dataset is unique in that by relying mostly on group dancing videos, the objects to track often have similar appearance, diverse motion, and extreme articulation. It features 40 videos for training, 25 videos for validation and 35 videos for testing, with a total of 105,855 frames captured at 20 FPS.

**BDD100K** The large-scale, diverse driving dataset BDD100K [23] contains 100,000 video sequences of dashcam driving footage. It contains several subsets with different types of annotations. We use the detection and tracking sets for training and the tracking set for evaluation. The tracking set annotates 8 categories for evaluation. It contains 1,400 videos (278k images) for training, 200 videos (40k images) for validation, and 400 videos (80k images) for testing. The detection set has 70,000 images for training. The images in the tracking set are annotated at 5 FPS.

**Waymo** Waymo open dataset [24] contains images from 5 cameras associated with 5 different directions: front, front left, front right, side left, and side right. There are 3,990 videos (790k images) for training, 1,010 videos (200k images) for validation, and 750 videos (148k images) for testing. It annotates 3 classes for evaluation. The videos are annotated at 10 FPS.

**TAO** TAO dataset [25] annotates 482 classes in total, which are a subset the classes annotated in the LVIS dataset [77]. It has 400 videos, 216 classes in the training set, 988 videos, 302 classes in the validation set, and 1419 videos, 369 classes in the test set. The classes in train, validation, and test sets may not overlap. The videos are annotated at 1 FPS. The annotated classes in TAO follow a long-tailed distribution, *e.g.*, half of the annotated instances are of class person and a sixth of the objects are of class car, while there are many classes with only few annotated instances.

### 4.2 Metrics

We use several well-established tracking metrics for evaluation.

**MOTA** The Multiple Object Tracking Accuracy (MOTA) [78] metric computes tracking accuracy in tandem with detection accuracy. It is defined as,

$$\text{MOTA} = 1 - \frac{\sum_t (m_t + f_t + e_t)}{\sum_t g_t}, \quad (8)$$

where  $t$  is the timestep,  $m_t$  is the number of misses,  $f_t$  is the number of false positives,  $e_t$  is the number of mismatches, and  $g_t$  is the number of objects. MOTA weighs detection performance more heavily than association performance. For tracking with multiple classes, we compute MOTA for each class independently then take an average over the number of classes (mMOTA).

**IDF1** The Identification  $F_1$  Score (IDF1) [79] matches ground truth and predictions on the trajectory level and computes a corresponding F1-score. It is defined as,

$$\text{IDF1} = \frac{|\text{IDTP}|}{|\text{IDTP}| + 0.5|\text{IDFN}| + 0.5|\text{IDFP}|}, \quad (9)$$

where IDTP, IDFN, and IDFP are the true positive, false negative, and false positive trajectories. IDF1 focuses on measuring association performance. Similar to MOTA, we compute an average over multiple classes for multi-class tracking (mIDF1).

**HOTA** Higher Order Tracking Accuracy (HOTA) [80] aims to fairly combine the evaluation of detection and association. Therefore, HOTA is composed of two accuracy scores, detection

TABLE 1: **Comparison to state-of-the-art on MOT Challenge benchmarks.** We benchmark our method against existing works on both MOT17 and MOT20 test sets with private detections. Datasets include CrowdHuman (CH) [62], CityPersons (CP) [63], and ETHZ [64]. † indicates using COCO pre-trained weights. † means higher is better.

Method	Detector	Base model	Datasets	MOTA †	IDF1 †	HOTA †	AssA †	DetA †	AssRe †	AssPr †	DetRe †	DetPr †
<b>MOT17</b>												
CenterTrack [4]	CenterNet [65]	DLA-34	MOT, CH	67.8	64.7	52.2	51.0	53.8	56.6	73.0	57.5	76.9
FairMOT [47] †	CenterNet [65]	DLA-34	MOT, CH, CP, ETHZ	73.7	72.3	59.3	58.0	60.9	63.6	76.3	66.0	78.5
ReMOT [66]	-	-	-	77.0	72.0	59.7	57.1	62.8	61.7	78.0	68.8	77.1
OC-SORT [44] †	YOLOX-X	Modified CSP	MOT, CH, CP, ETHZ	78.0	77.5	63.2	63.4	63.2	67.5	80.8	67.2	80.3
MAA [67]	CrowdDet [68]	R50-FPN	MOT, CH	79.4	75.9	62.0	60.2	64.2	67.3	74.0	70.9	76.4
StrongSORT [69] †	YOLOX-X	Modified CSP	MOT, CH, CP, ETHZ	79.6	79.5	64.4	64.4	64.6	71.0	78.7	70.2	78.3
ByteTrack [43] †	YOLOX-X	Modified CSP	MOT, CH, CP, ETHZ	80.3	77.3	63.1	62.0	64.5	68.2	76.0	70.1	78.1
QDTrack (Ours) †	FRCNN	R50-FPN	MOT, CH	77.2	72.2	58.8	56.2	61.8	62.6	74.1	66.6	78.1
	YOLOX-X	Modified CSP	MOT, CH	78.7	77.5	63.5	62.6	64.5	69.3	76.2	71.0	77.7
<b>MOT20</b>												
SGT [70] †	CenterNet [65]	DLA-34	MOT, CH	72.8	70.5	56.9	55.3	58.8	60.3	75.4	63.8	76.8
StrongSORT [69] †	YOLOX-X	Modified CSP	MOT, CH	73.8	77.0	62.6	64.0	61.3	69.6	80.0	65.3	81.2
MAA [67]	CrowdDet [68]	R50-FPN	MOT, CH	73.9	71.2	57.3	55.1	59.7	61.1	72.1	64.8	77.4
OC-SORT [44] †	YOLOX-X	Modified CSP	MOT, CH	75.7	76.3	62.4	62.5	62.4	67.4	79.6	66.9	80.4
ReMOT [66]	-	-	-	77.4	73.1	61.2	58.7	63.9	63.1	79.5	69.8	78.6
ByteTrack [43] †	YOLOX-X	Modified CSP	MOT, CH	77.8	75.2	61.3	59.6	63.4	66.2	74.6	69.1	78.4
QDTrack (Ours) †	YOLOX-X	Modified CSP	MOT, CH	74.7	73.8	60.0	58.9	61.4	65.7	74.8	66.4	79.1

TABLE 2: **Comparison to state-of-the-art on DanceTrack.** We compare our method to existing methods on the challenging DanceTrack test set. We use YOLOX-X [71] as our detector.

Method	HOTA †	DetA †	AssA †	MOTA †	IDF1 †
FairMOT [47]	39.7	66.7	23.8	82.2	40.8
CenterTrack [4]	41.8	78.1	22.6	86.8	35.7
TraDes [72]	43.3	74.5	25.4	86.2	41.2
TransTrack [51]	45.5	75.9	27.5	88.4	45.2
ByteTrack [43]	47.7	71.0	32.1	89.6	53.9
GTR [73]	48.0	72.5	31.9	84.7	50.3
MOTR [50]	54.2	73.5	40.2	79.7	51.5
OC-SORT [44]	55.1	80.3	38.3	92.0	54.6
QDTrack (Ours)	54.2	80.1	36.8	87.7	50.4

accuracy DetA and association accuracy AssA. DetA is defined as,

$$\text{DetA} = \frac{|\text{TP}|}{|\text{TP}| + |\text{FN}| + |\text{FP}|}, \quad (10)$$

where TP, FN, and FP are the true positive, false negative, and false positive detections. Additionally, detection recall DetRe and detection precision DetPr are used. AssA is defined as,

$$\text{AssA} = \frac{1}{|\text{TP}|} \sum_{a \in \text{TP}} \frac{|\text{TPA}(a)|}{|\text{TPA}(a)| + |\text{FNA}(a)| + |\text{FPA}(a)|}, \quad (11)$$

where TPA, FNA, and FPA are the true positive, false negative, and false positive associations. Similarly, association recall AssRe and association precision AssPr are used. HOTA is computed as a geometric mean of DetA and AssA.

### 4.3 Implementation details

Two-stage object detectors use a Region Proposal Network (RPN) to first generate a set of proposal bounding boxes, *i.e.*, Region of Interests (RoIs). We use the RoIs from the RPN for similarity learning. One-stage object detectors do not have a proposal stage and instead perform detection directly on the entire dense grid of bounding box locations. As our similarity learning protocol requires object regions, we generate them by simply using the dense detection outputs before post-processing. We follow the same box filtering procedure as the RPN [74], where we keep the most confident 1000 boxes then apply Non-maximum Suppression

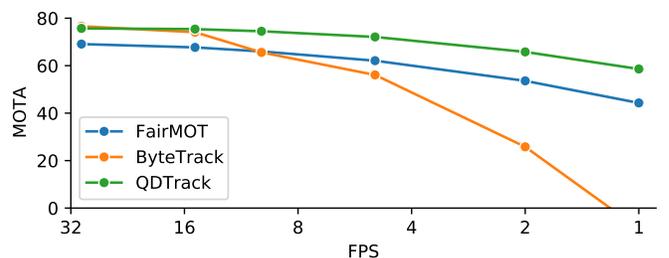


Fig. 3: **Ablation study on video frame rate.** We compare QDTrack to state-of-the-art tracking methods, namely ByteTrack [43] and FairMOT [47], on the MOT17 validation split at different video frame rates.

(NMS) with an IoU threshold of 0.7. We investigate Faster R-CNN [74] for two-stage detectors and RetinaNet [81] and YOLOX [71] for one-stage detectors in this work.

