

Multiple People Tracking Using Hierarchical Deep Tracklet Re-identification

Maryam Babae* Ali Athar* Gerhard Rigoll

Institute for Human-Machine Communication, Technical University of Munich
Arcisstrasse 21, Munich, Germany

{maryam.babae@tum.de, ali.athar@tum.de, rigoll@tum.de}

Abstract

The task of multiple people tracking in monocular videos is challenging because of the numerous difficulties involved: occlusions, varying environments, crowded scenes, camera parameters and motion. In the tracking-by-detection paradigm, most approaches adopt person re-identification techniques based on computing the pairwise similarity between detections. However, these techniques are less effective in handling long-term occlusions. By contrast, tracklet (a sequence of detections) re-identification can improve association accuracy since tracklets offer a richer set of visual appearance and spatio-temporal cues. In this paper, we propose a tracking framework that employs a hierarchical clustering mechanism for merging tracklets. To this end, tracklet re-identification is performed by utilizing a novel multi-stage deep network that can jointly reason about the visual appearance and spatio-temporal properties of a pair of tracklets, thereby providing a robust measure of affinity. Experimental results on the challenging MOT16 and MOT17 benchmarks show that our method significantly outperforms state-of-the-arts.

1. Introduction

Multi-object tracking (MOT) is a key problem in computer vision with many applications such as video surveillance, activity analysis, and abnormality detection [2, 1, 3]. It is challenging in unconstrained environments due to influencing factors such as illumination variance, camera motion, target interactions, and more importantly, lengthy occlusions.

Most existing multi-object tracking methods fall into the tracking-by-detection category, where the goal is to link detections in the video belonging to the same target. Recent tracking methods adopt person re-identification techniques based on pairwise similarity of detections [52, 34] for this data association. However, this can lead to wrong associa-

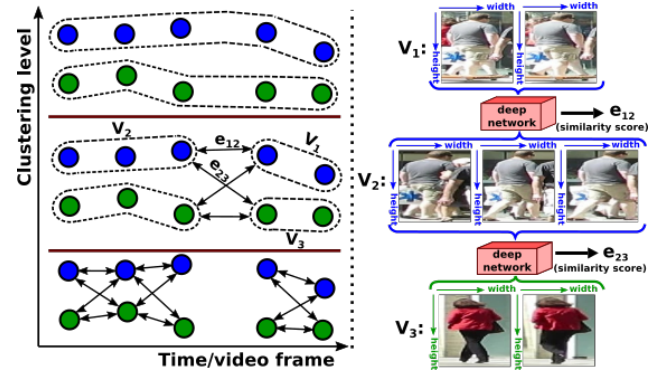


Figure 1: Overview of the proposed people tracking framework that hierarchically clusters tracklets and employs a deep network to evaluate tracklet similarity. Circles denote vertices in the graph and their colors reflect their person IDs.

tions, especially if there are lengthy occlusions. By contrast, considering a group of detections before and after an occlusion as a tracklet can improve the re-identification accuracy.

Moreover, we argue that these association errors can be further reduced if pedestrian tracking is formulated as a hierarchical clustering problem that iteratively merges detections into longer tracklets. This way, the association complexity increases gradually as opposed to a one-step approach that directly aims to obtain the final solution.

When evaluating possible associations between detections in crowded scenes where multiple pedestrians are closely located and/or overlapping in the image, it is essential to jointly reason about both their visual appearance and spatio-temporal properties. Recently, several works [47, 23, 52] use neural networks to process visual appearance, and separately compute hand-crafted features to incorporate spatio-temporal information from the bounding boxes. Logistic regression or some learning technique is then used to assign weights to these features in order to compute an overall similarity metric. Even though this mitigates the need to empirically set weights through trial-and-error, we argue that hand-crafted features nonetheless

*equal contribution

do not generalize well since they make certain assumptions about the underlying motion model, and as in [47], some of these features may have to be computed separately for each video sequence to account for the difference in camera parameters. Additionally, such approaches lack the ability to jointly reason about spatio-temporal and visual cues in a strong manner since the features are computed separately and combined only at the final step.

