

# ReST: A Reconfigurable Spatial-Temporal Graph Model for Multi-Camera Multi-Object Tracking

Cheng-Che Cheng<sup>1</sup> Min-Xuan Qiu<sup>1</sup>

<sup>1</sup>National Tsing Hua University, Taiwan

{chengche6230, maisiechiu}@gapp.nthu.edu.tw

Chen-Kuo Chiang<sup>2</sup> Shang-Hong Lai<sup>1</sup>

<sup>2</sup>National Chung Cheng University, Taiwan

ckchiang@cs.ccu.edu.tw lai@cs.nthu.edu.tw

## Abstract

*Multi-Camera Multi-Object Tracking (MC-MOT) utilizes information from multiple views to better handle problems with occlusion and crowded scenes. Recently, the use of graph-based approaches to solve tracking problems has become very popular. However, many current graph-based methods do not effectively utilize information regarding spatial and temporal consistency. Instead, they rely on single-camera trackers as input, which are prone to fragmentation and ID switch errors. In this paper, we propose a novel reconfigurable graph model that first associates all detected objects across cameras spatially before reconfiguring it into a temporal graph for Temporal Association. This two-stage association approach enables us to extract robust spatial and temporal-aware features and address the problem with fragmented tracklets. Furthermore, our model is designed for online tracking, making it suitable for real-world applications. Experimental results show that the proposed graph model is able to extract more discriminating features for object tracking, and our model achieves state-of-the-art performance on several public datasets. Code is available at <https://github.com/chengche6230/ReST>.*

## 1. Introduction

Multi-Object Tracking (MOT) is an important task in computer vision, which involves object detection and tracking multiple objects over time in an image sequence. It can be applied to several real-world scenarios, such as video surveillance, autonomous vehicles, and sports analysis. Despite numerous research methods proposed for MOT, the problem of fragmented tracklets or ID switching caused by frequent occlusion in crowded scenes remains a major challenge. One potential solution is to track objects under a multi-camera setting, which is called a Multi-Camera Multi-Object Tracking (MC-MOT) task. By leveraging information from multiple cameras, occluded objects in one view may become clearly visible in another view, allowing

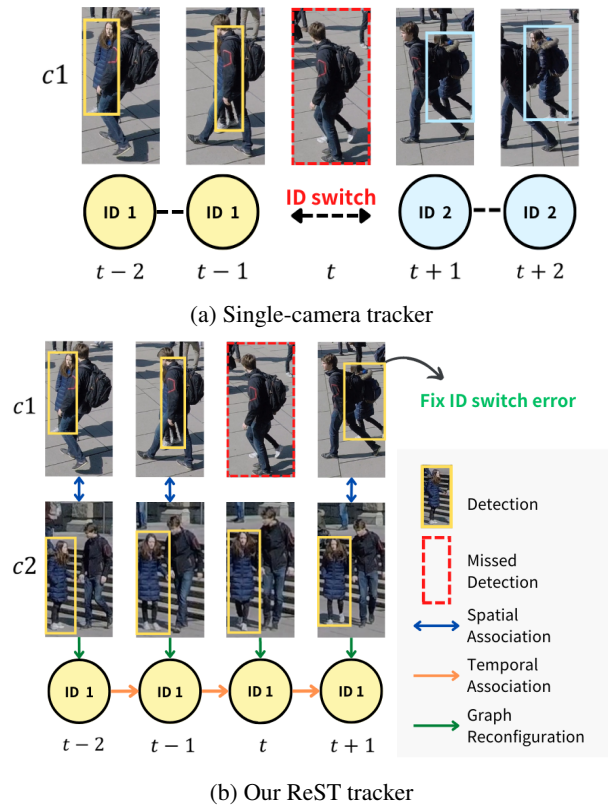


Figure 1: Example of handling object occlusion for MC-MOT. (a) When occlusion occurs at time  $t$  (red dotted box), the single-camera tracker generates fragmented tracklets and causes ID switch errors. (b) Our ReST tracker corrects object ID in  $c1$  via Spatial and Temporal Association, by leveraging spatial and temporal consistency.

for more accurate object tracking results.

Most of tracking-by-detection paradigms [1] adopt Kalman filter [17] in the data association stage. It serves as a motion model, predicting the next possible position and matching with previous detection. However, such method is usually deterministic and cannot adapt to the dynamically changing environment. In addition, the tracking results are difficult to achieve globally optimal, since the illumi-

nation, relative geometry distance, or sampling rate varies from dataset to dataset, which is common in real-world scenarios. Accordingly, there is another fashion reformulating the association problem into link prediction on graph [5, 18, 25, 27]. It allows a trainable model to determine how strong the connection is between two detections. Thus, objects can be dynamically associated depending on environmental conditions.

However, there still remains some issues in current graph-based models for MC-MOT. First of all, many approaches rely on single-camera tracker to generate the initial tracklets [13, 25, 27, 37]. Although many methods have been proposed to refine tracklets, tracking errors in single-view are often left unaddressed. Additionally, these methods do not fully leverage the rich spatial and temporal information that is crucial for MC-MOT task. Recently, spatial-temporal models have been employed to learn representative features for tracklets. However, the resulting graphs are usually complex and hard to optimize.

In this paper, we propose a novel Reconfigurable Spatial-Temporal graph model (ReST) for MC-MOT to overcome the problems mentioned above. The MC-MOT problem is re-formulated as two sub-tasks, Spatial Association and Temporal Association, in our approach. In Spatial Association, it focuses on matching objects across different views. Temporal Association exploits temporal information, such as speed and time, to build temporal graph which associates objects across frames. By splitting the problem into two sub-tasks, spatial and temporal consistency can be individually optimized to achieve better tracking results. In addition, the graph model becomes smaller and easy to optimize. To bridge two association stages, Graph Reconfiguration module is proposed to aggregate information from spatial and temporal graph models. The merits of involving graph reconfiguration are two-fold. Firstly, when the nodes of the same object are merged, the reconfigured graph becomes very compact. Secondly, the refinement of the graph model can be iteratively performed in each reconfiguration step during inference, leading to more representative feature extraction and better tracking results. As depicted in Figure 1a, when the girl is occluded, fragmented tracklets are produced, causing the ID switch problem. In Figure 1b, correct object ID can be retained by employing spatial and temporal consistency via Spatial Association, Temporal Association, and Graph Reconfiguration modules.