We select 128 RoIs from the key frame as training samples, and 256 RoIs from the reference frame with a positive-negative ratio of 1.0 as contrastive targets. We use IoU-balanced sampling [82] to sample negative RoIs, which better balances the sampling of hard negatives according to their IoU. We use *4conv-1fc* head with group normalization [83] to extract feature embeddings. The channel number of embedding features is set to 256 by default. We keep backdrops only from the previous frame. For association, we associate objects only when they are classified as the same category.

On MOT17 and MOT20, we follow the recent practice of [43], [44], [69] and train QDTrack with the popular YOLOX [71] detector on the union of CrowdHuman [62] and the respective MOT benchmark. On DanceTrack and BDD100K, we again follow [43] and use the same detector, but we only train on the respective dataset. For data augmentation, we follow [71] and utilize MixUp [84] and Mosaic augmentations. For our ablation studies we use Faster-RCNN in combination with ResNet-50 and FPN unless otherwise noted.

On Waymo, we use the original scale of the images for training and inference. We do not use any other data augmentation methods

TABLE 3: **Comparison to state-of-the-art on BDD100K.** We report results and compare with existing works on the BDD100K tracking validation and test set. † indicates using COCO pre-trained weights.

Split	Method	Detector	Base model	mMOTA ↑	mIDF1 ↑	MOTA ↑	IDF1 ↑	FN ↓	FP ↓	ID Sw. ↓	MT ↑	ML ↓
val	Yu <i>et al</i> [23]	FRCNN	DLA-34	25.9	44.5	56.9	66.8	122406	52372	8315	8396	3795
	DeepSORT [3]	FRCNN	R50-FPN	35.2	49.3	-	-	-	-	-	-	-
	TETer [19]	FRCNN	R50-FPN	39.1	53.3	-	-	-	-	-	-	-
	ByteTrack [43] †	YOLOX-X	Modified CSP	45.5	54.8	69.1	70.4	92805	34998	9140	9626	3005
	QDTrack (Ours)	FRCNN	R50-FPN	37.7	52.9	65.7	72.7	104861	41355	5640	9649	2874
	QDTrack (Ours) †	YOLOX-X	Modified CSP	42.1	54.3	68.2	73.3	83395	48798	8478	10925	2272
test	Yu <i>et al</i> [23]	FRCNN	DLA-34	26.3	44.7	58.3	68.2	213220	100230	14674	16299	6017
	DeepBlueAI	-	-	31.6	38.7	56.9	56.0	292063	35401	25186	10296	12266
	madamada	-	-	33.6	43.0	59.8	55.7	209339	76612	42901	16774	5004
	DeepSORT [3]	FRCNN	R50-FPN	34.0	50.2	-	-	-	-	-	-	-
	TETer [19]	FRCNN	R50-FPN	37.4	53.3	-	-	-	-	-	-	-
	ByteTrack [43] †	YOLOX-X	Modified CSP	40.1	55.8	69.9	71.3	169073	63869	15466	18057	5107
	QDTrack (Ours)	FRCNN	R50-FPN	38.7	54.1	66.5	74.0	185773	78068	10098	18167	4635
	QDTrack (Ours) †	YOLOX-X	Modified CSP	42.4	55.6	68.4	73.9	154797	89376	14282	19852	3924

TABLE 4: **Comparison to state-of-the-art on Waymo.** We show results of our method compared with existing methods on the Waymo Open tracking validation set using py-motmetrics library (top) <sup>1</sup> and test set using official evaluation (bottom). We use Faster R-CNN [74] as our detector. \* indicates methods using undisclosed detectors.

Split	Method	Category	MOTA ↑	IDF1 ↑	FN ↓	FP ↓	ID Sw. ↓	MT ↑	ML ↓	mAP ↑
val	IoU baseline [12]	Vehicle	38.3	-	-	-	-	-	-	45.8
	Tracktor++ [1], [12]	Vehicle	42.6	-	-	-	-	-	-	42.4
	RetinaTrack [12]	Vehicle	44.9	-	-	-	-	-	-	45.7
	QDTrack (Ours)	Vehicle	<b>55.6</b>	66.2	514548	214998	24309	17595	5559	<b>49.5</b>
		All	44.0	56.8	674064	264886	30712	21410	7510	40.1
Method	Split	Category	MOTA/L1 ↑	FP/L1 ↓	MisM/L1 ↓	Miss/L1 ↓	MOTA/L2 ↑	FP/L2 ↓	MisM/L2 ↓	Miss/L2 ↓
test	Tracktor [24], [75]	Vehicle	34.8	10.6	14.9	39.7	28.3	8.6	12.1	51.0
	CascadeRCNN-SORTv2*	All	50.2	7.8	2.7	39.3	44.2	<b>6.9</b>	2.4	46.5
	HorizonMOT*	All	51.0	7.5	2.4	<b>39.0</b>	<b>45.1</b>	7.1	2.3	<b>45.5</b>
	Ours (ResNet-50)	All	49.4	<b>7.4</b>	1.5	41.7	43.9	7.1	<b>1.3</b>	48.2
	Ours (ResNet-101 + DCN)	All	<b>51.2</b>	7.6	<b>1.5</b>	39.7	45.1	7.2	<b>1.3</b>	46.4

TABLE 5: **Comparison to state-of-the-art on TAO.** We evaluate and compare our method on the TAO challenge benchmark. We use Faster R-CNN [74] as our detection method. † indicates offline methods, ‡ indicates methods using additional data.

Split	Method	AP50	AP75	AP	AP50(S)	AP50(M)	AP50(L)
val	SORT_TAO [25]	13.2	-	-	-	-	-
	QDTrack (Ours)	16.1	5.0	7.0	2.4	4.6	9.6
	GTR [73] †	22.5	-	-	-	-	-
	AOA [76] ‡	25.8	-	-	-	-	-
test	SORT_TAO [25]	10.2	4.4	4.9	7.7	8.2	15.2
	QDTrack (Ours)	12.4	4.5	5.2	3.7	8.3	18.8
	GTR [73] †	20.1	-	-	-	-	-
	AOA [76] ‡	27.5	-	-	-	-	-

except random horizontal flipping and initialize the base model with ImageNet pre-trained weights for training. On TAO, we randomly select a scale between 640 to 800 and resize the shorter side of images during training. At inference time, the shorter side of the images are resized to 800. We use an LVIS [77] pre-trained model, consistent with the implementation of [25]. We freeze the detection model and only fine-tune the embedding head to extract instance representations as the annotations in TAO are incomplete. Since not all objects in the videos are annotated, fine-tuning the detection model will lead to worse performance. We provide a more detailed overview and analysis of our hyper-parameters in the appendix.

#### 4.4 Comparison to state-of-the-art

We compare our method to existing literature across five challenging multi-object tracking benchmarks.

**MOT** The official benchmark results with private detectors on MOT17 and MOT20 benchmarks are shown in Table 1. Our method achieves competitive performance on both benchmarks, despite only utilizing appearance cues for association. Notably, QDTrack obtains a high score of 63.5 HOTA on MOT17 and 60.0 HOTA on MOT20. Since the MOT benchmarks are captured at a relatively high frame rate and include only limited camera motion, 2D motion based association [43], [44], [69] works very well in this scenario. However, this only holds true for the high frame rate scenario. In Figure 3 we show that when reducing the video frame rate on the MOT17 validation split, the performance of ByteTrack [43] drops quickly and even completely fails at a frame rate of 1 FPS, while our tracker still achieves 58.6 MOTA at this frame rate. Furthermore, we show that our tracker also compares favorably to other appearance-based trackers in this regime, namely FairMOT [47], which drops to 44.3 MOTA maintaining only 64.1% of its original performance, while we maintain 77.4%.

**DanceTrack** The results on the benchmark are shown in Table 2. Surprisingly, while DanceTrack was specifically designed to provide a platform to develop MOT algorithms that rely less on visual appearance and more on motion analysis, we find that our appearance based tracker performs very well on this dataset, reaching a HOTA score only marginally behind the state-

TABLE 6: **Learning to track from static images.** We apply quasi-dense instance similarity learning on static images without tracking annotations. We use different data augmentation strategies to distort the image pair: horizontal flip (HF), multi-scale resize and crop (MS), color jittering (Color), and MixUp / Mosaic. We denote non-consistent augmentation parameters between key and reference images as ‘NC’. We evaluate performance on the BDD100K tracking validation set and compare to a model trained with video input and tracking annotations.

Input	Supervision	Augmentations					BDD100K				
		HF	HF-NC	MS-NC	Color-NC	MixUp / Mosaic	mMOTA $\uparrow$	mIDF1 $\uparrow$	MOTA $\uparrow$	IDF1 $\uparrow$	AP $\uparrow$
image	detection	✓	-	-	-	-	32.0	43.4	54.7	58.4	29.8
		✓	-	✓	-	-	35.2	47.3	62.1	66.8	32.8
		-	✓	✓	-	-	35.6	48.1	62.4	67.8	32.6
		-	✓	✓	✓	-	36.3	48.0	62.2	67.8	32.8
		-	✓	✓	✓	✓	35.2	47.7	62.4	67.5	33.0
image / video	detection	✓	-	-	-	-	32.4	45.4	56.8	60.6	32.0
		✓	-	✓	-	-	36.7	49.9	61.6	66.9	33.7
		-	✓	✓	-	-	35.1	50.4	61.8	68.0	33.5
		-	✓	✓	✓	-	36.1	50.4	62.1	68.3	34.2
		-	✓	✓	✓	✓	<b>37.7</b>	51.1	62.9	68.2	<b>35.2</b>
image / video	tracking	✓	-	-	-	-	36.6	50.8	63.5	71.5	33.0
		-	✓	✓	✓	-	36.9	52.1	64.1	71.7	34.1
		-	✓	✓	✓	✓	<b>37.7</b>	<b>52.9</b>	<b>65.7</b>	<b>72.7</b>	35.0

TABLE 7: **Ablation study on different training schedules, feature extractors, and detectors.** We train multiple QDTrack models with different detectors, feature extractors and training schedules and measure tracking performance on the BDD100K tracking validation set.