In this paper, we propose a multiple people tracking framework (illustrated in Fig. 1) that hierarchically merges tracklets to overcome occlusions and minimize association errors. Our main contributions are: (1) A novel end-to-end deep network for assessing tracklet similarity that can jointly reason about visual and spatio-temporal cues in a generalized manner without requiring hand-crafted features and/or tunable parameters; (2) an extension of Kernighan-Lin with Joins algorithm [27] that enables the tracklet clustering problem to be formulated as a constrained minimum-cost multicut graph problem, and; (3) a new state-of-the-art in the MOT Challenge [35].

2. Related Work

Most multi-object tracking approaches are based on the tracking-by-detection paradigm [36, 25, 19], where tracking is formulated as a data association problem between the detections extracted from a video using object detectors.

Data association can be performed either on individual detections [36, 6], or a set of confident and short tracklets [55, 33] which are generated by first performing a low level data association to group detections. A well-known representation of the tracking-by-detection paradigm is to present each detection as a node in a graph, with edges representing the likelihood that connected detections belong to the same person. This data association problem can be solved using Conditional Random Field inference [57], network flow optimization [60, 5], maximum multi-clique [13], greedy algorithms [43], or subgraph decomposition [50].

By learning discriminative feature representations, deep learning has enhanced many computer vision applications such as image classification [32], video background subtraction [4], and pedestrian detection [42]. In the context of tracking, Convolutional Neural Networks (CNN) have been utilized to learn feature representations of targets instead of using heuristic and hand-crafted features [54, 37, 56]. CNNs have also been utilized for modeling the similarity between a pair of detections [34, 52]. [48] models the appearance with temporal coherency by designing a quadruplet CNN. Adopting a different network structure, Milan *et al.* [41] propose an end-to-end Recurrent Neural Network (RNN) for the data association problem in online multi-target tracking. They use RNNs for target state prediction, and to determine a track’s birth/death in each frame.

Among other online multi-target tracking approaches which are based on tracklet-detection matching, [61] exploits structural invariance constraint and develops a probability frame that is able to jointly reason about both appearance and structure cues for an object-detection pair. In [62], the authors propose an online tracking method using dual matching attention networks with both spatial and temporal attention mechanisms. In [16], a temporal generative modeling framework is proposed that uses a recurrent autoregressive network to characterize the appearance and motion dynamics of multiple objects over time. In [38], a novel scoring function based on a fully convolutional network is presented to perform optimal selection from a large number of candidates in real-time. [12] utilizes the merits of single object trackers using shared CNN features and Region of Interest (ROI) pooling. In addition, a spatial-temporal attention mechanism was adopted to alleviate the problem of drift caused by frequent occlusions.

Recently, Ma *et al.* [39] presented a framework that employs a three step process in which tracklets are first created, then cleaved, and then reconnected using a combination of Siamese-trained CNNs, Bi-Gated Recurrent Unit (GRU) and LSTM cells. By contrast, our approach utilizes a hierarchical clustering mechanism with a single multi-stage network to compute tracklet similarity, thereby minimizing false associations in the first step and mitigating the need for tracklet cleaving and reconnection. In [45], a multi-stage network was proposed to model the appearance, motion and interaction of targets. Their network design is similar to ours, but with the key difference that our model computes the similarity between two tracklets, rather than between a tracklet and a single detection.

3. Approach

The proposed framework hierarchically merges tracklets to reduce association errors. It utilizes a deep network for tracklet re-identification that computes pairwise similarity scores between tracklets by jointly learning visual and spatio-temporal features. This network consists of a CNN that learns pairwise detection visual appearance, and two bidirectional RNNs that learn spatio-temporal features, and aggregate visual and spatio-temporal features, respectively. Hierarchical clustering is formulated as a series of constrained minimum cost multicut graph problems with vertices representing tracklets, and edges representing tracklet similarities as computed by the network.

