

TransTrack: Multiple-Object Tracking with Transformer

Peize Sun¹, Yi Jiang², Rufeng Zhang³, Enze Xie¹, Jinkun Cao⁴,
 Xinting Hu⁵, Tao Kong², Zehuan Yuan², Changhu Wang², Ping Luo¹

¹The University of Hong Kong ²ByteDance AI Lab ³Tongji University

⁴Carnegie Mellon University ⁵Nanyang Technological University

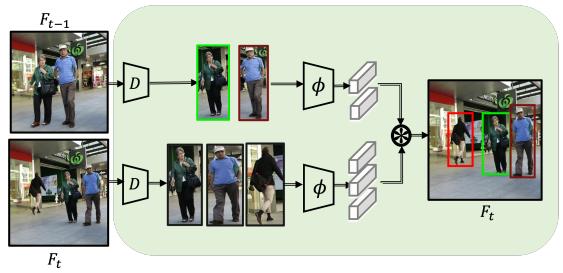
Abstract

Multiple-object tracking(MOT) is mostly dominated by complex and multi-step tracking-by-detection algorithm, which performs object detection, feature extraction and temporal association, separately. Query-key mechanism in single-object tracking(SOT), which tracks the object of the current frame by object feature of the previous frame, has great potential to set up a simple joint-detection-and-tracking MOT paradigm. Nonetheless, the query-key method is seldom studied due to its inability to detect new-coming objects.

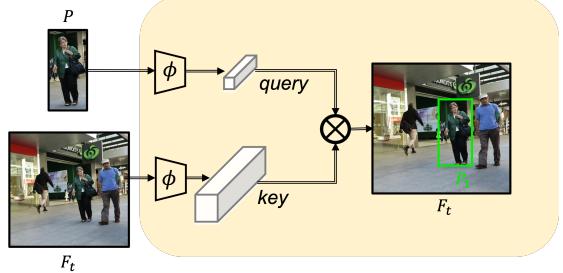
*In this work, we propose **TransTrack**, a baseline for MOT with Transformer. It takes advantage of query-key mechanism and introduces a set of learned object queries into the pipeline to enable detecting new-coming objects. TransTrack has three main advantages: (1) It is an online joint-detection-and-tracking pipeline based on query-key mechanism. Complex and multi-step components in the previous methods are simplified. (2) It is a brand new architecture based on Transformer. The learned object query detects objects in the current frame. The object feature query from the previous frame associates those current objects with the previous ones. (3) For the first time, we demonstrate a much simple and effective method based on query-key mechanism and Transformer architecture could achieve competitive 65.8% MOTA on the MOT17 challenge dataset. We hope TransTrack can provide a new perspective for multiple-object tracking. The code is available at: <https://github.com/PeizeSun/TransTrack>.*

1. Introduction

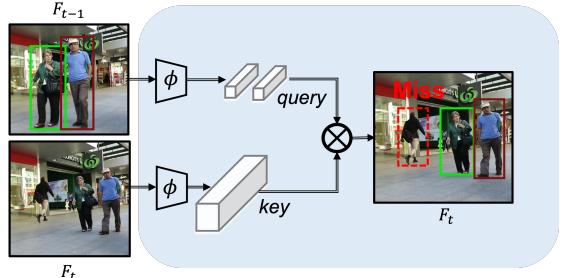
Video-based scene understanding and human behavior analysis are essential for current computer vision systems to understand the world at a high level. Aiming to estimate the trajectories of objects of interest in videos, **Object Tracking** is one significant task applied in many practical real applications, such as visual surveillance, public secu-



(a) Complex tracking-by-detection MOT pipeline.



(b) Simple query-key SOT pipeline.



(c) Query-key pipeline has great potential to setup a simple MOT method. However, it will miss new-coming objects.

Figure 1: Motivation of TransTrack. The dominant MOT method is the complex multi-step tracking-by-detection pipeline. Query-key mechanism in SOT pipeline is potential to set up a much simple MOT pipeline, however, it will miss new-coming objects. TransTrack is aimed to take advantage of query-key mechanism and to detect new-coming objects. The pipeline is shown in Figure 2.