Detector	Base model	Schedule	MOTA $\uparrow$	IDF1 $\uparrow$	AP $\uparrow$
FRCNN	R50-FPN	1x	65.7	72.7	35.0
	R50-FPN	2x	65.6	73.1	35.1
	R101-FPN	1x	66.2	73.1	35.3
	R101-FPN	2x	65.6	72.7	34.6
RetinaNet	R50-FPN	1x	60.8	69.5	32.1
YOLOX-X	Modified CSP	25 epochs	<b>68.2</b>	<b>73.3</b>	<b>38.9</b>

of-the-art method OC-SORT [44] ( $-0.9$  HOTA). We achieve this score without any bells and whistles, naively applying the same configuration as in our MOT17 experiments to train on the DanceTrack dataset, following [43]. This reinforces our argument that one can in fact build a robust tracking algorithm by relying on our quasi-dense instance similarity.

**BDD100K** The main results on the BDD100K tracking validation and testing sets are in Table 3. On the validation set, QDTrack with YOLOX-X achieves 42.1 mMOTA and 54.3 mIDF1, which are the second-best results behind ByteTrack [43]. Still, QDTrack achieves much better results in IDF1 (73.3 vs. 70.4). On the test set, QDTrack with YOLOX-X achieves a high score of 42.4 mMOTA, 55.6 mIDF1, and 73.9 IDF1, outperforming all other methods by a significant margin. In particular, QDTrack outperforms ByteTrack by 2.3 mMOTA and 2.6 IDF1. QDTrack with Faster R-CNN also achieves a competitive score of 38.7 mMOTA, 54.1 mIDF1, and 74.0 IDF1, outperforming other methods using the same detector. TETer [19] is an extension of QDTrack that employs a new association strategy designed for improving long-tailed object tracking. These results demonstrate that QDTrack can perform well even on a more challenging large-scale benchmark with a simple framework.

**Waymo** Table 4 shows our main results on Waymo open dataset. We report the results on the validation set following the setup of RetinaTrack [12], which only conduct experiments on the vehicle class. We also report the overall performance for future

comparison. We report the results on the test set via official rules. Our method outperforms all baselines on both validation set and test set. We obtain 44.0 MOTA and 56.8 IDF1 on the validation set and 49.4 MOTA/L1 and 43.9 MOTA/L2 on the test set. The performance of vehicle on the validation set is 10.7, 13.0, and 17.4 points higher than RetinaTrack [12], Tractor++ [1], [12], and IoU baseline [12], respectively. Our model with ResNet-101 and deformable convolution (DCN) [85] has state-of-the-art performance on the test benchmark, which is on par with the champion of Waymo 2020 2D Tracking Challenge (HorizonMOT) despite only using a simple single model.

**TAO** The results for TAO are shown in Table 5. We obtain 16.1 AP50 on the validation set and 12.4 AP50 on the test set. The results are 2.9 points and 2.2 points higher than TAO’s baseline. Although we only boost the overall performance by 2 to 3 points, we outperform the baseline by a large margin on frequent classes, *i.e.*, 38.6 points vs. 18.5 points on person. This improvement is not well represented in the standard evaluation metrics of TAO, since it averages per-class scores across hundreds of classes. GTR [73] and AOA [76] are recent methods proposed to tackle long-tail multi-object tracking. Although they outperform our method, GTR is an offline method and AOA utilizes separate ReID networks trained on additional data.

#### 4.5 Learning to track from static images

Since our quasi-dense instance similarity learning is agnostic to how the image pair is generated during training, we investigate how we can leverage static images where no association annotations are available. Inspired by recent literature in self-supervised representation learning [55], [56], we experiment with different data augmentations on static images to learn discriminative instance representations from static input. In particular, for a given training sample in a detection dataset, we generate two distorted images via data augmentation techniques. We find that random horizontal flip (HF), multi-scale resize and crop (MS), color jittering (Color), and MixUp / Mosaic augmentations are the most suitable for our use-case. If the augmentation parameters are not shared across the key and reference view, we denote it with ‘NC’ (non-consistent). We only use MixUp / Mosaic with consistent parameters in order to compose the same images between key and reference views and

TABLE 8: **Ablation study on quasi-dense matching and inference strategy.** We investigate the contribution of various components on the BDD100K tracking validation set. All models are comparable on detection performance. D. R. means duplicate removal. (P) means results of the class “pedestrian”.

Quasi-Dense		Metric	Matching candidates		MOTA $\uparrow$	IDF1 $\uparrow$	mMOTA $\uparrow$	mIDF1 $\uparrow$	MOTA(P) $\uparrow$	IDF1(P) $\uparrow$
one-positive	multi-positive		D. R.	Backdrops						
-	-	<i>cosine</i>	-	-	60.4	63.0	34.0	47.9	37.6	49.7
✓	-	<i>cosine</i>	-	-	61.5	66.8	35.5	50.0	40.5	52.7
-	✓	<i>cosine</i>	-	-	62.5	67.8	36.2	50.0	44.0	54.3
-	✓	<i>bi-softmax</i>	-	-	62.9	70.0	35.4	48.5	45.5	58.8
-	✓	<i>bi-softmax</i>	✓	-	63.2	70.1	36.4	50.4	45.5	58.3
-	✓	<i>bi-softmax</i>	✓	✓	<b>63.5</b>	<b>71.5</b>	<b>36.6</b>	<b>50.8</b>	<b>46.7</b>	<b>60.2</b>
					<b>+3.1</b>	<b>+8.5</b>	<b>+2.6</b>	<b>+2.9</b>	<b>+9.1</b>	<b>+10.5</b>

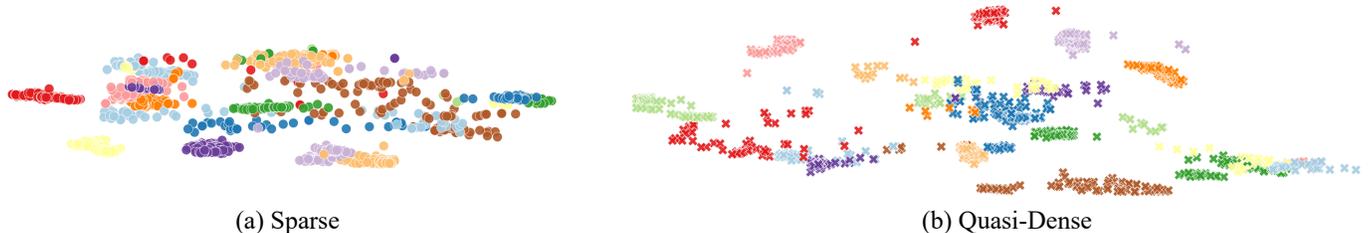


Fig. 4: **Instance embedding space visualization.** We visualize the instance embedding space learned via (a) sparse matching and (b) quasi-dense matching using t-SNE. We show ground truth embedding identities as color and plot embedding vectors sampled from a sequence in the BDD100K tracking validation set.

TABLE 9: **Ablation study on location and motion cues.** We investigate if our method benefits from using a range of motion priors on the BDD100K tracking validation set. We integrate bounding box IoU, a simple linear motion model, and displacement regression into the association procedure.

Appearance	IoU	Motion	Regression	mMOTA $\uparrow$	mIDF1 $\uparrow$
-	✓	-	-	26.3	36.0
-	✓	✓	-	27.7	38.5
-	✓	-	✓	28.6	39.3
✓	-	-	-	<b>36.6</b>	<b>50.8</b>
✓	✓	-	-	36.3	49.8
✓	✓	✓	-	36.4	49.9
✓	✓	-	✓	36.4	50.1

thus be able to match objects across them. To train our models, we utilize the detection (image) and tracking (video) splits of BDD100K. Note that the tracking split contains much more data, thus influencing the detection performance.

The results of our experiments are shown in Table 6. We use Faster-RCNN [74] as the detector with ResNet-50 [86] and FPN [58] as the base model and evaluate the tracking performance on the BDD100K tracking benchmark [23]. We observe that when we only apply consistent HF, the tracking performance is far behind the version trained with full tracking supervision. By adding in non-consistent augmentations and MixUp / Mosaic, we can narrow this gap and achieve comparable accuracy to the fully supervised model. In particular, we exceed the mMOTA of the fully supervised baseline trained without augmentations besides HF when training on the same amount of training data by a significant margin. This clearly shows that not only detection, but also association benefits greatly from the data augmentation, and that with proper data augmentation, our similarity learning scheme can track objects effectively while trained on static images alone. If we use the same amount of training data, we indeed rival the performance of the

best supervised model, shown by the small gap in mMOTA ( $-0.1$  points).

In addition, we observe that the data augmentation scheme can also benefit the supervised models, reaching a much higher score than in our initial work [17] without changing the network architecture ( $+1.2$  points in mMOTA,  $+2.1$  points in mIDF1). The increase in mIDF1 highlights the benefit of data augmentation to the robustness of instance similarity learning.