The proposed graph model is called reconfigurable because the vertex set and edge set of spatial and temporal graphs are reconfigured to construct a new graph at each time. Thus, it tends to adapt to dynamic scenes. Unlike existing methods, our model does not rely on the results from single-camera tracker. The tracking and association of the detected objects is accomplished through iteratively constructing spatial and temporal graphs. Our model is de-

signed for online object tracking since it does not use or rely on any information from future frames.

**Contributions** Our contributions can be summarized as follows. 1) The Multi-Camera Multi-Object Tracking problem is formulated as two sub-tasks in the proposed graph model, Spatial Association and Temporal Association. This enables the employment of spatial and temporal consistency and better model optimization. 2) Graph Reconfiguration module is proposed to leverage tracking results from two stages. This makes the object tracking apt to dynamic scene changes and online tracking scenarios. 3) Experimental results demonstrate that our model achieves state-of-the-art performance on Wildtrack and competitive results on other benchmark datasets.

## 2. Related Work

In recent years, a number of research works have focused on single-camera MOT. For example, [6, 9, 19, 41] focus on improving data association and precisely extracting motion. [3, 8, 29, 35, 39, 43] unify the object detection and association stage into an end-to-end model. Recently, MC-MOT has received significant attention and grown increasingly [7, 13, 14, 25, 26, 27, 37, 38, 40]. Although it contains more spatial-temporal information than single-camera tracking problem, the MC-MOT problem still presents several challenges that must be overcome, including varying environmental conditions and the lack of integration of spatial-temporal information.

**Spatial-Temporal Representation Learning** The spatial and temporal feature is a key factor in motion-related areas, such as human pose estimation, and MOT. [40] sets up an occupancy map to fuse cross-view spatial correlations, followed by a Deep Glimmse Network to capture temporal information. [37] formulates MC-MOT as a compositional structure optimization problem, associating tracklets by appearance, geometry, and motion consistency. Starting with a scene node, a Spatial-Temporal Attributed Parse Graph is then constructed in [38]. The scene nodes are then decomposed into several tracklet nodes, containing different types of semantic attributes, such as appearance and action.

**Graph-Based Methods** Graph Neural Networks (GNN) [12] and Graph Convolutional Networks (GCN) [21] have also been extensively studied for MOT [5, 9, 18, 25, 27] due to the flexibility of dynamical affinity association training. In [18], it performs a standard graph model, passing the messages using localized polar feature representation. Nonetheless, it makes the association difficult when the spatial and temporal features are ambiguous and edges are all mixed up in one graph. [9] first constructs a candidate graph, followed by Transformer [34] which serves as a feature encoder. However, they do not utilize spatial-temporal consistency in the graph. Combined with attention

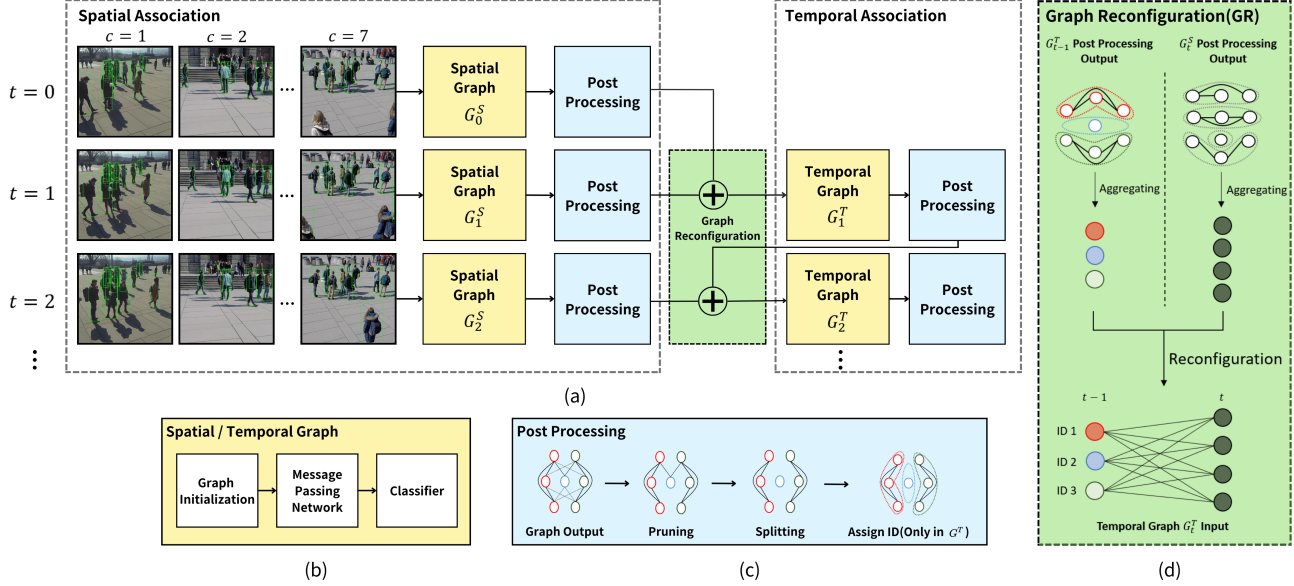


Figure 2: Proposed ReST framework. (a) Inference architecture: given input detection from all views at time  $t$ , Spatial Association is performed and followed by Graph Reconfiguration and Temporal Association. (b) Graph model learning: both  $G^S$  and  $G^T$  are trained following the steps: initialization, message passing network, and edge classification. (c) Post-processing module: output a refined graph by two schemes: pruning and splitting. (d) Graph Reconfiguration: aggregate graph nodes from  $G_{t-1}^T$  and  $G_t^S$  and reconfigure a new temporal graph  $G_t^T$ .

mechanism, the structural and temporal attention layers in [27] enable robust feature extraction for link prediction. It relies on single-camera tracker as input, leading to a sub-optimal solution in multi-views setting. The computation cost of dual attention layers is expensive. Graph Reconfiguration was not performed in [27]; instead, they aggregate objects in consecutive frames by directly adding new nodes and edges. In [25], multicut [30, 31, 32] is applied to MC-MOT to obtain the globally optimal solution. Likewise, their model relies on single-camera tracker to generate initial tracklets. Compared to our stage-wise optimization, [25] proposed a joint spatial-temporal optimization model. Mixing all spatial and temporal edges into one graph cannot fully leverage spatial and temporal consistency individually.