3.1. Deep Tracklet Re-identification

Before elaborating the network architecture, let us define the following nomenclature: a tracklet T_i^N , uniquely identified by i , is defined as a collection of N detections $\{D_i^1, D_i^2, \dots, D_i^{N-1}, D_i^N\}$ such that $N \in [1, N_{\max}]$, subject to the constraint that no more than one detection is allowed

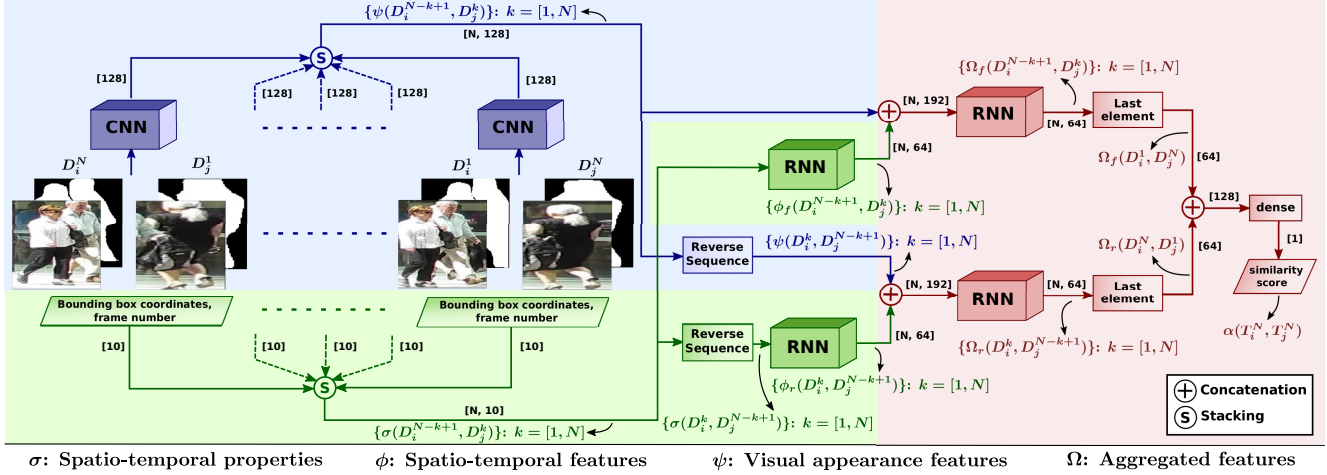


Figure 2: Block diagram of the end-to-end network. The region highlighted in blue is relevant to visual appearance features, the region in green to spatio-temporal features, and the region in red to the combined representation of both. All CNN blocks share the same parameters. The feature dimensions are written in square brackets.

in any given image frame of the video sequence. Let $F(D_a)$ denote the frame number in which detection D_a lies. Also assume that the detections of T_i^N are sorted in ascending order of frame number, i.e., $F(D_i^1) < F(D_i^2) < \dots < F(D_i^N)$. Geometrically, a detection D_a is a rectangular bounding box in the image plane that is described by the tuple $\sigma(D_a) = [X_a, Y_a, W_a, H_a, F(D_a)] \in \mathbb{R}^5$, where X_a and Y_a are the top-most and left-most pixel coordinates, respectively, and W_a and H_a are the width and height, respectively. Before any $\sigma(\cdot)$ is input to the network, the bounding box dimensions are normalized by the image dimensions, and offset by the coordinates of the first detection in a given tracklet pair. Similarly, the frame number of the first detection in a given tracklet pair is considered 0, and all subsequent frame numbers are normalized by the frame rate of the video.

Let us further define $\alpha(T_i^L, T_j^M) \in [0, 1]$ as the probability of tracklets T_i^L and T_j^M belonging to the same person. Here, we impose the constraint that $\alpha(T_i^L, T_j^M)$ can only be computed when tracklet T_i^L precedes T_j^M in the video sequence with no overlapping frames, i.e. $F(D_i^1) \leq F(D_i^L) < F(D_j^1) \leq F(D_j^M)$. Since the framework processes detections pairwise, the number of detections in both tracklets is reduced to $N = \min(L, M, N_{\max})$ by removing the first $L - N$ detections from T_i^L and the last $M - N$ detections from T_j^M . Let us refer to these pruned versions of T_i^L and T_j^M as T_i^N and T_j^N , respectively.