## 4.6 Ablation studies

We conduct ablation studies on the validation set of BDD100K [23], where we investigate the importance of the major model components for training and inference procedures.

**Different object detectors, feature extractors, and training schedules** We combine our method with different object detectors and feature extractors to verify the flexibility of our instance similarity learning scheme. In Table 7, we show the tracking performance of our method with ResNet-50, ResNet-101 [86], as well as the modified CSPNet [87] on the tracking validation set of BDD100K. We combine those feature extractors with a Faster-RCNN [74] detector and observe that ResNet-101 achieves the best performance with 66.2 MOTA, 73.1 IDF1, and 35.3 AP.

In addition, we apply our method on two more base object detection models, namely RetinaNet [81] and YOLOX [71]. Both methods produce reasonable results, and the YOLOX model achieves the best overall scores with 68.2 MOTA, 73.3 IDF1, and 38.9 AP. It shows that our method can work independent of feature extractor or base detection model. Finally, we also experiment with different training schedules. We investigate the effect of longer training, increasing the epochs from 12 (1x schedule) to 24 (2x schedule) and 25. Note that we use the extensive data augmentation techniques presented in section 4.5 in this ablation study to counteract overfitting when training with longer schedules. We find that increasing the number of epochs does not help the

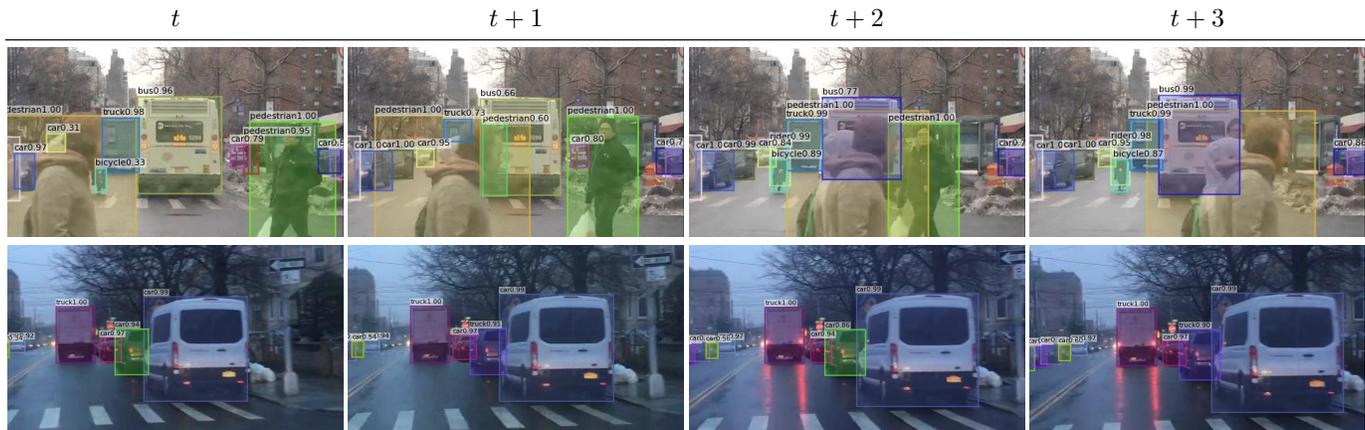


Fig. 5: **Illustration of failure cases.** We illustrate the two most common failure cases of our method (best viewed digitally). In the top, we can see that the bus (light green and violet) switches identity due to extreme occlusion by pedestrians. In the bottom, we observe that the pickup truck (green and purple) switches identity when the class prediction changes between ‘truck’ and ‘car’. Note that we still re-identify the pickup truck once the class predictions match.

TABLE 10: **Extension of our method to segmentation tracking.** We show segmentation tracking results of our method on the BDD100K segmentation tracking validation set. I: ImageNet. C: COCO. S: Cityscapes. B: BDD100K.

Method	Pretraining	mMOTSA $\uparrow$	mIDF1 $\uparrow$	ID sw. $\downarrow$
SORT [2]	I, C, S	10.3	21.8	15951
MaskTrackRCNN [46]	I, C, S	12.3	26.2	9116
STEm-Seg [88]	I, C, S	12.2	25.4	8732
QDTrack (ours)	I, B	25.6	45.2	980
PCAN [18]	I, B	27.4	45.1	876

smaller ResNet models, but is beneficial for training very large models like YOLOX-X.

**Importance of quasi-dense matching** The results are presented in the top sub-table of Table 8. We use a Faster R-CNN detector with ResNet-50 base model. MOTA and IDF1 are calculated over all instances without considering categories as overall evaluations. We use cosine distance to calculate the similarity scores during the inference procedure. Compared to learning with sparse ground truths, quasi-dense tracking improves the overall IDF1 by 4.8 points (63.0% to 67.8%). The significant improvement on IDF1 indicates quasi-dense tracking greatly improves the feature embeddings and enables more accurate associations.

We then analyze the improvements in detail. In the table, we can observe that when we match each training sample to more negative samples and train the feature space with Eq. (1), the IDF1 is significantly improved by 3.4 points. This improvement contributes 70% to the total improved 4.8 points IDF1. This experiment shows that more contrastive targets, even most of them are negative samples, can improve the feature learning process. The multiple-positive contrastive learning following Equation (4) further improves the IDF1 by 1 point (66.8% to 67.8%).

**Importance of bi-softmax** We investigate how different inference strategies influence the performance. As shown in the bottom of Table 8, replacing cosine similarity by bi-softmax improves overall IDF1 by 2.2 points and the IDF1 of pedestrian by 4.5 points. This experiment also shows that the one-to-one constraint further strengthens the estimated similarity.

**Importance of matching candidates** Duplicate removal and

backdrops improve IDF1 by 1.5 points. Overall, our training and inference strategies improve the IDF1 by 8.5 points (63.0% to 71.5%). The total number of ID switches is decreased by 30%. Especially, the MOTA and IDF1 of pedestrian are improved by 9.1 points and 10.5 points respectively, which further demonstrate the power of quasi-dense contrastive learning.

**Combinations with motion and location** Finally, we try to add location and motion priors to understand whether they are still helpful when we have good feature embeddings for measuring similarity. These experiments follow the procedures in Tractor [1] and use the same detector for fair comparisons. As shown in Table 9, without appearance features, the tracking performance is consistently improved with the introduction of additional information. However, these cues barely enhance the performance of our approach. Our method yields the best results when only using appearance embeddings. The results indicate that our instance feature embeddings are sufficient for multiple object tracking with the effective quasi-dense matching, which greatly simplifies the inference pipeline.

**Inference speed** To understand the runtime efficiency, we profile our method on a single NVIDIA RTX 3090 graphics card. Because it only adds a lightweight embedding head to the detector, our method only causes marginal overhead in inference speed. With an input size of  $1296 \times 720$  and a Faster R-CNN detector with ResNet-50 base model on BDD100K, the inference time is 61 ms, equating to 16.3 FPS. However, the embedding extractor consumes 3 ms, representing only 5% of the total runtime.

## 4.7 Embedding visualizations

We use t-SNE to visualize the embeddings trained with sparse matching and our quasi-dense matching and show them in Figure 4. The instances are selected from a video in BDD100K tracking validation set. The same instance is shown with the same color. We observe that it is easier to separate objects in the feature space of quasi-dense matching.

## 4.8 Segmentation tracking

Owing to the simplicity of our method, we can extend it to instance segmentation tracking in a straightforward manner. To do so, we

simply add a Mask R-CNN [59] mask prediction head to the existing network architecture and use a pre-trained QDTrack model trained on the BDD100K tracking set to fine-tune the mask head on MOTSA data. In particular, BDD100K provides a subset for the segmentation tracking task. There are 154 videos in the training set, 32 videos in the validation set, and 37 videos in the test set. Table 10 shows the results on the BDD100K segmentation tracking task compared to other methods. QDTrack achieves 25.6 mMOTSA and 45.2 mIDF1. PCAN [18] is an extension of QDTrack that utilizes a prototypical appearance module to further improve segmentation. We observe that QDTrack based models achieve much better performance than previous methods.

#### 4.9 Limitations

While our method gains in simplicity and generality by solely relying on instance similarity learning, we also identify certain challenges that arise with this paradigm. In particular, we observe that our model struggles with rapid changes in object appearance, e.g., through partial occlusion. Also, since our model relies on discrete class labels to aid the matching process, we observe that classification errors can lead to truncated object tracks. These cases are illustrated in Figure 5. In addition, inaccurate object localization can lead to difficulties in association when regions within a bounding box cover background and/or other objects, thus impeding accurate instance embedding extraction. For more detailed failure case and oracle analysis, please refer to the appendix.

## 5 CONCLUSION

We present QDTrack, a tracking method based on quasi-dense instance similarity learning. The key idea behind our method is to utilize all object regions in an image for similarity learning, in contrast to previous methods that only use sparse ground-truth regions as similarity supervision. We observe that the feature embedding space we learn from quasi-dense matches is much better suited to discriminate instances, allowing for a simple tracking framework that associates objects via nearest neighbor search in the embedding space without bells and whistles. Our method can be easily coupled with most existing object detectors and feature extractors for end-to-end training, and learns effective instance similarity even without video input or tracking annotations.