### 3. Proposed Method

Our method follows the tracking-by-detection paradigm [1]. In contrast to prior methods [13, 25, 27, 37], our model does not rely on the off-the-shelf single-camera trackers trained by massive single-camera MOT datasets. Instead, the tracking and association problem is formulated into a link prediction problem on graph. The proposed ReST framework divides the MC-MOT process into two sub-tasks, Spatial Association and Temporal Association. Spatial Association focuses on matching objects across different views. The spatial graph concentrates on building spatial correlation between nodes. Temporal Association exploits temporal information, such as speed and time. The temporal graph extracts temporal features to associate ob-

jects across frames. These two modules take turns constructing spatial and temporal graphs for each frame. In addition, a novel Graph Reconfiguration module is proposed to reconfigure current spatial graph and temporal graph from previous timestamp as a new temporal graph. By this setting, the proposed graph model associates all detected objects frame by frame via extracting robust spatial and temporal-aware features. For the inference of each graph, we perform a post-processing module, which contains pruning and splitting, to fix association errors in time in each iteration. Compared with one single spatial-temporal graph model, two expert graph models can be trained separately, making it focus on extracting specific features and reduce the ambiguity of spatial-temporal correlations. The system framework of ReST is depicted in Figure 2.

#### 3.1. Problem Formulation

In MC-MOT, the goal is to track multiple objects across frames and views. Assume there are  $C$  synchronous and static cameras that have overlapping fields of view (FoV). Define a graph model  $G_t = (V_t, E_t)$  at time  $t$ , where  $V_t$  is the vertex set and  $E_t$  is the edge set. Each node  $v_i \in V_t$  represents one input detection and it contains the following information: camera ID  $c_{v_i} \in \mathbb{R}^1$ , timestamp  $t_{v_i} \in \mathbb{R}^1$ , object ID  $o_{v_i} \in \mathbb{R}^1$ , i.e. ground truth label, bounding-box position  $b_{v_i} \in \mathbb{R}^4$ , appearance feature  $d_{v_i} \in \mathbb{R}^{512}$ , geometry position  $p_{v_i} \in \mathbb{R}^2$ , and speed information  $s_{v_i} \in \mathbb{R}^2$ . Let  $I_{v_i}$  be the cropped image of  $b_{v_i}$ . The appearance feature  $d_{v_i}$

can be obtained by

$$d_{v_i} = f_{ReID}(I_{v_i}), \quad (1)$$

where  $f_{ReID}$  is an off-the-shelf Re-Identification (ReID) model. Denote  $P_c$  as the projection function of camera  $c$ . It projects the foot point of the bounding-box from its camera view to a common ground plane (Section A Appendix). The geometry position  $p_{v_i}$  and the speed information  $s_{v_i}$  for the reference node  $v_j$  can be calculated by

$$p_{v_i} = P_{c_{v_i}}(x + \frac{w}{2}, y + h), \quad s_{v_i} = \frac{p_{v_i} - p_{v_j}}{t_{v_i} - t_{v_j}}, \quad (2)$$

where  $(x, y, w, h)$  represents the bounding box  $b_{v_i}$ , and  $t_{v_i} > t_{v_j} \geq 0$ .

The relative distance of geometry position, appearance feature, and speed between any pair of nodes  $v_i$  and  $v_j$  can be defined by

$$\begin{aligned} \Delta d_{ij} &= [\|d_{v_i} - d_{v_j}\|_1, 1 - \text{cosine\_similarity}(d_{v_i}, d_{v_j})], \\ \Delta p_{ij} &= [\|p_{v_i} - p_{v_j}\|_1, \|p_{v_i} - p_{v_j}\|_2], \\ \Delta s_{ij} &= [\|s_{v_i} - s_{v_j}\|_1, \|s_{v_i} - s_{v_j}\|_2]. \end{aligned} \quad (3)$$

Following [24],  $\Delta p_{ij}$ ,  $\Delta d_{ij}$ , and  $\Delta s_{ij}$  are used as the initial edge features.

### 3.2. Reconfigurable Spatial-Temporal Graph

To better associate objects across views and frames, a novel reconfigurable graph framework is proposed for the inference stage. Our model follows the pipeline to achieve object tracking and association: perform Temporal Association to construct temporal graph at time  $t-1$ , Spatial Association to build spatial graph at time  $t$  and then Graph Reconfiguration is applied to reconfigure a new temporal graph at time  $t$ .

**Spatial Association** Objects from different views at the current frame are first associated. In this stage, only spatial information is required to construct the spatial graph. The spatial graph is denoted as  $G^S$ . It concentrates on extracting spatial features for cross-view association. We denote all detected objects from camera  $c$  at time  $t$  as  $B_c^t$ . Given detected objects from all cameras at time  $t$ , the vertex set  $V_t^S$  of spatial graph  $G_t^S$  can be defined by

$$V_t^S = \bigcup_{i=1}^C B_i^t. \quad (4)$$

The vertex set  $V_t^S$  is composed of all detected objects across cameras at time  $t$ . We denote the adjacency matrix of  $G_t^S$  as  $A_t^S = [a_{ij}^S]$ . The initial edge construction is defined by

$$a_{ij}^S = \begin{cases} 1, & \text{if } c_{v_i} \neq c_{v_j} \\ 0, & \text{otherwise} \end{cases}. \quad (5)$$

That is, there is an edge between node  $v_i$  and  $v_j$  if both nodes are from different cameras. Once the spatial graph is constructed, the initial node feature  $h_{v_i}^0$  and initial edge feature  $h_{e_{ij}}^0$  can be defined by

$$h_{v_i}^0 = f_{FE}^v(d_{v_i}),$$

$$h_{e_{ij}}^0 = f_{FE}^e([\Delta p_{ij}, \Delta d_{ij}]), \quad (6)$$

where  $f_{FE}^v$  is a node feature encoder and  $f_{FE}^e$  is an edge feature encoder. Both are implemented by Multi-Layer Perceptron (MLP). The operator  $[\cdot, \cdot]$  denotes the concatenation of two terms. Thus,  $h_{v_i}^0$  and  $h_{e_{ij}}^0$  are used as input of Message Passing Network (MPN) to extract edge features. When final edge features are extracted, link prediction is performed to construct the final spatial graph as the result of object association. Followed by post-processing module, the spatial graph can be further refined. Line 3 to line 9 in Algorithm 1 present the steps of spatial association. The details of MPN and link prediction will be described in subsection 3.4.1 and 3.4.2.

**Temporal Association** In this stage, objects are associated from different frames by time, without using any camera information. In other words, the temporal graph is view-invariant and cares more about temporal correlation. Given temporal graph  $G_t^T$  at time  $t$ , the initial node and edge features can be computed by

$$h_{v_i}^0 = f_{FE}^v([d_{v_i}, p_{v_i}]),$$

$$h_{e_{ij}}^0 = f_{FE}^e([\Delta p_{ij}, \Delta d_{ij}, \Delta s_{ij}]). \quad (7)$$

Note that we append an extra speed term in the input of edge feature to capture relative motion and direction between nodes. Following similar steps in  $G_t^S$ , MPN, link prediction, and post-processing module are performed, as indicated from line 12 to 17 in Algorithm 1. For temporal graph  $G_t^T$ , tracklet ID is assigned to each node within the same connected component as the tracking results at time  $t$ .