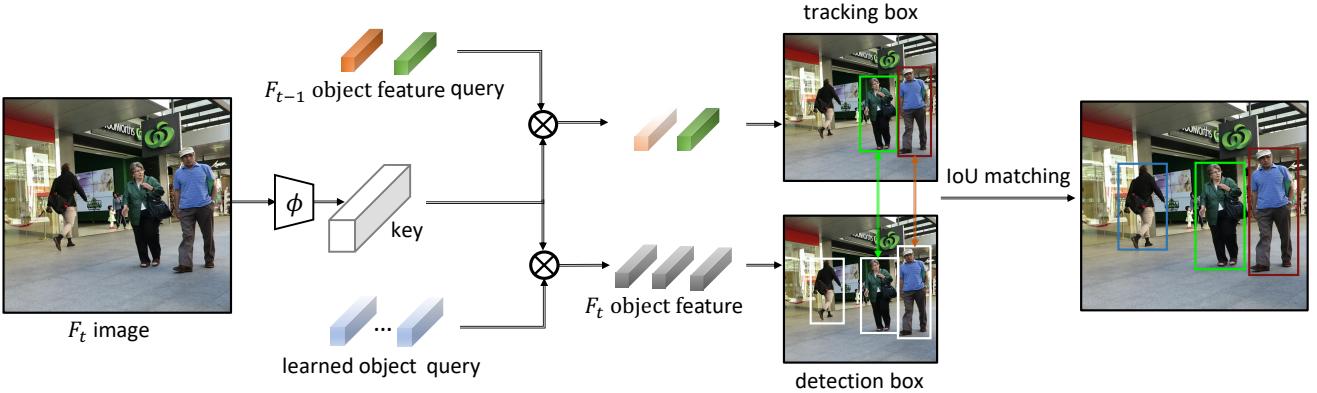


Figure 2: **Pipeline of TransTrack.** Both object feature query from the previous frame and learned object query are taken as input query. The image feature map are common key. The learned object query detects objects in the current frame. The object feature query from the previous frame associates objects of the current frame with the previous ones. This process is performed sequentially over all adjacent frames and finally completes multiple-object tracking task.

rity, sports video analysis, and human-computer interaction.

According to the number of tracked objects, Object Tracking can be divided into two directions: **Single-Object Tracking (SOT)** and **Multiple-Object Tracking (MOT)**. In recent years, SOT has made great progress because of the emerging of deep siamese networks [3, 35, 21, 20], where the correlation between the object target and image regions are captured and trained in a much simple and effective way. However, the current MOT methods have been suffering from the model complexity and computational cost due to the multi-stage pipeline [42, 34, 39], as shown in Figure 1a. Object detection and re-identification are performed separately. They can not benefit each other and bring challenges in either unordered object pairs between two consecutive frames or incomplete detection objects in each frame. To tackle these problems in MOT, a joint-detection-and-tracking framework is needed.

Reviewing SOT, we emphasize that siamese network is actually **Query-Key** mechanism, where the object target is the query and the image regions are the keys, as shown in Figure 1b. For the same object, its feature in different frames is highly similar, which enables query-key mechanism to output ordered object sets. Aiming to borrow the merits from SOT, an intuitive strategy is to introduce the query-key mechanism in MOT, e.g., objects feature of the previous frame as query and image feature of the current frame as the key, as shown in Figure 1c. However, merely transferring the vanilla query-key mechanism in SOT into the MOT task leads to obviously poor performance, especially the FN metric. The reason is that when a new object comes into the current frame, its feature is not in the query, leading to the missing of the new-coming objects. Therefore, a natural question to ask is: Is it possible to design an MOT framework that is based on query-key mechanism

to output ordered object sets; meanwhile, detect the new-coming objects.

In this paper, we present a new MOT framework for joint-detection-and-tracking, referred to **TransTrack**, which takes advantage of the query-key mechanism to track pre-existing objects in the current frame and detect new-coming objects. The overall pipeline is indicated as Figure 2. TransTrack is built on Transformer architecture [36], a widely-used entity of query-key mechanism. The input key is the feature map of the current frame. The input query is both object feature from the previous frame and a set of learned object queries. The learned object query is a set of parameters, trained together with all other parameters in the network. It is used to detect new-coming objects in the current frame and output *detection boxes*. The object feature from the previous frame is the object feature vectors generated in the detection process of the previous frame. It is used to locate the pre-existing objects in the current frame and output *tracking boxes*. After simple matching between detection boxes and tracking boxes, the final results are output.

Our method is simple, straightforward, and easy to implement. Both the tracking boxes and detection boxes can be viewed as object detection of the current frame. It allows us to simultaneously train these two sub-networks, rather than separately optimize detection and re-identification network, as in tracking-by-detection methods [42, 34]. On the challenging MOT dataset [25], TransTrack achieves 65.8 MOTA, comparable performance with state-of-the-art frameworks. Our contributions are as follows:

- We introduce an online joint-detection-and-tracking MOT pipeline based on query-key mechanism, simplifying complex and multi-step components in the previous methods.

- We perform a detailed analysis of complicated tracking scenarios, demonstrating that both learned object query and object feature from the previous frame can be used as query input of Transformer architecture to detect and track simultaneously.
- We present that our method achieves comparable performance with state-of-the-art models without bells and whistles. We hope our work could provide a new perspective for multiple-object tracking.

2. Related Work

As our work is to introduce query-key mechanism into the multi-object tracking model, we first review the applications of query-key mechanism in object detection and single-object tracking, then dive into related works on multi-object tracking.

Query-Key mechanism in Object Detection. Query-key mechanism has been successfully applied in object detection areas for its entities of self-attention and cross-attention [36], *i.e.*, Relation Network [14], DETR [5], Deformable DETR [48]. Among them, DETR reasons about the relations of the object queries and the global image context to directly output the final set of predictions in parallel. DETR streamlines the detection pipeline, effectively removing the need for non-maximum suppression procedures and anchor generation.

We notice that these object detection frameworks can be intuitively applied to multiple-object tracking pipeline to provide object detection.