## REFERENCES

- [1] P. Bergmann, T. Meinhardt, and L. Leal-Taixé, “Tracking without bells and whistles,” *arXiv preprint arXiv:1903.05625*, 2019.
- [2] A. Bewley, Z. Ge, L. Ott, F. T. Ramos, and B. Uprocft, “Simple online and realtime tracking,” in *International Conference on Image Processing*, 2016.
- [3] N. Wojke, A. Bewley, and D. Paulus, “Simple online and realtime tracking with a deep association metric,” in *International Conference on Image Processing*, 2017.
- [4] X. Zhou, V. Koltun, and P. Krähenbühl, “Tracking objects as points,” in *European Conference on Computer Vision*, 2020.
- [5] Z. Wang, L. Zheng, Y. Liu, and S. Wang, “Towards real-time multi-object tracking,” *arXiv preprint arXiv:1909.12605*, 2019.
- [6] D. Ramanan and D. A. Forsyth, “Finding and tracking people from the bottom up,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2003.
- [7] E. Bochinski, V. Eiselein, and T. Sikora, “High-speed tracking-by-detection without using image information,” in *IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*, 2017.
- [8] A. Milan, S. Roth, and K. Schindler, “Continuous energy minimization for multitarget tracking,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2013.
- [9] W. Choi and S. Savarese, “Multiple target tracking in world coordinate with single, minimally calibrated camera,” in *European Conference on Computer Vision*, 2010.
- [10] C. Feichtenhofer, A. Pinz, and A. Zisserman, “Detect to track and track to detect,” in *IEEE International Conference on Computer Vision*, 2017.
- [11] J. Peng, C. Wang, F. Wan, Y. Wu, Y. Wang, Y. Tai, C. Wang, J. Li, F. Huang, and Y. Fu, “Chained-tracker: Chaining paired attentive regression results for end-to-end joint multiple-object detection and tracking,” in *European Conference on Computer Vision*, 2020.
- [12] Z. Lu, V. Rathod, R. Votel, and J. Huang, “Retinatrack: Online single stage joint detection and tracking,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [13] A. v. d. Oord, Y. Li, and O. Vinyals, “Representation learning with contrastive predictive coding,” *arXiv preprint arXiv:1807.03748*, 2018.
- [14] R. Hadsell, S. Chopra, and Y. LeCun, “Dimensionality reduction by learning an invariant mapping,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2006.
- [15] K. Sohn, “Improved deep metric learning with multi-class n-pair loss objective,” in *Advances in Neural Information Processing Systems*, 2016.
- [16] Z. Wu, Y. Xiong, S. X. Yu, and D. Lin, “Unsupervised feature learning via non-parametric instance discrimination,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [17] J. Pang, L. Qiu, X. Li, H. Chen, Q. Li, T. Darrell, and F. Yu, “Quasi-dense similarity learning for multiple object tracking,” in *IEEE/CVF Conference on Computer Vision and Pattern Recognition*, June 2021.
- [18] L. Ke, X. Li, M. Danelljan, Y.-W. Tai, C.-K. Tang, and F. Yu, “Prototypical cross-attention networks for multiple object tracking and segmentation,” *Advances in Neural Information Processing Systems*, vol. 34, 2021.
- [19] S. Li, M. Danelljan, H. Ding, T. E. Huang, and F. Yu, “Tracking every thing in the wild,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, Oct 2022.
- [20] H.-N. Hu, Y.-H. Yang, T. Fischer, T. Darrell, F. Yu, and M. Sun, “Monocular quasi-dense 3d object tracking,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.
- [21] A. Milan, L. Leal-Taixé, I. Reid, S. Roth, and K. Schindler, “Mot16: A benchmark for multi-object tracking,” *arXiv preprint arXiv:1603.00831*, 2016.
- [22] P. Sun, J. Cao, Y. Jiang, Z. Yuan, S. Bai, K. Kitani, and P. Luo, “Dancetrack: Multi-object tracking in uniform appearance and diverse motion,” *arXiv preprint arXiv:2111.14690*, 2021.
- [23] F. Yu, H. Chen, X. Wang, W. Xian, Y. Chen, F. Liu, V. Madhavan, and T. Darrell, “Bdd100k: A diverse driving dataset for heterogeneous multitask learning,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [24] P. Sun, H. Kretzschmar, X. Dotiwalla, A. Chouard, V. Patnaik, P. Tsui, J. Guo, Y. Zhou, Y. Chai, B. Caine, V. Vasudevan, W. Han, J. Ngiam, H. Zhao, A. Timofeev, S. Ettinger, M. Krivokon, A. Gao, A. Joshi, Y. Zhang, J. Shlens, Z. Chen, and D. Anguelov, “Scalability in perception for autonomous driving: Waymo open dataset,” 2019.
- [25] A. Dave, T. Khurana, P. Tokmakov, C. Schmid, and D. Ramanan, “Tao: A large-scale benchmark for tracking any object,” in *European Conference on Computer Vision*, 2020.
- [26] L. Zhang, Y. Li, and R. Nevatia, “Global data association for multi-object tracking using network flows,” in *IEEE Conference on Computer Vision and Pattern Recognition*, 2008.
- [27] A. Ess, B. Leibe, K. Schindler, and L. Van Gool, “A mobile vision system for robust multi-person tracking,” in *2008 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2008, pp. 1–8.
- [28] W. Choi, “Near-online multi-target tracking with aggregated local flow descriptor,” in *IEEE International Conference on Computer Vision*, 2015.
- [29] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg, “Multiple hypothesis tracking revisited,” in *IEEE International Conference on Computer Vision*, 2015.
- [30] B. Xiao, H. Wu, and Y. Wei, “Simple baselines for human pose estimation and tracking,” in *European Conference on Computer Vision*, 2018.
- [31] F. Yu, W. Li, Q. Li, Y. Liu, X. Shi, and J. Yan, “POI: multiple object tracking with high performance detection and appearance feature,” in *European Conference on Computer Vision Workshop*, 2016.
- [32] D. Held, S. Thrun, and S. Savarese, “Learning to track at 100 FPS with deep regression networks,” in *European Conference on Computer Vision*, 2016.
- [33] A. Osep, W. Mehner, M. Mathias, and B. Leibe, “Combined image- and world-space tracking in traffic scenes,” in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017.
- [34] S. Sharma, J. A. Ansari, J. K. Murthy, and K. M. Krishna, “Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking,” in *IEEE International Conference on Robotics and Automation*, 2018.