**Graph Reconfiguration** After the association stage, several connected components representing the same object can be obtained. In order to bridge two association stages, our model reconfigures two graphs into a new Temporal Graph.  $G_t^T$  is denoted as the temporal graph at time  $t$ .  $G_t^T$  is reconfigured from spatial graph  $G_t^S$  and temporal graph  $G_{t-1}^T$ . From  $G_t^S$  and  $G_{t-1}^T$ , all nodes in the same connected components are aggregated into one node. After the aggregation, a new vertex set is formed for  $G_t^T$ . Let  $H(G) = \{H_1, \dots, H_n\}$  be the set of connected components in  $G$ . All node information within one connected component is averaged and serves as the initial features, i.e.

$$d_v = \frac{\sum_{v \in H_i} d_v}{|H_i|},$$

$$p_v = \frac{\sum_{v \in H_i} p_v}{|H_i|}. \quad (8)$$

where  $H_i \in H(G_t^S) \cup H(G_{t-1}^T)$ . Once the vertex set of  $G_t^T$  is determined, the edge of  $G_t^T$  can be defined by

$$a_{ij}^T = \begin{cases} 1, & \text{if } t_{v_i} \neq t_{v_j}, \\ 0, & \text{otherwise} \end{cases}, \quad (9)$$

where  $a_{ij}^T \in A_t^T$  is adjacency matrix of  $G_t^T$ . The edges exist in  $G_t^T$  if two nodes are from different time frames. The complete steps in the inference are given in Algorithm 1.

---

#### Algorithm 1 Inference Algorithm

---

**Input:** temporal graph  $G_{t-1}^T$ , detection set  $B_c^t$  from all  $C$  views at time  $t$ .

**Output:** temporal graph  $G_t^T$ , tracking result at time  $t$ .

- 1: construct spatial graph  $G_t^S$  through  $B_c^t$
  - 2: compute initial feature  $h_{v_i}^0, h_{e_{ij}}^0$  for  $G_t^S$
  - 3: **for**  $l = 1$  to  $L$  **do**
  - 4:    $h_{v_i}^l, h_{e_{ij}}^l = \text{MPN}(G_t^S, h_{v_i}^{l-1}, h_{e_{ij}}^{l-1})$
  - 5: **end for**
  - 6:  $\hat{y}_{e_{ij}}^L = f_{CLS}(h_{e_{ij}}^L)$
  - 7:  $G_t^S = \text{post-processing}(G_t^S, \hat{y}_{e_{ij}}^L)$
  - 8: **if**  $t > 0$  **then**
  - 9:    $G_t^T = \text{reconfiguration}(G_t^S, G_{t-1}^T)$
  - 10:   compute initial feature  $h_{v_i}^0, h_{e_{ij}}^0$  for  $G_t^T$
  - 11:   **for**  $l = 1$  to  $L$  **do**
  - 12:      $h_{v_i}^l, h_{e_{ij}}^l = \text{MPN}(G_t^T, h_{v_i}^{l-1}, h_{e_{ij}}^{l-1})$
  - 13:   **end for**
  - 14:    $\hat{y}_{e_{ij}}^L = f_{CLS}(h_{e_{ij}}^L)$
  - 15:    $G_t^T = \text{post-processing}(G_t^T, \hat{y}_{e_{ij}}^L)$
  - 16: **end if**
- 

---

#### Algorithm 2 Post-Processing Algorithm

---

**Input:**  $G_t$ : graph at time  $t$ ,  $\hat{y}_{e_{ij}}^L$ : confidence score.

**Output:**  $G_t$ : refined graph at time  $t$ .

- 1:  $G_t = \text{pruning}(G_t, \hat{y}_{e_{ij}}^L)$
  - 2:  $G_t = \text{splitting}(G_t)$
  - 3: **if**  $G_t$  is temporal graph **then**
  - 4:   assigning tracklet ID
  - 5: **end if**
- 

### 3.3. Post-Processing

In post-processing, the objective is to refine the graph output. Since the vertices in the same connected component represent objects of the same ID, the connected component may contain vertices with different object IDs. Additional constraints can be included to reduce incorrect ID assignments. The post-processing is divided into three steps: pruning, splitting, and assigning object ID.

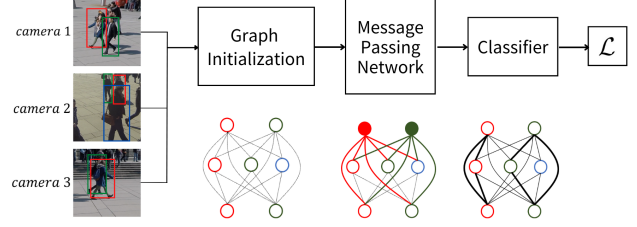


Figure 3: Graph model training. Given different input detection and constraints, we construct and train  $G^S$  and  $G^T$ , respectively. This figure shows an example of input detection for spatial graph.

**Pruning** Confidence score of each edge predicted by our model is used to prune the graph. If the score is greater than a given threshold  $\varepsilon$ , the edge is kept. Otherwise, it is removed. After pruning, edges with weak confidence are removed to improve the correctness of connected components.

**Splitting** Similar to [24], a few physical constraints and assumptions can be employed to further optimize each connected component. In spatial graph, assume that an object can only appear in each camera once, there are at most  $C$  nodes in each connected component and each node can be connected to at most  $C - 1$  nodes. Therefore, the constraints of spatial graph can be defined as

$$|V(H_i^S)| \leq C, \quad (10)$$

where  $H_i^S \in H(G^S)$ . For every node  $v \in H_i^S$ , we have

$$\text{degree}(v) \leq C - 1. \quad (11)$$

Each node in temporal graph can be connected to at most  $M - 1$  nodes, where  $M$  is the temporal window size.

$$\forall v \in H_i^T \in H(G^T), \text{degree}(v) \leq M - 1. \quad (12)$$

For any connected component violating Eq.(10)-(12), the edge with the lowest confidence score is removed. The operation is performed recursively until all constraints are satisfied.

**Assigning tracklet ID** After the post-processing for temporal graph is finished, tracklet IDs are assigned to nodes in the current frame. In practice, one node inherits the ID from nodes that are already in the same connected component, otherwise it is assigned a new ID (Section B Appendix). The post-processing steps are given in Algorithm 2.