Query-Key mechanism in Single-Object Tracking. Recently, siamese-network-based [35, 3, 21, 20, 50, 38] single-object trackers have received significant attention for their excellent tracking accuracy and efficiency. These trackers formulate visual tracking as a cross attention correlation problem, hoping to better leverage the merits of deep networks from the end-to-end learning process without needing any post-procedure. Siamese-network-based trackers consist of two network branches, one for the object template and the other for the image search region. Then the trackers fuse the feature map of the two network branches and produce a similarity map. Following the basic idea, SINT [35] and SiamFC [3] adopt to learn the similarity between the object target and candidate image patches in an offline manner. SiamRPN [21] and DaSiamRPN [50] continue to improve with a region proposal network. With the predefined anchor boxes, SiamRPN can capture the scale changes of objects effectively. SiameseRPN++ [20] proposes a new model architecture with deeper layers to perform layer-wise and depthwise aggregations, which not only further improves the accuracy but also reduces the model size. SiamMask [38] presents a new architecture that

performs both visual object tracking and semi-supervised video object segmentation simultaneously.

The object template and the image search region are a query-key pair, so these mentioned related works could be categorized as query-key mechanism in SOT task. The wide application proves the effectiveness of object feature on the previous frame as the query to locate its own position on the following frame. This makes simultaneous detection and association possible. However, in MOT scenario, there are always new objects coming into the image view. These new-coming objects cannot be detected since they have no corresponding feature query. Only using object feature from the beginning frame as query will make all objects born on the following frames missed. This is why query-key mechanism in MOT trails. Instead, tracking-by-detection is the mainstream method in MOT.

Tracking-by-detection. Most state-of-the-art multi-object trackers [43, 39, 7, 42, 47, 41, 4] follow the tracking-by-detection paradigm. The paradigm is first using the object detectors such as [22, 23, 28] to localize all target objects in the images by several boxes, and secondly cropping the images according to the detected boxes. Tracking is then a problem of bounding box association. Re-ID features and Intersection over Unions (IoU) of the bounding boxes are usually used for the box association. First, the IOU or feature-space-based distance is computed for the boxes, then Kalman Filter [40] and Hungarian algorithm [19] are used to accomplish the box association task. SORT [4] tracks bounding boxes using a Kalman Filter [40] and associates each bounding box with its highest overlapping detection in the current frame using the Hungarian algorithm. DeepSORT [41] uses the appearance features from deep convolutional networks to compute the association cost in SORT. Lifted-Multicut[34] leverage person-re identification features and human pose features.

The advantage of these methods is that they use the most suitable model for each task, respectively. Besides, they crop the image patches according to the detected bounding boxes and resize them to the same size before feeding into deep networks. In this way, they can reduce the scale variations of objects. These approaches have achieved the best performance on the public dataset [42]. However, these methods have two drawbacks. First, the detector and deep appearance models are trained separately, and thus the detector and deep appearance models can not take advantage of each other to get better performance. Second, two separate networks in existing methods greatly increase the model complexity and computational cost. To tackle these problems, joint-detection-and-tracking methods are required.

Joint-detection-and-tracking. Recently joint-detection-and-tracking frameworks have begun to attract

more attention. D&T [9] uses a siamese network with the current and past frame as input and predict inter-frame offsets between bounding boxes. Integrated-Detection [44] uses tracked bounding boxes as additional region proposals to enhance detection, followed by bipartite-matching-based bounding-box association. Tracktor [1] directly uses the previous frame tracking results as region proposals and then applies the bounding box regression to provide tracking results, thus eliminating the box association procedure. JDE [39] and FairMOT [43] learn the object detection task and appearance embedding task from a shared neural network backbone. CenterTrack [45] is a simultaneous detection and tracking algorithm which localizes objects and predicts their offsets to the previous frame. ChainedTracker [27] chains paired bounding boxes regression results estimated from overlapping nodes, of which each node covers two adjacent frames. In video object detection, FGFA [49] uses optical flow to warp intermediate features from the previous frames to accelerate inference. T-CNN [16] feeds stacked consecutive frames into the network and performed detection for a whole video segment.

Our method is intuitively a joint-detection-and-tracking pipeline. The difference is that all of these previous works adopt anchor-based [28] or point-based [46] detection framework, where the tracked boxes were used as proposals or points. Instead, we build the pipeline based on query-key mechanism and the tracked object is used as query.

3. TransTrack

We hypothesize that a desirable tracking model outputs object sets which are **complete** and **ordered**. To this end, TransTrack takes both learned object query and object feature from the previous frame as input query. The learned object query is decoded into detection boxes on each frame to provide common object detection results. The object feature from the previous frames is decoded into tracking boxes. TransTrack performs tracking associations based on the tracking boxes and detection boxes on the same frame. This enables a simple box IoU matching strategy to associate two consecutive frames. The architecture details of TransTrack is shown in Figure 3.

3.1. Pipeline

TransTrack includes one encoder to generate composite feature map and two parallel decoders to perform object detection and object propagation. Given the detection boxes and tracking boxes on the same frame, box IoU matching is used to obtain the final tracking result.