- [35] D. Mitzel and B. Leibe, "Taking mobile multi-object tracking to the next level: People, unknown objects, and carried items," in *European Conference on Computer Vision*. Springer, 2012, pp. 566–579.
- [36] J. Luiten, T. Fischer, and B. Leibe, "Track to reconstruct and reconstruct to track," *IEEE Robotics and Automation Letters*, 2020.
- [37] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg, "Multiple hypothesis tracking revisited," in *IEEE International Conference on Computer Vision*, 2015.
- [38] L. Leal-Taixé, C. Canton-Ferrer, and K. Schindler, "Learning by tracking: Siamese CNN for robust target association," in *IEEE Conference on Computer Vision and Pattern Recognition Workshop*, 2016.
- [39] B. Yang and R. Nevatia, "An online learned CRF model for multi-target tracking," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2012.
- [40] J. Son, M. Baek, M. Cho, and B. Han, "Multi-object tracking with quadruplet convolutional neural networks," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [41] A. Sadeghian, A. Alahi, and S. Savarese, "Tracking the untrackable: Learning to track multiple cues with long-term dependencies," in *IEEE International Conference on Computer Vision*, 2017.
- [42] A. Milan, S. H. Rezatofighi, A. R. Dick, I. D. Reid, and K. Schindler, "Online multi-target tracking using recurrent neural networks," in *The AAAI Conference on Artificial Intelligence*, 2017.
- [43] Y. Zhang, P. Sun, Y. Jiang, D. Yu, Z. Yuan, P. Luo, W. Liu, and X. Wang, "Bytetrack: Multi-object tracking by associating every detection box," *arXiv preprint arXiv:2110.06864*, 2021.
- [44] J. Cao, X. Weng, R. Khirodkar, J. Pang, and K. Kitani, "Observation-centric sort: Rethinking sort for robust multi-object tracking," *arXiv preprint arXiv:2203.14360*, 2022.
- [45] D. Held, J. Levinson, and S. Thrun, "Precision tracking with sparse 3d and dense color 2d data," in *2013 IEEE International Conference on Robotics and Automation*. IEEE, 2013, pp. 1138–1145.
- [46] L. Yang, Y. Fan, and N. Xu, "Video instance segmentation," in *IEEE International Conference on Computer Vision*, 2019.
- [47] Y. Zhang, C. Wang, X. Wang, W. Zeng, and W. Liu, "Fairmot: On the fairness of detection and re-identification in multiple object tracking," *arXiv preprint arXiv:2004.01888*, 2020.
- [48] A. Hermans, L. Beyer, and B. Leibe, "In defense of the triplet loss for person re-identification," *arXiv preprint arXiv:1703.07737*, 2017.
- [49] T. Meinhardt, A. Kirillov, L. Leal-Taixé, and C. Feichtenhofer, "Trackformer: Multi-object tracking with transformers," *arXiv preprint arXiv:2101.02702*, 2021.
- [50] F. Zeng, B. Dong, T. Wang, X. Zhang, and Y. Wei, "Motr: End-to-end multiple-object tracking with transformer," *arXiv preprint arXiv:2105.03247*, 2021.
- [51] P. Sun, J. Cao, Y. Jiang, R. Zhang, E. Xie, Z. Yuan, C. Wang, and P. Luo, "Transtrack: Multiple-object tracking with transformer," *arXiv preprint arXiv: 2012.15460*, 2020.
- [52] P. Bachman, R. D. Hjelm, and W. Buchwalter, "Learning representations by maximizing mutual information across views," in *Advances in Neural Information Processing Systems*, 2019.
- [53] O. J. Hénaff, A. Srinivas, J. De Fauw, A. Razavi, C. Doersch, S. Eslami, and A. v. d. Oord, "Data-efficient image recognition with contrastive predictive coding," *arXiv preprint arXiv:1905.09272*, 2019.
- [54] Y. Tian, D. Krishnan, and P. Isola, "Contrastive multiview coding," *arXiv preprint arXiv:1906.05849*, 2019.
- [55] K. He, H. Fan, Y. Wu, S. Xie, and R. Girshick, "Momentum contrast for unsupervised visual representation learning," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2020.
- [56] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton, "A simple framework for contrastive learning of visual representations," *arXiv preprint arXiv:2002.05709*, 2020.
- [57] Y. Sun, C. Cheng, Y. Zhang, C. Zhang, L. Zheng, Z. Wang, and Y. Wei, "Circle loss: A unified perspective of pair similarity optimization," *arXiv preprint arXiv:2002.10857*, 2020.
- [58] T.-Y. Lin, P. Dollár, R. B. Girshick, K. He, B. Hariharan, and S. J. Belongie, "Feature pyramid networks for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [59] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *IEEE International Conference on Computer Vision*, 2017.
- [60] L. Bertinetto, J. Valmadre, J. F. Henriques, A. Vedaldi, and P. H. Torr, "Fully-convolutional siamese networks for object tracking," in *European Conference on Computer Vision*, 2016.
- [61] P. Dendorfer, H. Rezatofighi, A. Milan, J. Shi, D. Cremers, I. Reid, S. Roth, K. Schindler, and L. Leal-Taixé, "Mot20: A benchmark for multi object tracking in crowded scenes," *arXiv preprint arXiv:2003.09003*, 2020.
- [62] S. Shao, Z. Zhao, B. Li, T. Xiao, G. Yu, X. Zhang, and J. Sun, "Crowdhuman: A benchmark for detecting human in a crowd," *arXiv preprint arXiv:1805.00123*, 2018.
- [63] S. Zhang, R. Benenson, and B. Schiele, "Citypersons: A diverse dataset for pedestrian detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 3213–3221.
- [64] A. Ess, B. Leibe, K. Schindler, and L. Van Gool, "A mobile vision system for robust multi-person tracking," in *2008 IEEE Conference on Computer Vision and Pattern Recognition*, 2008, pp. 1–8.
- [65] X. Zhou, D. Wang, and P. Krähenbühl, "Objects as points," in *arXiv preprint arXiv:1904.07850*, 2019.
- [66] F. Yang, X. Chang, S. Sakti, Y. Wu, and S. Nakamura, "Remot: A model-agnostic refinement for multiple object tracking," *Image and Vision Computing*, vol. 106, p. 104091, 2021.
- [67] D. Stadler and J. Beyerer, "Modelling ambiguous assignments for multi-person tracking in crowds," in *2022 IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*, 2022.
- [68] X. Chu, A. Zheng, X. Zhang, and J. Sun, "Detection in crowded scenes: One proposal, multiple predictions," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [69] Y. Du, Y. Song, B. Yang, and Y. Zhao, "Strongsort: Make deepsort great again," *arXiv preprint arXiv:2202.13514*, 2022.
- [70] J. Hyun, M. Kang, D. Wee, and D.-Y. Yeung, "Detection recovery in online multi-object tracking with sparse graph tracker," *arXiv preprint arXiv:2205.00968*, 2022.
- [71] Z. Ge, S. Liu, F. Wang, Z. Li, and J. Sun, "Yolox: Exceeding yolo series in 2021," *arXiv preprint arXiv:2107.08430*, 2021.
- [72] J. Wu, J. Cao, L. Song, Y. Wang, M. Yang, and J. Yuan, "Track to detect and segment: An online multi-object tracker," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2021, pp. 12 352–12 361.
- [73] X. Zhou, T. Yin, V. Koltun, and P. Krähenbühl, "Global tracking transformers," *arXiv preprint arXiv:2203.13250*, 2022.
- [74] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," in *Advances in Neural Information Processing Systems*, 2015.
- [75] C. Kim, F. Li, and J. M. Rehg, "Multi-object tracking with neural gating using bilinear LSTM," in *European Conference on Computer Vision*, 2018.
- [76] F. Du, B. Xu, J. Tang, Y. Zhang, F. Wang, and H. Li, "1st place solution to eccv-tao-2020: Detect and represent any object for tracking," *arXiv preprint arXiv:2101.08040*, 2021.
- [77] A. Gupta, P. Dollar, and R. Girshick, "Lvis: A dataset for large vocabulary instance segmentation," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [78] K. Bernardin and R. Stiefelhagen, "Evaluating multiple object tracking performance: The clear mot metrics," *EURASIP Journal on Image and Video Processing*, vol. 2008, 01 2008.
- [79] E. Ristani, F. Solera, R. Zou, R. Cucchiara, and C. Tomasi, "Performance measures and a data set for multi-target, multi-camera tracking," in *European conference on computer vision*. Springer, 2016, pp. 17–35.
- [80] J. Luiten, A. Osep, P. Dendorfer, P. Torr, A. Geiger, L. Leal-Taixé, and B. Leibe, "Hota: A higher order metric for evaluating multi-object tracking," *International journal of computer vision*, 2021.
- [81] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *IEEE International Conference on Computer Vision*, 2017.
- [82] J. Pang, K. Chen, J. Shi, H. Feng, W. Ouyang, and D. Lin, "Libra r-cnn: Towards balanced learning for object detection," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [83] Y. Wu and K. He, "Group normalization," in *European Conference on Computer Vision*, 2018.
- [84] H. Zhang, M. Cisse, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *arXiv preprint arXiv:1710.09412*, 2017.
- [85] J. Dai, H. Qi, Y. Xiong, Y. Li, G. Zhang, H. Hu, and Y. Wei, "Deformable convolutional networks," in *IEEE International Conference on Computer Vision*, 2017.
- [86] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2016.
- [87] C.-Y. Wang, H.-Y. Mark Liao, Y.-H. Wu, P.-Y. Chen, J.-W. Hsieh, and I.-H. Yeh, "Cspnet: A new backbone that can enhance learning capability of cnn," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, 2020, pp. 390–391.
- [88] A. Athar, S. Mahadevan, A. Osep, L. Leal-Taixé, and B. Leibe, "Stem-seg: Spatio-temporal embeddings for instance segmentation in videos," in *European Conference on Computer Vision*. Springer, 2020, pp. 158–177.

## APPENDIX A ADDITIONAL DETAILS

We detail the training settings and hyper-parameters used for each benchmark investigated in the main paper. We also provide details regarding the augmentations used during training.

**Training setting** On MOT17 and MOT20, we train QDTrack with YOLOX [71] on the union of CrowdHuman [62] and the respective MOT benchmark for 80 epochs, following recent practice [43], [44], [69]. We use an image scale of  $1440 \times 800$  for MOT17 and  $1600 \times 896$  for MOT20. We use a batch size of 32 with a learning rate of 0.0005, and we use a cosine annealing learning rate schedule, ending at a learning rate 0.05 times the original. We use an exponential learning rate warm-up for one epoch. For augmentations, we turn off MixUp and Mosaic for the last ten epochs. We also use an exponential moving average (EMA), as done in [71]. On DanceTrack, we use the same training setup as in MOT17, except for training for 12 epochs. On BDD100K, we use mostly the same training setup as in MOT17, except for training for 25 epochs, using a batch size of 48, learning rate of 0.00075, and not turning off augmentations.

**Hyper-parameters** The detailed hyper-parameters are shown in Table 11. As our object association only relies on appearance, it is robust to different motion patterns in different datasets. The experiments share similar tracking parameters except for TAO, since TAO uses 3D mAP instead of CLEAR MOT and HOTA metrics for evaluation. On TAO,  $\beta_{\text{new}}$  and  $\beta_{\text{obj}}$  are set to 0.0001 to obtain a high recall. Considering the numerous tracks that results using these thresholds, we do not maintain backdrops.

**Augmentation parameters** For Mosaic, we sample the center of the mosaic image in the range (0.5, 1.5). For the main results on the benchmarks, we also use random affine transformation with a rotation degree in the range (-10.0, 10.0), translation factor in the range (-0.1, 0.1), scale factor in the range (0.5, 1.5), and shear degree in the range (-2.0, 2.0). For MixUp, we jitter the additional image by a factor in the range (0.5, 1.5) and flip it with a probability of 0.5. We also flip the resulting combined image with a probability of 0.5. We additionally apply random resizing with a scale range (0.5, 1.5) while maintaining the aspect ratio and random cropping. Finally, we apply color jitter with a brightness factor in the range (0.875, 1.125), contrast factor in the range (0.5, 1.5), saturation factor in the range (0.5, 1.5), and hue shift in the range  $(-0.2\pi, 0.2\pi)$ .