### 3.4. Model Training

In this subsection, the details of Message Passing Network, link prediction, and training scheme are described, as depicted in Figure 3.

### 3.4.1 Message Passing Network

Given the initial feature  $h_{v_i}^0$  and  $h_{e_{ij}}^0$  of a graph, we follow the standard framework of MPN [5, 18, 24] and perform a fixed number of graph updates to obtain enhanced feature representation. Specifically, there are two steps in a graph update iteration; namely, Edge Update and Node Update.

**Edge Update** At each message passing iteration  $l = 1, \dots, L$ , edge feature is firstly updated by aggregating its source node feature and destination node feature as:

$$h_{e_{ij}}^l = f_{ME}^e([h_{v_i}^{l-1}, h_{e_{ij}}^{l-1}, h_{v_j}^{l-1}]), \quad (13)$$

where  $f_{ME}^e$  is an edge message encoder. MLP is exploited to encode the original message into a high-dimensional feature space.

**Node Update** After updating the edge feature, we then update each node by the messages sent from its neighbor nodes as:

$$h_{v_i}^l = \sum_{j \in N(v_i)} m_{ij}^l, \quad (14)$$

where  $N(v_i)$  denotes the neighbor nodes of  $v_i$ , and the message term can be computed by

$$m_{ij}^l = f_{ME}^v([h_{v_j}^{l-1}, h_{e_{ij}}^l]), \quad (15)$$

where  $f_{ME}^v$  is a node message encoder similar to  $f_{ME}^e$ .

### 3.4.2 Link Prediction

After MPN, enhanced edge features can be obtained for link prediction. It aims to decide whether an edge should be kept or removed in a graph. Specifically, a binary classifier is cascaded to MPN. Given the edge feature from iteration  $l$ , the classifier outputs a confidence score:

$$\hat{y}_{e_{ij}}^l = f_{CLS}(h_{e_{ij}}^l), \quad (16)$$

where  $f_{CLS}$  is a binary classifier implemented by MLP followed by a softmax layer. In the inference stage, the confidence score  $\hat{y}_{e_{ij}}^L$  at the last iteration is used for pruning.

### 3.4.3 Training Scheme

In model training, the spatial graph and temporal graph are trained independently to learn spatial and temporal-aware feature representation. To train the spatial graph, training input is all detections from different views at the same time frame. For temporal graph, training input only contains detections from different frames of the same camera  $c$  as

$$V_t^T = \bigcup_{i=0}^{M-1} B_c^{t-i}. \quad (17)$$

For both graphs, the ReID model  $f_{ReID}$  is frozen during the training process.  $f_{FE}^v, f_{FE}^e, f_{ME}^v, f_{ME}^e$ , and  $f_{CLS}$

are trainable MLPs (Section C Appendix). Focal Loss [23] is exploited to calculate the loss between ground-truth label and predicted label at each message passing iteration  $l$ , given by

$$\mathcal{L} = \sum_{l=1}^L \sum_{e_{ij} \in E^S \cup E^T} FL(\hat{y}_{e_{ij}}^l, y_{e_{ij}}), \quad (18)$$

where  $y_{e_{ij}}$  is ground truth label and its value equals 1 if  $v_i$  and  $v_j$  have the same object ID, i.e.  $o_{v_i} = o_{v_j}$ . Otherwise, it is 0.

In this way, our graph model can effectively learn how to associate two nodes spatially and temporally. Compared with other single spatial-temporal graph methods [5, 18, 25, 27], our graph can focus on learning more discriminating spatial and temporal features to cope with challenging multi-object tracking scenarios.

## 4. Experimental Results

In this section, we demonstrate our model performance on several benchmark datasets. Detailed implementation settings and ablation studies are presented. For evaluation of MC-MOT methods, ID score [28], i.e. IDF1, and the standard CLEAR MOT metrics [4], including MOTA, MOTP, Mostly Tracked (MT), and Mostly Lost (ML), are employed for a fair comparison. Experimental comparisons with the state-of-the-art MC-MOT methods are also presented in this section.

### 4.1. Datasets

Our experiments are conducted on three multi-view multi-object tracking datasets under diverse environmental conditions, such as illumination, density, and detection quality. All video sequences have synchronous and calibrated cameras with a certain ratio of overlapping FoV.

**Wildtrack [7]** It is considered the most challenging dataset with 7 cameras, having the most occlusion problems and the highest density. Specifically, there are about 25 people standing and walking around at each frame on average. We follow the common setting as [7, 10, 15, 16], which is trained on the first 360 frames and tested on the last 40 frames.

**CAMPUS [37]** All of the sequences in CAMPUS are reported in our results. People are doing all kinds of sports in *Garden 1*, which means the capability to capture diverse motion is crucial. *Garden 2* is a relatively sparse sequence. People are often occluded by cars in the *Parkinglot* sequence, making it hard to recover from different views. *Auditorium* is a sequence recorded in two scenes that we use to validate our model's ability to apply in a non-overlapping FoV scenario.



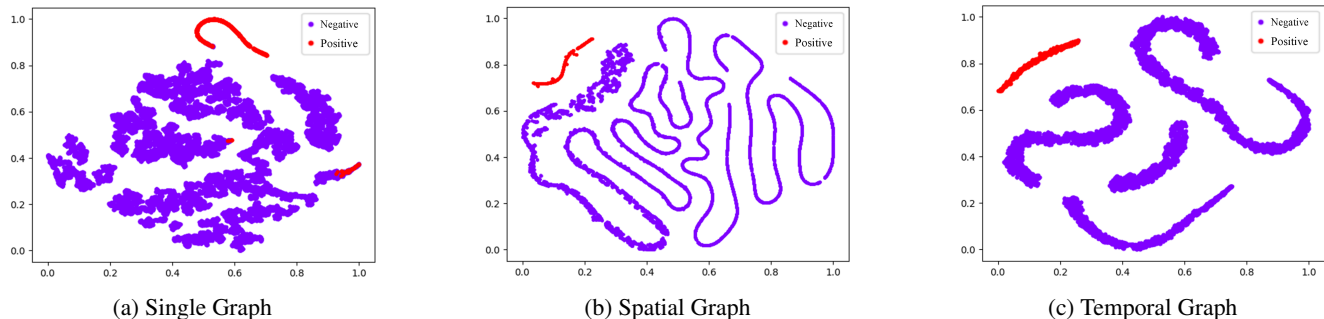


Figure 4: Edge feature clustered by t-SNE[33]. (a) Edge feature of single graph model. (b) Edge feature of  $G^S$ . (c) Edge feature of  $G^T$ . We sample 20 graphs and plot all edge features on figures. Positive and negative are ground truth labels determined by whether the edge comes from the same object. There is a more clear boundary in both spatial graph and temporal graph than single graph, leading to more accurate classification results.