Architecture. TransTrack is built on Transformer [36], a widely-used entity of query-key mechanism. It includes

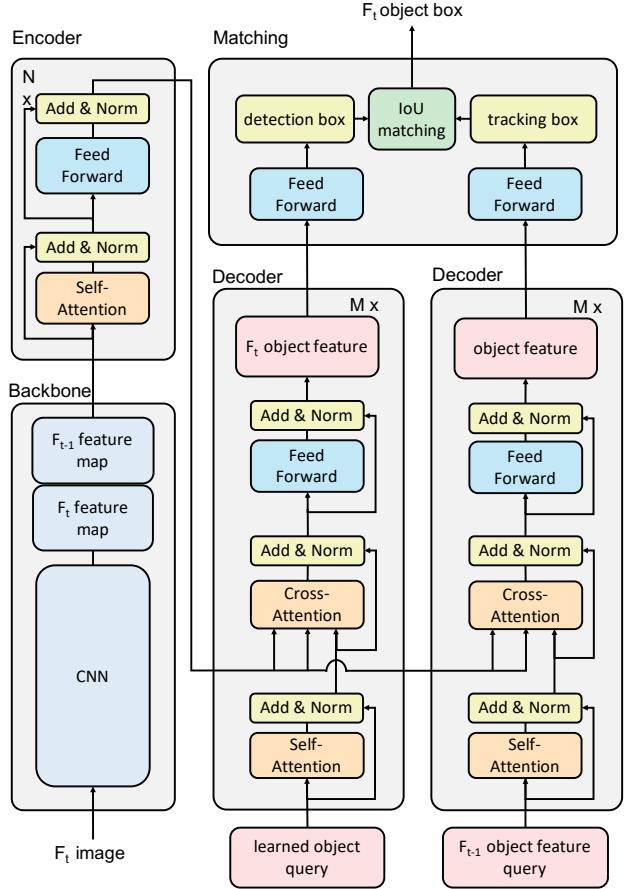


Figure 3: **The architecture details of TransTrack.** First, the current frame image is input to CNN backbone to extract feature map. Then, both the current frame feature map and the previous one are fed into encoder to generate composite feature. Next, learned object query is decoded into detection boxes and object feature of the previous frame is decoded into tracking boxes. Finally, IoU matching is employed to associate detection boxes to tracking boxes.

encoder and decoder, both of which are composed of stacked multi-head attention layers and point-wise fully connected layers. Multi-head attention is called self-attention if the input query and the input key are the same, otherwise, cross-attention. The point-wise fully connected layers are called feed-forward networks, consisting of linear transformations and non-linear activation functions. In transformer architecture, the encoder generates keys and the decoder takes as input task-specific queries. This query-key attention mechanism makes it suitable for sequence tasks and achieves excellent performance, such as natural lan-

guage processing and video understanding [36, 10].

The encoder of TransTrack takes the composed feature maps of two consecutive frames as input to catch useful correlations, as shown in the encoder block of Figure 3. To avoid duplicated computation, the extracted features of the current frame are temporarily saved and then re-used for the next frame.

Two parallel decoders are employed in TransTrack. Feature maps generated from the encoder are used as common keys by the two decoders. The two decoders are designed to perform object detection and object propagation, respectively. Specifically, a decoder takes learned object query as input and predicts *detection boxes* on the current frame. The other decoder takes the object feature from previous frames as input and predicts the locations of the corresponding objects on the current frame, namely, *tracking boxes*.

Object Detection. TransTrack leverages the concept of learned object query to perform object detection in each frame. The learned object query is firstly proposed in DETR [5], a new object detector based on query-key mechanism. The learned object query is a set of learnable parameters, trained together with all other parameters in the network. During detection, the key is the global feature maps generated from the input image and the learned object query looks up objects of interest in the image and outputs the final detection predictions, termed as “detection boxes”. The left-hand decoder block in Figure 3 illustrates the object detection stage of TransTrack.

Object Propagation. Given detected objects in the previous frame, TransTrack propagates these objects to the current frame by the way of object propagation, shown in the right-hand decoder block in Figure 3. The decoder has basically the same architecture as the left-hand one but it takes object feature from previous frames as input query instead. This inherited object feature conveys the appearance and location information of previously seen objects, so this decoder could well locate the position of the corresponding object on the current frame and output “tracking boxes”.

Box Association. Provided the detection boxes and tracking boxes on the same frame, TransTrack uses box IoU matching method to get the final tracking result, as shown in Figure 3. Since both detection boxes and tracking boxes are the locations of objects in the same frame, there is only slight offset between them. It allows a simple matching strategy, namely box IoU matching, to associate the two sets of boxes. Applying the Kuhn-Munkres (KM) algorithm [19] to IoU similarity of detection boxes and tracking boxes, detection boxes are matched to tracking boxes. Those unmatched detection boxes are added as new objects.

3.2. Training

Training Data. The training data of TransTrack could be the same as most other tracking methods, in which two consecutive frames or two randomly selected frames from a short sequence are used as training samples. Furthermore, training data could also be the static image [45], where the adjacent frame is simulated by randomly scaling and translating the static image.

Training Loss. Both tracking boxes and detection boxes can be viewed as object detection of the current frame. It allows us to simultaneously train two decoders by the same training loss.