## APPENDIX B ADDITIONAL DETAILS FOR MOT CHALLENGE

For challenging scenes with heavy occlusions (*i.e.*, MOT17 and MOT20), object association with only appearance cues can be very challenging, as there are large overlaps between objects. We utilize several additional techniques to mitigate these issues. First, object appearance cues immediately after re-appearance can be unreliable for association, leading to a higher number of ID switches and lower ASSRe score. We address this by introducing a near-online merging strategy. For each object that did not match to any previous tracks and initialized a new track, in each of the  $t$  subsequent frames we merge its current track with a vanished track if their matching score is higher than a threshold  $\beta_{\text{merge}}$ . We use distance thresholding with distance  $d_{\text{merge}}$  to ignore objects that are too far away. This enables us to utilize more reliable appearance features moments after re-appearance for matching. We use  $t = 10$ ,  $\beta_{\text{merge}} = 0.5$ ,  $d_{\text{merge}} = 50$ .

TABLE 11: **Hyper-parameters used in each benchmark.** We include both tracking and detection parameters.

Parameter	MOT17	MOT20	DanceTrack	BDD100K	Waymo	TAO
$\beta_{\text{new}}$	0.75	0.75	0.8	0.5	0.8	0.0001
$\beta_{\text{obj}}$	0.3	0.3	0.6	0.35	0.5	0.0001
$\beta_{\text{match}}$	0.5	0.5	0.5	0.5	0.5	0.5
$K$	30	30	20	10	10	10
$L$	1	1	1	1	1	-
$m$	0.5	0.5	0.8	0.8	0.8	0.8
Det. confidence	0.1	0.001	0.1	0.1	0.05	0.0001
Det. NMS threshold	0.7	0.7	0.7	0.65	0.7	0.5

TABLE 12: **Ablation study of association strategies used for MOT Challenge.** We evaluate on the MOT17 validation set.

Distance Threshold	Tracklet Merging	Linear Interpolation	MOTA $\uparrow$	IDF1 $\uparrow$	HOTA $\uparrow$
-	-	-	76.1	73.6	63.5
✓	-	-	76.2	74.5	64.0
✓	✓	-	76.0	76.0	64.4
✓	✓	✓	<b>76.8</b>	<b>76.2</b>	<b>64.8</b>

Additionally, due to the high frame rate and low object motion in the MOT benchmarks, we use distance thresholding to reduce ambiguities during association by ignore matching candidates that are greater than a distance  $d$  away. We use  $d = 50$ . Following [43], we also perform linear interpolation to recover bounding boxes of fully-occluded objects.

We provide an ablation study of the aforementioned techniques on the MOT benchmarks, MOT17 and MOT20. The results on the MOT17 validation set are shown in Table 12. Using distance thresholding can improve IDF1 from 73.6 to 74.5 (+0.9). Performing tracklet merging can improve IDF1 from 74.5 to 76.0 (+1.5). Linear interpolation can further improve all metrics.

## APPENDIX C ADDITIONAL ABLATION STUDIES

**Momentum of the embeddings** Assume there is an existing track and its embedding is  $E_0$ . This track is associated to an object on the current frame and its embedding is  $E_1$ . The new embedding of this track will be  $m * E_1 + (1 - m) * E_0$ , where  $m$  is the momentum. The momentum does not improve the results too much but it considers the history of embeddings. We show the ablation studies of different values of momentum in Table 13. The models for this table are re-trained so the results are slightly different from the results in the main paper.

**Sensitivity of  $\gamma_1$  and  $\gamma_2$  in Eq. 6** We found  $\gamma_2$  does not change the final results while  $\gamma_1$  does. If  $\gamma_1$  is higher than 0.5, the performance will drop, but does not matter if it is lower than 0.5.

## APPENDIX D ORACLE ANALYSIS

We investigate the performances of two types of oracles on the BDD100K tracking validation set: detection oracle and tracking oracle. For the detection oracle, we directly extract feature embeddings of the ground truth objects in each frame and associate them using our method. For the tracking oracle, we use ground truth tracking labels to associate the detected objects.

**Detection oracle** The results are shown in Table 14. We can observe that all MOTAs are higher than 94%, and some of them are

TABLE 13: **Ablation study of momentum of the embeddings.** We evaluate on the BDD100K tracking validation set.

Momentum	mMOTA $\uparrow$	mIDF1 $\uparrow$	MOTA $\uparrow$	IDF1 $\uparrow$
0.6	37.0	50.9	63.3	71.4
0.7	37.0	50.9	63.3	71.3
0.8	37.0	50.7	63.3	71.1
0.9	37.0	50.6	63.3	70.8
1.0	37.0	50.5	63.3	70.5

even close to 100%. This is because we use the ground truth boxes directly so that the number of false negatives and false positives are close to 0.

The metric IDF1 and ID Switches can measure the performance of identity consistency. The average IDF1 over the 8 classes is 88.8%, which is 38 points higher than our result. The gaps on classes “car” and “pedestrian” are only 11.1 points and 19.3 points between oracle results and our results respectively, while gaps on other classes are exceeding 30 points. These results show that if highly accurate detection results are provided, our method can obtain robust feature embeddings and associate objects effectively. However, the huge performance gaps also indicate the demand of promoting detection algorithms in the video domain. We also notice that the total number of ID switches in the oracle experiment is higher than ours. This is due to the high object recalls in the oracle experiments, as more detected instances may introduce more ID switches accordingly.

**Tracking oracle** The results are shown in Table 15. We can observe that when associating object directly with tracking labels, the mIDF1 is only boosted by 4.3 points. This promising oracle analysis shows the effectiveness of our method and indicates that our method is bounded more by detection performance than tracking performance.

## APPENDIX E FAILURE CASE ANALYSIS

Our method can distinguish different instances even they are similar in appearance. However, there are still some failure cases. We show them below with figures, in which we use yellow color to represent false negatives, red color to represent false positives, and cyan color to represent ID switches. The float number at the corner of each box indicates the detection score, while the integer indicates the object identity number. We use green dashed box to highlight the objects we want to emphasize.

**Object classification** Inaccurate classification confidence is the main distraction for the association procedure because false negatives and false positives destroy the one-to-one matching constraint. As shown in Figure 6, the false negatives are mainly small objects or occluded objects under crowd scenes. The false positives are objects that have similar appearances to annotated objects, such as persons in the mirror or advertising board, etc.

Inaccurate object category is a less frequent distraction caused by classification. The class of the instance may switch between different categories, which mostly belong to the same super-category. Figure 7 shows an example. The category of the highlighted object changes from “rider” to “pedestrian” when the bicycle is occluded. Our method fails in this case because we require the associated objects have the same category.

These failure cases caused by object classification suggest that improvements could be achieved via leveraging video object

detection algorithms, *i.e.* exploiting temporal information to improve the detector, thus obtaining better tracking performance.

**Object truncation/occlusion** Object truncation/occlusion causes inaccurate object localization. As shown in Figure 8, the highlighted objects are truncated by other objects. The detector detects two objects. One of them is a false positive box that only covers a part of the object. The other one is a box with a lower detection score but covers the entire object. This case may influence the association process if the two boxes have similar feature embeddings.

An instance may have totally different appearances before and after occlusion that result in low similarity scores. As shown in Figure 9, only the front of the car appears before occlusion, while only the rear of the car appears after occlusion. Our method can associate two boxes if they cover the same discriminative regions of an object, not necessarily the exact same region. However, if two boxes cover totally different regions of the object, they will have a low matching score.

Another corner case is the extreme high-level truncation. As shown in Figure 10, the highly truncated objects only appear a little when they just enter or leave the camera view. We cannot distinguish different instances effectively according to the limited appearance information.

## APPENDIX F VISUALIZATIONS

We show the visualizations of different instance patches during the testing procedure in Figure 11. The detected objects in each frame are matched to prior objects via bi-directional softmax. The prior objects include tracks in the consecutive frame, vanished tracks, and backdrops. We annotate them with different colors. Each detected object is enclosed by the same color of its matched object. We can observe that most false positives in the current frame are matched to backdrops, which demonstrates keeping backdrops during the matching procedure helps reduce the number of false positives.

## APPENDIX G QUALITATIVE RESULTS

We show some qualitative results of our method on BDD100K dataset and MOT17 dataset in Figure 12 and Figure 13, respectively. The results are sampled from a certain interval for illustrative purposes.

TABLE 14: **Detection oracle analysis.** The numbers in the round brackets mean the gaps between oracle results and our results.

Category	MOTA $\uparrow$	IDF1 $\uparrow$	MOTP $\uparrow$	FN $\downarrow$	FP $\downarrow$	ID Sw. $\downarrow$	MT $\uparrow$	ML $\downarrow$
Pedestrian	94.3	79.5 (+19.3)	99.8	1	1	3226	3506	0
Rider	95.8	88.5 (+40.4)	99.9	0	0	107	134	0
Car	97.7	86.1 (+11.1)	99.9	0	0	7716	13189	0
Bus	99.2	93.0 (+31.2)	100.0	0	0	72	196	0
Truck	98.8	90.3 (+33.8)	100.0	0	0	340	726	0
Bicycle	88.2	79.5 (+31.8)	98.7	8	8	470	243	0
Motorcycle	97.0	94.5 (+37.8)	99.8	0	0	27	44	0
Train	99.4	98.7 (+98.7)	100.0	0	0	2	6	0
All	96.3	88.8 (+38.0)	99.8	9	9	11960	18044	0

TABLE 15: **Tracking oracle analysis.** The numbers in the round brackets mean the gaps between oracle results and our results.