**PETS-09 [11]** The results of S2.L1 sequence are reported for a complete comparison with other methods. Similar to CAMPUS, there are less than 10 people at each frame on average. Although it is not as dense as Wildtrack, the low video quality, e.g. various illumination and cameras far away from people, is the most challenging part.

## 4.2. Implementation Details

OSNet [42] is exploited as ReID model to extract appearance feature, which outputs a 512-D feature vector. Input image for OSNet is object bounding-boxes cropped from the original frame and resized to  $256 \times 128$ . In graph model, the dimension of node feature is 32-D, while edge feature is 6-D. We run  $L = 4$  message passing iterations in all experiments, and then output edge feature for link prediction classifier, which is also a 6-D feature vector.

For model training, we use ground truth detection as input and set the temporal window size  $M=3$  for  $G^T$ . The model is trained by the following settings: Adam optimizer [20] is employed to run 100 epochs. Warm-up learning rate is set starting from 0 to 0.01 in the first 10 epochs. We randomly drop detection for data augmentation to mimic false negative cases. In the inference stage, we use detection from MVDeTr [15] in Wildtrack and provided detection in the other datasets for a fair comparison. The weights of graph model with the highest validation performance are used for testing. The pruning threshold  $\varepsilon$  is set to 0.9 to retain high confident edges only.

## 4.3. Results of MC-MOT

We report our model performance on Wildtrack, CAMPUS, and PETS-09, in Table 1, Table 2, and Table 3, respectively. On Wildtrack, we compare with other online approaches using detector as input. Our results achieve state-of-the-art performance with 3.8% and 7.0% higher than the second place on IDF1 and MOTA. Our experimental results on CAMPUS outperform other approaches on most metrics. One can notice that, for MT and ML, our results on

Method	IDF1 $\uparrow$	MOTA $\uparrow$	MOTP $\uparrow$	MT $\uparrow$	ML $\downarrow$
KSP-DO [7]	73.2	69.6	61.5	28.7	25.1
KSP-DO-ptrack [7]	78.4	72.2	60.3	42.1	14.6
GLMB-YOLOv3 [26]	74.3	69.7	73.2	79.5	21.6
GLMB-DO [26]	72.5	70.1	63.1	<b>93.6</b>	22.8
T-Glimpse [40]	77.8	72.8	79.1	61.0	4.9
T-Glimpse Stack [40]	81.9	74.6	78.9	65.9	4.9
Ours	<b>85.7</b>	<b>81.6</b>	<b>81.8</b>	79.4	<b>4.7</b>

Table 1: Evaluation results on Wildtrack. We achieve state-of-the-art performance with 3.8% and 7.0% progress on IDF1 and MOTA.

Sequence	Method	MOTA $\uparrow$	MOTP $\uparrow$	MT $\uparrow$	ML $\downarrow$
Garden 1	HCT [37]	49	71.9	31.3	6.3
	STP [38]	57	75	-	-
	TRACTA [13]	58.5	74.3	30.6	1.6
	DyGLIP [27]	71.2	91.6	31.3	<b>0.0</b>
	LMGP [25]	76.9	95.9	62.9	1.6
	Ours	<b>77.6</b>	<b>99.1</b>	<b>100.0</b>	<b>0.0</b>
Garden 2	HCT [37]	25.8	71.6	33.3	11.1
	STP [38]	30	75	-	-
	TRACTA [13]	35.5	75.3	16.9	11.3
	DyGLIP [27]	<b>87.0</b>	98.4	66.7	<b>0.0</b>
	Ours	86.0	<b>99.9</b>	<b>100.0</b>	<b>0.0</b>
Parkinglot	HCT [37]	24.1	66.2	6.7	26.6
	STP [38]	28	68	-	-
	TRACTA [13]	39.4	74.9	15.5	10.3
	DyGLIP [27]	72.8	98.6	26.7	<b>0.0</b>
	LMGP [25]	<b>78.1</b>	97.3	62.1	<b>0.0</b>
Ours	77.7	<b>99.8</b>	<b>100.0</b>	<b>0.0</b>	
Auditorium	HCT [37]	20.6	69.2	33.3	11.1
	STP [38]	24	72	-	-
	TRACTA [13]	33.7	73.1	37.3	20.9
	DyGLIP [27]	<b>96.7</b>	<b>99.5</b>	<b>95.2</b>	<b>0.0</b>
	Ours	81.2	98.8	92.1	<b>0.0</b>

Table 2: Evaluation results on CAMPUS. Our model perfectly tracked all people most of time (Figure 1b), leading to perfect scores on MT and ML and competitive results on MOTA.

overlapping FoV sequences achieve 100 and 0, respectively. Even in the sequence with non-overlapping FoV, our model still performs well. This is because our method properly leverages spatial and temporal consistency, which allows

Method	Online	MOTA $\uparrow$	MOTP $\uparrow$	MT $\uparrow$	ML $\downarrow$
KSP [2]		80	57	-	-
TRACTA [13]		87.5	79.2	-	-
DyGLIP [27]	✓	93.5	94.7	-	-
STVH [36]		95.1	79.8	<b>100.0</b>	<b>0.0</b>
MLMRF [22]		96.8	79.9	<b>100.0</b>	<b>0.0</b>
LMGP [25]		<b>97.8</b>	82.4	<b>100.0</b>	<b>0.0</b>
Ours	✓	92.3	<b>99.7</b>	<b>100.0</b>	<b>0.0</b>

Table 3: Evaluation results on PETS-09 sequence S2.L1. Check mark indicates the online method. When compared with other offline methods, our method still achieves very competitive performance.

Appearance	Projection	Speed	IDF1 $\uparrow$	MOTA $\uparrow$
✓		✓	61.5	77.2
✓	✓		86.9	94.8
	✓	✓	89.2	95.0
✓	✓	✓	<b>91.6</b>	<b>97.0</b>

Table 4: Tracking performance between different combination of input feature on Wildtrack. We use ground-truth detection to focus on the impact on association stage.

steady tracking on each object. The results indicate that our method is suitable for handling fragmented tracklet due to occlusion. On PETS-09, our method is competitive compared with other offline methods.

#### 4.4. Ablation Study

To validate the robustness of our model, we conduct several ablation studies in this section.