For training detection boxes, TransTrack applies set prediction loss [5, 48, 33, 32, 37] on the object set of predictions of classification and box coordinates. Set-based loss produces an optimal bipartite matching between predictions and ground truth objects. Following [5, 48, 33, 32, 37], the matching cost is defined as follows:

$$\mathcal{L} = \lambda_{cls} \cdot \mathcal{L}_{cls} + \lambda_{L1} \cdot \mathcal{L}_{L1} + \lambda_{giou} \cdot \mathcal{L}_{giou} \quad (1)$$

where \mathcal{L}_{cls} is focal loss [23] of predicted classifications and ground truth category labels, \mathcal{L}_{L1} and \mathcal{L}_{giou} are L1 loss and generalized IoU loss [29] between normalized center coordinates and height and width of predicted boxes and ground truth box, respectively. λ_{cls} , λ_{L1} and λ_{giou} are coefficients of each component. The training loss is the same as the matching cost except that only performed on matched pairs. The final loss is the sum of all pairs normalized by the number of objects inside the training batch.

For training tracking boxes, optimal bipartite matching is removed and the matching index is directly from detection boxes in the previous frame. The training loss is the same as detection boxes.

3.3. Inference

TransTrack first performs object detection in the first frame, where the composite feature maps are two copies of the feature map of the first frame. Then TransTrack operates object propagation and box association from the first frame to the second frame. This process is performed sequentially over all adjacent frames and finally completes the multiple-object tracking task.

Track Rebirth. We introduce track rebirth in the inference procedure of TransTrack in order to enhance robustness to occlusions and short-term disappearing [1, 45, 27]. Specifically, if a tracking box is unmatched, it keeps as an “inactive” tracking box until it remains unmatched for K consecutive frames. Inactive tracking boxes can be matched to detection boxes and regain their ID. Similar to [45], we choose $K = 32$.

4. Experiments

4.1. Datasets and evaluation metrics

We conduct experiments on MOT17 dataset [25], which contains 7 training sequences and 7 test sequences. Only pedestrians are annotated. The MOT dataset does not provide an official split. Similar to [45], we split each training sequences into two halves, and use the first half for training, second for validation, in ablation study. Benchmark evaluation is training on the whole training set and evaluating on test set.

Tracking performance is measured by the widely-used MOT metrics [2], including Multiple-Object Tracking Accuracy(MOTA), Multiple-Object Tracking Precision(MOTP), the total number of False Negatives (FN), False Positives (FP), Identity Switches (IDs), and the percentage of Mostly Tracked Trajectories (MT), Mostly Lost Trajectories (ML). ID F1 Score (IDF1) is also used to measure the trajectory identity accuracy. Among them, MOTA is the primary metric to measure the overall performance both in detection and tracking:

$$\text{MOTA} = 1 - \frac{\sum_t (\text{FP}_t + \text{FN}_t + \text{IDs}_t)}{\sum_t \text{GT}_t} \quad (2)$$

where GT_t is the number of ground-truth bounding boxes in frame t.

Running-time is not reported since different methods run on different computing platforms. For inference speed of Transformer architecture, we refer to [5, 48].

4.2. Implementation details

We use ResNet-50 [12] as the network backbone. The optimizer is AdamW [24] with batch size 16, initial transformer’s learning rate 2×10^{-4} , the backbone’s 2×10^{-5} , and weight decay 10^{-4} . All transformer weights are initialized with Xavier init [11], and the backbone ImageNet-pretrained [8] model with frozen batch-norm layers [15]. We use data augmentation including random horizontal, random crop, scale augmentation, resizing the input images such that the shortest side is at least 480 and at most 800 pixels while the longest at most 1333. We train the networks for 150 epochs and the learning rate drops by a factor of 10 at the 100th epoch, unless otherwise noted.

Following [45], we pre-train our network on CrowdHuman [31]. More external data [43, 30] may further increase the performance but is not the focus of this work. The effect of external training data is shown in Table 1.

4.3. Ablation study

We first ablate the effect of Transformer architecture. The architectures differ mainly on the input feature. Results are shown in Table 2.

| Training data | MOTA↑ | FP↓ | FN↓ | IDs↓ |
|----------------------|-------|-------|-------|------|
| Only CrowdHuman [31] | 53.8 | 13.0% | 32.3% | 1.0% |
| Only MOT | 61.6 | 3.4% | 34.2% | 0.9% |
| CrowdHuman + MOT | 65.4 | 4.0% | 29.7% | 0.9% |

Table 1: **Ablation study on external training data.** 1st row is the model trained only on CrowdHuman dataset. 2nd row is the model trained only on split training set of MOT dataset. 3rd row is the model trained on CrowdHuman dataset first and then on split training set of MOT dataset. All models are tested on split validation set of MOT dataset.

| Method | MOTA↑ | FP↓ | FN↓ | IDs↓ |
|-----------------------------|-------|------|-------|------|
| Transformer [5] | 55.4 | 7.4% | 35.2% | 2.0% |
| Transformer-DC5 [5] | 59.0 | 5.2% | 34.0% | 1.8% |
| Transformer-P3 | 59.3 | 5.1% | 33.8% | 1.8% |
| Deformable Transformer [48] | 65.4 | 4.0% | 29.7% | 0.9% |

Table 2: **Ablation study on Transformer architecture.** Original transformer suffers from low feature resolution. Deformable DETR with multi-scale feature input achieves best performance.

Transformer. Following [5], original Transformer architecture is built on feature map of res5 stage [12]. This design requires more training epochs, and we train the networks for 500 epochs and the learning rate drop by a factor of 10 at the 400th epoch. However, the final performance is limited, only 55.4 MOTA. The main reason is the low feature resolution of res5 is negative to detection and tracking of small objects.