Category	MOTA $\uparrow$	IDF1 $\uparrow$	MOTP $\uparrow$	FN $\downarrow$	FP $\downarrow$	ID Sw. $\downarrow$	MT $\uparrow$	ML $\downarrow$
Pedestrian	54.7	71.2 (+11.0)	77.6	14990	10095	755	1835	367
Rider	31.4	52.6 (+4.5)	76.6	1390	242	115	16	56
Car	74.3	82.9 (+7.9)	84.1	54585	31014	2309	8759	1141
Bus	38.2	65.8 (+4.0)	86.1	3532	2031	57	61	41
Truck	37.0	60.9 (+4.4)	84.7	12719	4259	247	149	239
Bicycle	30.6	55.6 (+7.9)	75.4	2031	714	125	60	58
Motorcycle	14.6	51.7 (-5.0)	76.4	443	292	35	10	18
Train	-0.6	0.0 (+0.0)	0.0	308	2	0	0	6
All	35.0	55.1 (+4.3)	70.1	89998	48649	3643	10890	1926

Fig. 6: **Failure cases caused by inaccurate classification confidences.** The objects enclosed by yellow rectangles are false negatives, and the objects enclosed by red rectangles are false positives.Fig. 7: **Failure case caused by inaccurate object category.** The category of the highlighted object changes from “rider” to “pedestrian” due to the occlusion of the bicycle. They cannot be associated because they do not satisfy the category consistency.



Fig. 8: **Inaccurate object localization caused by truncation.** The red false positive box only covers part of the object, while the yellow box covers the entire object. They may have similar feature embeddings thus influencing the association procedure.

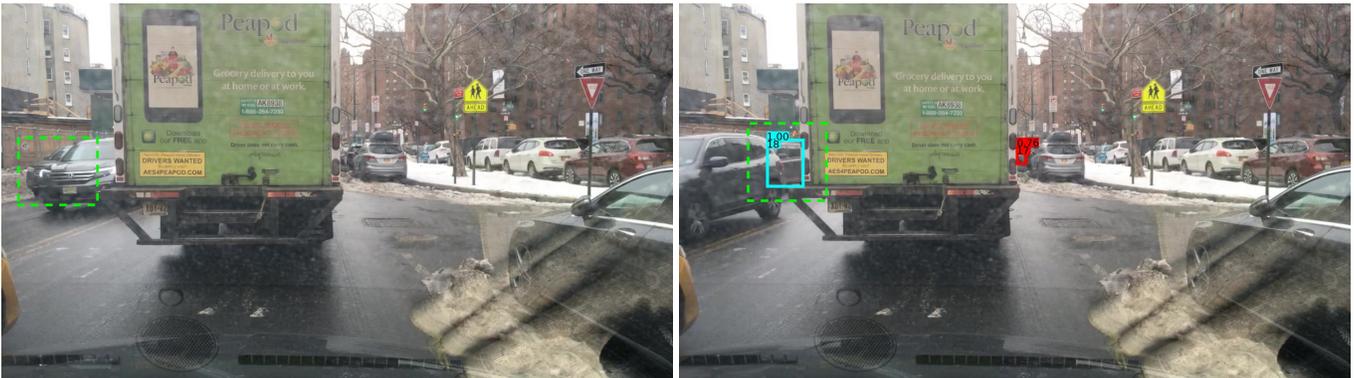


Fig. 9: **Occlusion in different regions of the same object.** Two detected objects in different frames cover totally different regions of the object thus having low appearance similarity.



Fig. 10: **Extreme high-level truncation.** Our method cannot distinguish different instances effectively according to the limited appearance information in highly truncated objects.

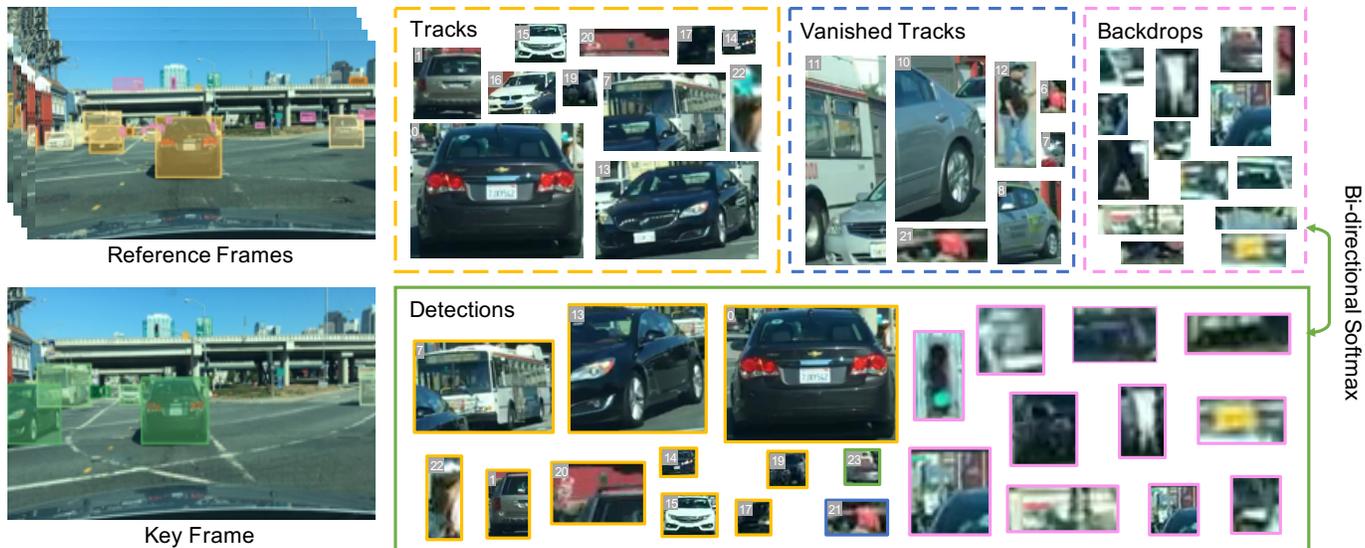


Fig. 11: **Visualizations of different instance patches during inference.** The detected objects in the current frame are matched to tracklets in the consecutive frame, vanished tracklets, and backdrops via bi-directional softmax.

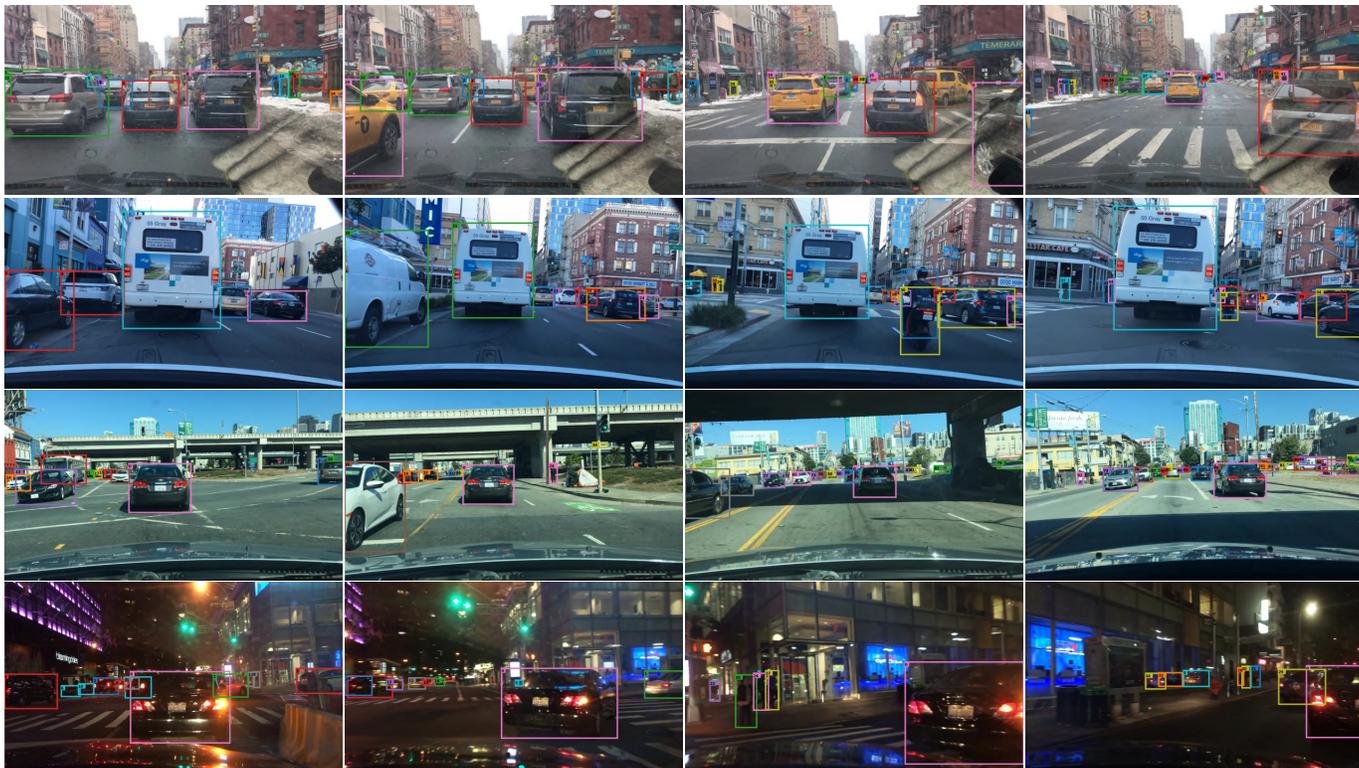


Fig. 12: Qualitative results of our method on BDD100K.



Fig. 13: Qualitative results of our method on MOT17.