**Design of Input Feature** Table 4 presents the impact of different input features. If appearance, projection, or speed feature is removed for both nodes and edges, the performance drop goes up. In this study, the projection term is crucial to our method. As for the speed term, it is important for Temporal Association since it provides motion information. There is only a small drop in performance if the appearance feature is removed. Although ReID feature helps to associate object with its appearance, our model does not heavily rely on it. Our method presents better generalization when the illumination or appearance changes drastically across datasets.

**Separated vs. Unified Graph Models** To extract spatial and temporal-aware feature for better association, spatial graph and temporal graph are trained separately with different training input. To compare with the feature learning by one unified spatial-temporal graph, a unified spatial-temporal graph model is trained without any edge constraints and the input data is chunks of frames containing detection across views and frames. Edge features with positive and negative labels are visualized by applying t-SNE [33]. As depicted in Figure 4(a), the negative edge features are messed with positive edge features in some areas, lead-

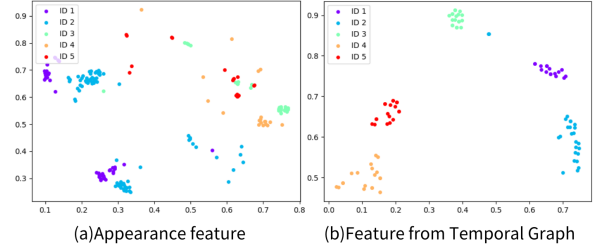


Figure 5: Node feature clustered by t-SNE. (a) Appearance feature extracted by [42]. (b) Feature of ReST model. Each node represents one detection at certain frame. We sample 5 people from 20 consecutive frames in both figures.

Testing Sequence	Training Sequence	MOTA $\uparrow$	MOTP $\uparrow$
Garden 1	Garden 1	77.6	<b>99.1</b>
	Wildtrack	<b>93.1</b>	90.9
Garden 2	Garden 2	86.0	99.9
	Wildtrack	<b>90.2</b>	<b>100.0</b>
Parkinglot	Parkinglot	77.7	<b>99.8</b>
	Wildtrack	<b>92.4</b>	<b>99.8</b>
Auditorium	Auditorium	81.2	98.8
	Wildtrack	<b>95.8</b>	<b>98.9</b>

Table 5: Cross-dataset testing on CAMPUS. We train the model on Wildtrack and test each sequence in CAMPUS.

ing to more failure cases in MC-MOT. In Figure 4(b)(c), a clear boundary between positive and negative edges can be observed. That is, with separate spatial and temporal graphs and their own constraints, the learned edge features are more effective for data association.

**Robust Feature Representation** The quality of feature representation significantly affects the association accuracy in MC-MOT. Precise tracklets can be predicted given better feature representation. In Figure 5, node feature embedding is visualized via t-SNE. In contrast to features extracted by ReID model [42], features learned by ReST have better between-class separation and within-class aggregation. Therefore, we can better discriminate and associate objects by using ReST.

**Cross-dataset Testing** We conduct cross-dataset testing to validate the generalization ability of our model. Specifically, we load model weights trained on Wildtrack and perform inference on CAMPUS dataset. As shown in Table 5, there is a significant improvement in MOTA up to 15.5%. Although Wildtrack has fewer frames than others, it supplies more abundant information for model training, e.g. more difficult occlusion cases and more diverse motion patterns. In addition, lower frame rate makes it extract non-static speed information between frames than other scenes. Furthermore, our model does not heavily rely on appearance feature as shown in Table 4, leading to better generalization results.



## 5. Conclusion

In this paper, we propose a novel reconfigurable graph model for MC-MOT. A two-stage association scheme is proposed via Spatial Association and Temporal Association. It first associates objects across different views at the same frame using spatial graph. Followed by Graph Reconfiguration module which aggregates the nodes within the same connected component to simplify the graph and reconfigures it into a new temporal graph. Lastly, Temporal Association is applied to match objects across frames to accomplish online tracking. The spatial graph and temporal graph are independently trained to concentrate on spatial and temporal-domain feature learning, respectively. As shown in the experimental results, we can learn more discriminating features for object association, leading to state-of-the-art performance on Wildtrack and competitive results on other datasets compared with other offline methods. In the future, we plan to investigate more flexible graph reconfiguration of spatial/temporal or spatial-temporal graph models for MC-MOT.

## 6. Acknowledgements

This work was supported in part by the National Science and Technology Council, Taiwan under grants NSTC-111-2221-E-007-106-MY3, NSTC-111-2634-F-007-010, and NSTC-111-2634-F-194-003-. We also thank National Center for High-performance Computing in Taiwan for providing computational and storage resources.