Transformer-DC5. To increase the feature resolution, we apply dilation convolution to res5 stage and remove a stride from the first convolution of this stage, called Transformer-DC5 [5]. This design yields an obvious 3.6 MOTA improvement. However, it also leads to the drawback of dilation convolution, such as big memory usage.

Transformer-P3. Feature pyramid network(FPN) [22] is a widely-used architecture for increasing feature resolution. Here we adopt P3 layer of FPN as the input feature map. Encoder is directly removed from the whole pipeline since the memory limitation. After removing encoder, the learning rate of the backbone could be raised to the same as transformers. The final performance is similar to Transformer-DC5. Although the feature resolution of Transformer-P3 is bigger than that of DC5, the absence of encoder blocks further improving performance.



Figure 4: **Visualization of TransTrack with different input query.** 1st row is **only learned object query**. 2nd row is **only object feature query from the previous frame**. 3rd row is **both learned object query and object feature query from the previous frame**. Only learned object query or object feature query from the previous frame causes ID switch case or missing object case. TransTrack takes both as input query and exhibits best detection and tracking performance.

Deformable Transformer. Deformable Transformer [48] is proposed to solve the issue of limited feature resolution in Transformer. Within plausible memory usage, it fuses multiple-scale feature into the whole encoder-decoder pipeline and achieves excellent performance in general object detection dataset. We introduce it into our method, and the performance is boosted up by a significant 5.6 MOTA, up to 64.9 MOTA. This is a very competitive performance among published methods. We use Deformable Transformer as our baseline setting.

Next, we ablate the effect of input query. Experimental and visual results are shown in Table 3 and Figure 4.

Only learned object query. When the input query is only learned object query, we introduce an extremely naive pipeline, in which each frame outputs detection boxes separately and detection results are associated according to their index in the output set. In fact, such a naive implementation can achieve not bad performance, both in detection metric

| Method | MOTA \uparrow | FP \downarrow | FN \downarrow | IDs \downarrow |
|-------------------------------------|-----------------|-----------------|-----------------|------------------|
| Only learned object query | 58.3 | 4.0% | 29.7% | 8.0% |
| Only F_{t-1} object feature query | - | 15.6% | 93.8% | 0.3% |
| Both | 65.4 | 4.0% | 29.7% | 0.9% |

Table 3: **Ablation study on input query.** Only learned object query obtains limited association performance. Only object feature query from the previous frame leads to numerous FN since it misses new-coming objects. Both achieve to best detection and tracking performance.

and association metric, shown in the first row of Table 3. This is because each object query predicts the object in the certain area, and most objects just move around a small distance in the video sequence. However, solely relying on the index in the output set leads to non-negligible wrong matching, especially when the object moves through a long distance. When the object moves around a wide range, it is correlated to different object query and its index in the

| Public Detection | | | | | | | | | |
|-------------------|-------------------------|-------------|-------------|-------------|--------------|--------------|--------------|---------------|-------------|
| Process | Method | MOTA↑ | IDF1↑ | MOTP↑ | MT↑ | ML↓ | FP↓ | FN↓ | IDS↓ |
| Offline | MHT-bLSTM [18] | 47.5 | 51.9 | 77.5 | 18.2% | 41.7% | 25981 | 268042 | 2069 |
| | EDMT [6] | 50.0 | 51.3 | 77.3 | 21.6% | 36.3% | 32279 | 247297 | 2264 |
| | JCC [17] | 51.2 | 54.5 | 75.9 | 20.9% | 37.0% | 25937 | 247822 | 1802 |
| | FWT [13] | 51.3 | 47.6 | 77.0 | 21.4% | 35.2% | 24101 | 247921 | 2648 |
| Online | DMAN [47] | 48.2 | 55.7 | 75.9 | 19.3% | 38.3% | 26218 | 263608 | 2194 |
| | MOTDT [7] | 50.9 | 52.7 | 76.6 | 17.5% | 35.7% | 24069 | 250768 | 2474 |
| | Tracktor[1] | 53.5 | 52.3 | 78.0 | 19.5% | 36.6% | 12201 | 248047 | 2072 |
| Private Detection | | | | | | | | | |
| Process | Method | MOTA↑ | IDF1↑ | MOTP↑ | MT↑ | ML↓ | FP↓ | FN↓ | IDS↓ |
| Online | Tracktor+CTdet [1] | 54.4 | 56.1 | 78.1 | 25.7% | 29.8% | 44109 | 210774 | 2574 |
| | DeepSORT [41] | 60.3 | 61.2 | 79.1 | 31.5% | 20.3% | 36111 | 185301 | 2442 |
| | TubeTK [26] | 63.0 | 58.6 | 78.3 | 31.2% | 19.9% | 27060 | 177483 | 4137 |
| | CenterTrack [45] | 67.8 | 64.7 | 78.4 | 34.6% | 24.6% | 18489 | 160332 | 3039 |
| | ChainedTracker [27] | 66.6 | 57.4 | 78.2 | 32.2% | 24.2% | 22284 | 160491 | 5529 |
| | TransTrack(ours) | 65.8 | 56.9 | 78.8 | 32.2% | 21.8% | 24000 | 163683 | 5355 |

Table 4: **Evaluation on MOT17 test sets.** We list published results of both public and private detection and compare TransTrack with methods of private detection. TransTrack performs excellent performance in terms of MOTP and FN, proving the success of introducing learned object query into the pipeline. IDs of TransTrack is comparable with ChainedTracker, which shows the effectiveness of object feature query in associating two adjacent frames.

output set will change. A visualization case is shown in the first row of Figure 4.