3.1.1 Visual Appearance Feature Learning

To learn visual appearance features, we employ a CNN based on the ResNet-50 architecture [22] which compares a pair of detections and outputs the probability of those de-

tections belonging to the same person. The input to this network is a pair of RGB detection images, along with a binary body mask for both detections that is active at pixel locations occupied by persons. The motivation behind incorporating the body mask is to focus the CNN’s attention on the relevant part of the image so that it becomes more sensitive to changes in the person’s appearance and learns to ignore background changes. These masks are generated using a pre-trained Mask-RCNN [21]. The RGB images and body masks are resized to 128x128. Since this dimension size is roughly half that used in [22], the first convolutional filter of our CNN is 5x5 instead of 7x7, and a stride of 1 is used when applying this filter instead of 2.

The input tensor dimensions are thus 128x128x8 (2x3 RGB image channels and 2x1 binary body masks). The output from the convolutional layers is flattened and input to a dense layer which reduces the feature size to 128. Given the input detection pair (D_a, D_b) , let us refer to this feature vector as $\psi(D_a, D_b) \in \mathbb{R}^{128}$. $\psi(D_a, D_b)$ is then input to a classification layer that contains a single neuron with sigmoid activation that outputs the probability of detections D_a and D_b belonging to the same person.

To compute $\alpha(T_i^L, T_j^M)$, the first step is to apply the CNN to compute the pairwise similarity features for the following sequence of detection image pairs in T_i^N and T_j^N : $\{(D_i^{N-k+1}, D_j^k)\}$ for $k \in [1, N]$. This results in a sequence of N feature vectors $\{\psi(D_i^{N-k+1}, D_j^k)\}$ for $k \in [1, N]$. An illustration of the detection pairs in a pair of tracklets is given in Fig. 3 using black, curved arrows.

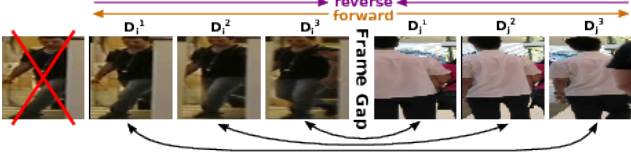


Figure 3: Illustration of the pairwise detections and their ordering for two example tracklets T_i^L (left) and T_j^M (right) with $L = 4, M = 3$ separated by a frame gap. The first detection of T_i^L is pruned since $N = \min(4, 3) = 3$ (ignoring N_{\max}).

3.1.2 Spatio-temporal Feature Learning

Separately, the sequence of spatio-temporal properties of the detection pairs belonging to tracklets T_i^N and T_j^N is input to a bidirectional RNN. Formally speaking, the sequence $\{\sigma(D_i^{N-k+1}) \oplus \sigma(D_j^k)\}$ for $k \in [1, N]$ (where \oplus denotes concatenation of two vectors), and its reversed version $\{\sigma(D_i^k) \oplus \sigma(D_j^{N-k+1})\}$ for $k \in [1, N]$, are input to two separate series of GRU cells of size 64, resulting in the two output sequences $\{\phi_f(D_i^{N-k+1}, D_j^k)\}$ and $\{\phi_r(D_i^k, D_j^{N-k+1})\}$ for $k \in [1, N]$, respectively. Intuitively, these sequences encode the spatio-temporal features of the pairwise combinations of detection bounding boxes in tracklets T_i^N and T_j^N . Note that while the visual appearance features are independent of the sequence in which the pairs of detections occur, spatio-temporal features are not. Fig. 3 illustrates the direction of the forward and backward sequences.