## References

- [1] Mykhaylo Andriluka, Stefan Roth, and Bernt Schiele. People-tracking-by-detection and people-detection-by-tracking. In *IEEE Conference on Computer Vision and Pattern Recognition(CVPR)*, pages 1–8, 2008.
- [2] Jerome Berclaz, Francois Fleuret, Engin Turetken, and Pascal Fua. Multiple object tracking using k-shortest paths optimization. *IEEE transactions on pattern analysis and machine intelligence(TPAMI)*, 2011.
- [3] Philipp Bergmann, Tim Meinhardt, and Laura Leal-Taixe. Tracking without bells and whistles. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [4] Keni Bernardin and Rainer Stiefelhagen. Evaluating multiple object tracking performance: The clear mot metrics. *EURASIP Journal on Image and Video Processing*, 2008, 01 2008.
- [5] Guillem Braso and Laura Leal-Taixe. Learning a neural solver for multiple object tracking. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [6] Jinkun Cao, Xinshuo Weng, Rawal Khirodkar, Jiangmiao Pang, and Kris Kitani. Observation-centric sort: Rethinking sort for robust multi-object tracking. *arXiv preprint arXiv:2203.14360*, 2022.
- [7] Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. Wildtrack: A multi-camera hd dataset for dense unscripted pedestrian detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition(CVPR)*, pages 5030–5039, 2018.
- [8] Peng Chu and Haibin Ling. Famnet: Joint learning of feature, affinity and multi-dimensional assignment for online multiple object tracking. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [9] Peng Chu, Jiang Wang, Quanzeng You, Haibin Ling, and Zicheng Liu. Transmot: Spatial-temporal graph transformer for multiple object tracking. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 4870–4880, 2023.
- [10] Martin Engilberge, Weizhe Liu, and Pascal Fua. Multi-view tracking using weakly supervised human motion prediction. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 1582–1592, 2023.
- [11] James Ferryman and Ali Shahrokni. Pets2009: Dataset and challenge. In *2009 Twelfth IEEE international workshop on performance evaluation of tracking and surveillance(IEEE)*, pages 1–6, 2009.
- [12] Marco Gori, Gabriele Monfardini, and Franco Scarselli. A new model for learning in graph domains(ieee). In *Proceedings of the IEEE International Joint Conference on Neural Networks*, pages 729–734 vol. 2, 2005.
- [13] Yuhang He, Xing Wei, Xiaopeng Hong, Weiwei Shi, and Yihong Gong. Multi-target multi-camera tracking by tracklet-to-target assignment. *IEEE Transactions on Image Processing(TIP)*, 29:5191–5205, 2020.
- [14] Martin Hofmann, Daniel Wolf, and Gerhard Rigoll. Hypergraphs for joint multi-view reconstruction and multi-object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition(CVPR)*, pages 3650–3657, 2013.
- [15] Yunzhong Hou and Liang Zheng. Multiview detection with shadow transformer (and view-coherent data augmentation). In *Proceedings of the 29th ACM International Conference on Multimedia (MM '21)*, 2021.
- [16] Yunzhong Hou, Liang Zheng, and Stephen Gould. Multiview detection with feature perspective transformation. In *ECCV*, 2020.
- [17] Rudolph Emil Kalman and Others. A new approach to linear filtering and prediction problems. *Journal of basic Engineering*, 82(1):35–45, 1960.
- [18] Aleksandr Kim, Guillem Brasó, Aljoša Ošep, and Laura Leal-Taixé. Polarmot: How far can geometric relations take us in 3d multi-object tracking? In *European Conference on Computer Vision (ECCV)*, pages 41–58, 2022.
- [19] Chanh Kim, Fuxin Li, Arridhana Ciptadi, and James M. Rehg. Multiple hypothesis tracking revisited. In *IEEE International Conference on Computer Vision (ICCV)*, pages 4696–4704, 2015.
- [20] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In Yoshua Bengio and Yann LeCun,

- editors, *3rd International Conference on Learning Representations (ICLR)*, 2015.
- [21] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks. In *International Conference on Learning Representations (ICLR)*, 2017.
- [22] Long Lan, Xinchao Wang, Gang Hua, Thomas S Huang, and Dacheng Tao. Semi-online multi-people tracking by re-identification. *International Journal of Computer Vision (IJCV)*, pages 1937–1955, 2020.
- [23] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *IEEE International Conference on Computer Vision (ICCV)*, pages 2999–3007, 2017.
- [24] Elena Luna, Juan C. SanMiguel, José M. Martínez, and Pablo Carballeira. Graph neural networks for cross-camera data association. *arXiv preprint arXiv:2201.06311*, 2022.
- [25] Duy M. H. Nguyen, Roberto Henschel, Bodo Rosenhahn, Daniel Sonntag, and Paul Swoboda. LMGP: Lifted multicut meets geometry projections for multi-camera multi-object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 8866–8875, 2022.
- [26] Jonah Ong, Ba-Tuong Vo, Ba-Ngu Vo, Du Yong Kim, and Sven Nordholm. A bayesian filter for multi-view 3d multi-object tracking with occlusion handling. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, pages 2246–2263, 2020.
- [27] Kha Gia Quach, Pha Nguyen, Huu Le, Thanh-Dat Truong, Chi Nhan Duong, Minh-Triet Tran, and Khoa Luu. DyGLIP: A dynamic graph model with link prediction for accurate multi-camera multiple object tracking. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13784–13793, 2021.
- [28] Ergys Ristani, Francesco Solera, Roger Zou, Rita Cucchiara, and Carlo Tomasi. Performance measures and a data set for multi-target, multi-camera tracking. In Gang Hua and Hervé Jégou, editors, *Computer Vision – ECCV 2016 Workshops*, pages 17–35, Cham, 2016. Springer International Publishing.
- [29] Shijie Sun, Naveed Akhtar, Xiangyu Song, Huansheng Song, Ajmal Mian, and Mubarak Shah. Simultaneous detection and tracking with motion modelling for multiple object tracking. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 626–643, 2020.
- [30] Siyu Tang, Bjoern Andres, Miykhaylo Andriluka, and Bernt Schiele. Subgraph decomposition for multi-target tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
- [31] Siyu Tang, Bjoern Andres, Mykhaylo Andriluka, and Bernt Schiele. Multi-person tracking by multicut and deep matching. In Gang Hua and Hervé Jégou, editors, *Computer Vision – ECCV 2016 Workshops*, pages 100–111, 2016.
- [32] 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)*, 2017.
- [33] Laurens van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of Machine Learning Research*, 9(86):2579–2605, 2008.
- [34] 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*, volume 30. Curran Associates, Inc., 2017.
- [35] Zhongdao Wang, Liang Zheng, Yixuan Liu, Yali Li, and Shengjin Wang. Towards real-time multi-object tracking. In *Computer Vision–ECCV 2020*, pages 107–122. Springer, 2020.
- [36] Longyin Wen, Zhen Lei, Ming-Ching Chang, Honggang Qi, and Siwei Lyu. Multi-camera multi-target tracking with space-time-view hyper-graph. *International Journal of Computer Vision (IJCV)*, pages 313–333, 2017.
- [37] Yuanlu Xu, Xiaobai Liu, Yang Liu, and Song-Chun Zhu. Multi-view people tracking via hierarchical trajectory composition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4256–4265, 2016.
- [38] Yuanlu Xu, Xiaobai Liu, Lei Qin, and Song-Chun Zhu. Cross-view people tracking by scene-centered spatio-temporal parsing. In *Proceedings of the AAAI Conference on Artificial Intelligence (AAAI)*, 2017.
- [39] Yihong Xu, Aljosa Osep, Yutong Ban, Radu Horaud, Laura Leal-Taixe, and Xavier Alameda-Pineda. How to train your deep multi-object tracker. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [40] Quanzeng You and Hao Jiang. Real-time 3d deep multi-camera tracking. *arXiv preprint arXiv:2003.11753*, 2020.
- [41] Yifu Zhang, Peize Sun, Yi Jiang, Dongdong Yu, Fucheng Weng, Zehuan Yuan, Ping Luo, Wenyu Liu, and Xinggang Wang. Bytetrack: Multi-object tracking by associating every detection box. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022.
- [42] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Omni-scale feature learning for person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019.
- [43] Xingyi Zhou, Vladlen Koltun, and Philipp Krähenbühl. Tracking objects as points. In *Computer Vision – ECCV 2020*, pages 474–490, 2020.