Only object feature query from the previous frame. When the input query is only object feature of the previous frame, each object feature predicts its own position in the current frame, in line with the query-key mechanism. This happens to be simultaneous detection and association. The visualization in the second row of Figure 4 shows that this method is capable to associate the object with a large range of motion. Nevertheless, the disadvantage of this implementation is obvious. Only the object that appears in the first frame can be tracked successively. For the whole video sequence, most of the objects will be missed and FN metric falls off, shown in the second row of Table 3.

Both. From above ablation study, we conclude that a desirable tracking model requires both learned object query and object feature query from the previous frame. This is the baseline setting of TransTrack. TransTrack outputs complete and ordered object sets. Both quantitative and visualization results are best among other settings.

4.4. Benchmark evaluation

We compare TransTrack with other methods on MOT17 test dataset in Table 4. Running-time is not reported here since different methods run on different computing platforms. TransTrack is designed as joint-detection-and-tracking method, so the “private detector” protocol is adopted.

TransTrack achieves comparable results with the current

state-of-the-art methods, especially in terms of MOTP and FN. The excellent MOTP demonstrates TransTrack can precisely locate objects in the image. The good FN represents that most objects are successfully detected. Those prove the success of introducing learned object query into query-key pipeline.

As for ID-switch, TransTrack is comparable with a state-of-the-art model, ChainedTracker [27], which proves the effectiveness of object feature query to associate adjacent frames. Although ID-switch is somewhat inferior to other methods, we explain that as the first work to introduce query-key mechanism into MOT, we do not apply complex operations in order to keep the simplicity and originality of tracking methods based on query-key mechanism. We believe that it is a promising direction to further improve the overall performance of TransTrack.

5. Conclusion

We set up a simple joint-detection-and-tracking MOT pipeline, TransTrack, based on query-key mechanism. The image feature maps are common keys among queries. The learned object query detects objects in the current frame and object feature query from the previous frame associates objects in the current frame with the previous ones. Our method achieves competitive 65.8% MOTA on the MOT17 challenge dataset. Query-key mechanism is widely used in the field of SOT, but is seldom studied in MOT. For the first time, we demonstrate query-key mechanism could serve as an effective and strong baseline for MOT.

References

- [1] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *ICCV*, pages 941–951, 2019. 4, 5, 8
- [2] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: the clear mot metrics. *EURASIP Journal on Image and Video Processing*, 2008:1–10, 2008. 6
- [3] Luca Bertinetto, Jack Valmadre, João F. Henriques, Andrea Vedaldi, and Philip H. S. Torr. Fully-convolutional siamese networks for object tracking, 2016. 2, 3
- [4] Alex Bewley, Zongyuan Ge, Lionel Ott, Fabio Ramos, and Ben Upcroft. Simple online and realtime tracking. In *2016 IEEE International Conference on Image Processing (ICIP)*, pages 3464–3468, 2016. 3
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-End object detection with transformers. In *ECCV*, 2020. 3, 5, 6
- [6] Jiahui Chen, Hao Sheng, Yang Zhang, and Zhang Xiong. Enhancing detection model for multiple hypothesis tracking. In *PCVPRW*, pages 18–27, 2017. 8
- [7] Long Chen, Haizhou Ai, Zijie Zhuang, and Chong Shang. Real-time multiple people tracking with deeply learned candidate selection and person re-identification. In *ICME*, 2018. 3, 8
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. 6
- [9] Christoph Feichtenhofer, Axel Pinz, and Andrew Zisserman. Detect to track and track to detect, 2018. 4
- [10] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2019. 5
- [11] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010. 6
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 6
- [13] Roberto Henschel, Laura Leal-Taixé, Daniel Cremers, and Bodo Rosenhahn. Fusion of head and full-body detectors for multi-object tracking. In *CVPRW*, pages 1428–1437, 2018. 8
- [14] Han Hu, Jiayuan Gu, Zheng Zhang, Jifeng Dai, and Yichen Wei. Relation networks for object detection. In *CVPR*, 2018. 3
- [15] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015. 6
- [16] Kai Kang, Hongsheng Li, Junjie Yan, Xingyu Zeng, Bin Yang, Tong Xiao, Cong Zhang, Zhe Wang, Ruohui Wang, Xiaogang Wang, and et al. T-cnn: Tubelets with convolutional neural networks for object detection from videos. *IEEE Transactions on Circuits and Systems for Video Technology*, 28(10):2896–2907, Oct 2018. 4
- [17] Margret Keuper, Siyu Tang, Bjoern Andres, Thomas Brox, and Bernt Schiele. Motion segmentation & multiple object tracking by correlation co-clustering. *TPAMI*, 42(1):140–153, 2018. 8
- [18] Chanho Kim, Fuxin Li, and James M Rehg. Multi-object tracking with neural gating using bilinear lstm. In *ECCV*, pages 200–215, 2018. 8
- [19] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 3, 5
- [20] Bo Li, Wei Wu, Qiang Wang, Fangyi Zhang, Junliang Xing, and Junjie Yan. Siamrpn++: Evolution of siamese visual tracking with very deep networks, 2018. 2, 3
- [21] B. Li, J. Yan, W. Wu, Z. Zhu, and X. Hu. High performance visual tracking with siamese region proposal network. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8971–8980, 2018. 2, 3
- [22] Tsung-Yi Lin, Piotr Dollar, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017. 3, 6
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection, 2018. 3, 5
- [24] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. 6
- [25] Anton Milan, Laura Leal-Taixé, Ian Reid, Stefan Roth, and Konrad Schindler. Mot16: A benchmark for multi-object tracking. *arXiv preprint arXiv:1603.00831*, 2016. 2, 6
- [26] Bo Pang, Yizhuo Li, Yifan Zhang, Muchen Li, and Cewu Lu. Tubekt: Adopting tubes to track multi-object in a one-step training model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6308–6318, 2020. 8
- [27] Jinlong Peng, Changan Wang, Fangbin Wan, Yang Wu, Yabiao Wang, Ying Tai, Chengjie Wang, Jilin Li, Feiyue Huang, and Yanwei Fu. Chained-tracker: Chaining paired attentive regression results for end-to-end joint multiple-object detection and tracking. *arXiv preprint arXiv:2007.14557*, 2020. 4, 5, 8
- [28] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks, 2016. 3, 4
- [29] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *CVPR*, 2019. 5
- [30] Chaobing Shan, Chunbo Wei, Bing Deng, Jianqiang Huang, Xian-Sheng Hua, Xiaoliang Cheng, and Kewei Liang. Fgagt: Flow-guided adaptive graph tracking. *arXiv preprint arXiv:2010.09015*, 2020. 6
- [31] Shuai Shao, Zijian Zhao, Boxun Li, Tete Xiao, Gang Yu, Xiangyu Zhang, and Jian Sun. Crowdhuman: A benchmark for detecting human in a crowd. *arXiv preprint arXiv:1805.00123*, 2018. 6