3.1.3 Feature Aggregation

The visual and spatio-temporal features of the tracklets are then concatenated, and input to another bidirectional RNN. Formally speaking, the sequences $\{\psi(D_i^{N-k+1}, D_j^k) \oplus \phi_f(D_i^{N-k+1}, D_j^k)\}$ and $\{\psi(D_i^k, D_j^{N-k+1}) \oplus \phi_r(D_i^k, D_j^{N-k+1})\}$ for $k \in [1, N]$ are input to two series of GRU cells of size 64, resulting in the two output sequences $\{\Omega_f(D_i^{N-k+1}, D_j^k)\}$ and $\{\Omega_r(D_i^k, D_j^{N-k+1})\}$ for $k \in [1, N]$, respectively. Intuitively, these features offer a combined representation of the visual and spatio-temporal features of the detections in the tracklets. Since we are interested in a single similarity score that considers both sequences in their entirety, we retain only the last two elements of the output, i.e., $\Omega_f(D_i^1, D_j^N)$ and $\Omega_r(D_i^N, D_j^1)$. Even though it is not reflected in the notation used, these two features actually incorporate information from the entire sequence of detections, because the input to an RNN cell consists of the input at the current time-step, as well as the output from the cell at the previous time-step. Finally, $\Omega_f(D_i^1, D_j^N)$ and $\Omega_r(D_i^N, D_j^1)$ are concatenated and input to a single neuron with sigmoid activation that outputs the final similarity score $\alpha(T_i^L, T_j^M)$. A block diagram of the complete network is provided in

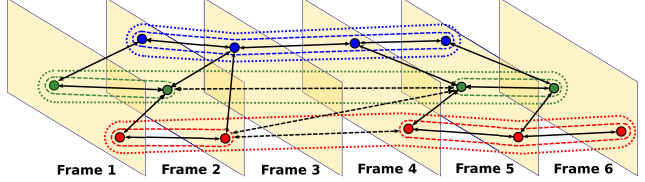


Figure 4: Illustration of the hierarchical clustering process. Detections are denoted by circles, and each color corresponds to a person ID. In the first iteration, edges between adjacent detections are created (shown by solid black arrows). The resulting clusters shown by the dashed colored lines form the vertices for the second iteration, in which longer edges shown by the dashed black arrows are created. The clustering result of the second iteration is shown by the outer-most dotted lines.

Fig. 2.

3.2. Hierarchical Clustering

The task of clustering tracklets globally given their pairwise similarities is formulated as a minimum cost multicut graph problem (MP) [11, 15]. Given a graph $G = (V, E)$, tracklets are modeled as vertices V in the graph, and undirected edges E allow pairs of tracklets to be checked for similarity and merged. Letting c_e denote the cost of an edge e , the MP can be defined as:

$$Y_E^* = \min_{y \in Y_E} \sum_{e \in E} c_e y_e \quad (1)$$

$$\text{s.t. } \forall C \in \text{cycles}(G), \forall e \in C : y_e \leq \sum_{e' \in C \setminus \{e\}} y_{e'} \quad (2)$$

where $y_e \in \{0, 1\}$ is a 0/1-label assigned to edge e . A zero indicates a 'join', i.e., the vertices connected by the edge belong to the same component, whereas a one indicates a 'cut', i.e., the vertices belong to separate components. The output of the MP, Y_E^* , defines a valid decomposition of the graph into one or more disjoint components, such that the sum of costs of the cut edges is minimized (note that the number of resulting components does not have to be specified in advance). Eq. (2) defines transitivity constraints which guarantee that the decomposition is well-defined. It follows that if edge costs denote the similarity between vertices, then the MP can be directly applied to the tracklet clustering problem.

Initially, all detections in the video sequence are assumed to be separate tracklets/vertices (both terms will be used interchangeably from here onward). Edges are then created between vertices in adjacent frames, and their costs are computed using the similarity score obtained from the visual appearance matching network described in Sec. 3.1.1. We then apply the Constrained Kernighan-Lin with Joins algorithm (described ahead) to compute a feasible decomposition of this graph; all detections that belong to the same

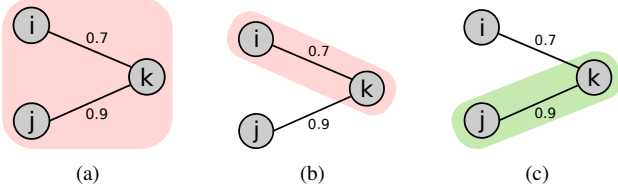


Figure 5: Illustration of the motivation behind extending the KLJ algorithm. The circles represent the three tracklets with single detections. The probability scores have been directly shown as edge weights here for ease of understanding.

component are subsequently merged into a single, longer tracklet. For the next iteration, edges are created between these newly merged tracklets, and the graph is again decomposed using the minimum cost multicut algorithm. Likewise, this process repeats until no more tracklets can be merged. As the tracklets become longer in subsequent iterations, we allow longer edges to be created that span over increasingly larger frame gaps in order to overcome occlusions. Moreover, for tracklets containing more than one detection, the complete network is used to compute the similarity score. This design choice will be justified in Sec. 4.3, but the underlying idea is that once tracklets become longer, the complete network offers improved performance since the RNN is able to leverage sequential patterns in the learned features. An abstract example of the clustering process is provided in Fig. 4. Lastly, note that the tracklet similarity scores $\in [0, 1]$ output by the network are mapped onto the range $[-\infty, \infty]$ by the following function to obtain the edge costs. This results in dissimilar edges having negative costs, encouraging the algorithm to cut them.

$$f(x) = \log\left(\frac{x}{1-x}\right) \quad (3)$$

Constrained Kernighan-Lin with Joins

Since the MP is NP-hard [7, 14], it is normally not feasible to compute a globally optimal solution. In [27], a generalization of the MP, namely the minimum cost Lifted Multicut Problem (LMP), is proposed, and an extension of the Kernighan-Lin algorithm [26], called Kernighan-Lin with Joins (KLJ), is presented to solve the problem. In this work, we propose a straightforward extension of the KLJ algorithm called Constrained Kernighan-Lin with Joins (CKLJ), and employ it to solve the MP.

The motivation behind extending the algorithm is that when KLJ is applied to the tracklet clustering problem, it often outputs invalid results where multiple tracklets are assigned to the same component even though some of their detections lie in the same frame. As a basic example, consider three tracklets T_i^1 , T_j^1 and T_k^1 that contain only a sin-

gle detection. Suppose that T_i^1 and T_j^1 lie in the first frame of the video whereas T_k^1 lies in the second frame. Now, if we create edges (T_i^1, T_k^1) , and (T_j^1, T_k^1) , and if, for some reason (appearance similarity or spatio-temporal proximity), the similarity scores $\alpha(T_i^1, T_k^1)$ and $\alpha(T_j^1, T_k^1)$ are both high, then the KLJ algorithm will not cut either edge, resulting in all three detections being assigned to the same component, as illustrated in Fig. 5a. Note that this happened even though there is no direct edge between T_i^1 and T_j^1 . To overcome this, CKLJ accepts a set of constraint pairs as input, where each pair (a, b) defines a constraint that tracklets T_a and T_b cannot be assigned to the same component. Since the KLJ algorithm reduces the total cost by greedily merging/splitting components and swapping vertices between them, such constraints can be easily incorporated by imposing a conditional check prior to executing these transformations. For the current example, we would thus provide (i, j) as a constraint to the algorithm.