- [32] Peize Sun, Yi Jiang, Enze Xie, Zehuan Yuan, Changhu Wang, and Ping Luo. Onenet: Towards end-to-end one-stage object detection. *arXiv preprint arXiv:2012.05780*, 2020. 5
- [33] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, and Ping Luo. Sparse r-cnn: End-to-end object detection with learnable proposals. *arXiv preprint arXiv:2011.12450*, 2020. 5
- [34] Siyu Tang, Mykhaylo Andriluka, Bjoern Andres, and Bernt Schiele. Multiple people tracking by lifted multicut and person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. 2, 3
- [35] Ran Tao, Efstratios Gavves, and Arnold W. M. Smeulders. Siamese instance search for tracking, 2016. 2, 3
- [36] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008, 2017. 2, 3, 4, 5
- [37] Jianfeng Wang, Lin Song, Zeming Li, Hongbin Sun, Jian Sun, and Nanning Zheng. End-to-end object detection with fully convolutional network. *arXiv preprint arXiv:2012.03544*, 2020. 5
- [38] Qiang Wang, Li Zhang, Luca Bertinetto, Weiming Hu, and Philip H. S. Torr. Fast online object tracking and segmentation: A unifying approach, 2019. 3
- [39] Zhongdao Wang, Liang Zheng, Yixuan Liu, and Shengjin Wang. Towards real-time multi-object tracking. *arXiv preprint arXiv:1909.12605*, 2019. 2, 3, 4
- [40] Greg Welch, Gary Bishop, et al. An introduction to the kalman filter, 1995. 3
- [41] Nicolai Wojke, Alex Bewley, and Dietrich Paulus. Simple online and realtime tracking with a deep association metric. In *ICIP*, pages 3645–3649. IEEE, 2017. 3, 8
- [42] Fengwei Yu, Wenbo Li, Quanquan Li, Yu Liu, Xiaohua Shi, and Junjie Yan. Poi: Multiple object tracking with high performance detection and appearance feature. In *European Conference on Computer Vision*, pages 36–42. Springer, 2016. 2, 3
- [43] Yifu Zhan, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. A simple baseline for multi-object tracking. *arXiv preprint arXiv:2004.01888*, 2020. 3, 4, 6
- [44] Zheng Zhang, Dazhi Cheng, Xizhou Zhu, Stephen Lin, and Jifeng Dai. Integrated object detection and tracking with tracklet-conditioned detection, 2018. 4
- [45] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Tracking objects as points, 2020. 4, 5, 6, 8
- [46] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points, 2019. 4
- [47] Ji Zhu, Hua Yang, Nian Liu, Minyoung Kim, Wenjun Zhang, and Ming-Hsuan Yang. Online multi-object tracking with dual matching attention networks. In *ECCV*, pages 366–382, 2018. 3, 8
- [48] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3, 5, 6, 7
- [49] Xizhou Zhu, Yujie Wang, Jifeng Dai, Lu Yuan, and Yichen Wei. Flow-guided feature aggregation for video object detection. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, Oct 2017. 4
- [50] Zheng Zhu, Qiang Wang, Bo Li, Wei Wu, Junjie Yan, and Weiming Hu. Distractor-aware siamese networks for visual object tracking, 2018. 3