Applying such constraints, however, gives rise to a new problem. Referring to the same example again, suppose that $\alpha(T_i^1, T_k^1) = 0.7$ and $\alpha(T_j^1, T_k^1) = 0.9$. Naturally, we would want the first edge to be cut, and the second one to be retained, but with the existing KLJ implementation, this may not happen, because when the algorithm tries to lower the total cost by merging two clusters or swapping vertices between them (in the "update_bipartition" function in Alg. 2 in [27]), it iterates through the neighboring vertices in an undefined order. Therefore, it may happen that the algorithm encounters the edge (T_i^1, T_k^1) first, joins it, and then later it is forced to cut T_j^1 and T_k^1 because joining it would violate a constraint (Fig. 5b illustrates this case). As a simple remedy to this problem, we make the CKLJ algorithm greedy by first sorting the neighboring components in descending order of the edge cost. This ensures that high similarity edges are joined first and the aforementioned scenario is avoided, as shown in Fig. 5c.

4. Experiments

4.1. CNN Architecture

To assess the performance of various CNN architectures for visual appearance matching, a test set was created using ground truth data from the MOT Benchmark 2015¹. This set contains 7836 samples, each containing a detection pair belonging to either the same person or to different persons, and separated by different frame gaps. A CNN based on VGG-16 [46] was applied, and we also trained our ResNet-50 based network on detection images without the binary body mask. As an alternative to stacking detection image pairs and inputting them to the network, we also trained a triple network [53, 24] to learn discriminative visual embeddings of size 128 from the detection images individually. To

¹https://motchallenge.net/data/2D_MOT_2015

Network	Body Mask	Accuracy (%)
Triple Network	Yes	84.9
VGG-16	Yes	85.4
ResNet-50	No	86.3
ResNet-50	Yes	87.7

Table 1: Comparison of CNNs for detection matching accuracy.

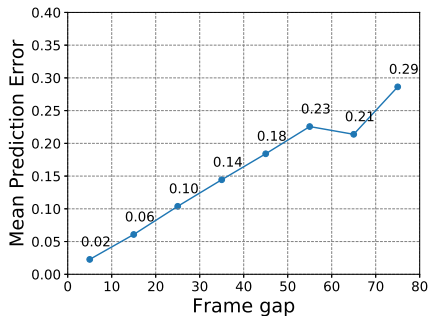


Figure 6: CNN detection matching error against frame gap.

this end, a network based on ResNet-50 was used to extract the embeddings, and the hinge loss function defined in [53] with a margin of 0.2 was used to train the network. An on-line smart mining strategy was used to create suitable training triplets based on approximate nearest neighbor search, as described in [20]. To assign a binary similarity label to each sample, we computed the normalized L_2 -distance between the feature embeddings of both detections, and selected an optimal threshold for classification across all samples (this came out to be 0.43). For the other networks, the predicted labels were obtained by thresholding the output at 0.5. The obtained results are presented in Table 1. It is evident that the ResNet-50 based network with stacked detection images and body masks gave the highest accuracy. It also converged faster during training, and had a lower validation loss than the other networks.

4.2. Effect of Temporal Distance on CNN Accuracy

One motivation behind comparing tracklets instead of detections is that the accuracy of detection matching networks deteriorates as the frame gap between them increases. To demonstrate this, we created another test set from the MOT Benchmark 2015, and plotted the mean prediction error (the average of the absolute difference between the predicted label $\in [0, 1]$ output by the network and the true label $\in \{0, 1\}$ for all test samples) against the frame gap in Fig. 6. This test set contains only same-person detection pair samples, since the visual (dis)similarity for different-person samples is largely independent of the frame gap. The results show an almost exactly linear relationship between the two parameters, which supports our claim.

4.3. Ablation Study

To justify the use of RNNs in our tracklet matching network, we performed an ablation study in which the effect of the network’s design, and the presence of visual and spatio-temporal features towards the final classification accuracy is analyzed for various track lengths and frame gaps. The different network configurations used in this study are:

- Spatio-temporal and visual (**ST+V**): This is the baseline configuration which uses both spatio-temporal and visual features as described in Sec. 3.1.
- Spatio-temporal (**ST**): The visual features are omitted, i.e., only the spatio-temporal properties are input to the RNN.
- Visual Sequence (**V**): The spatio-temporal features are omitted, i.e., only the visual features output by the CNN are fed to the RNN.
- Visual (**CNN only**): The RNN is omitted entirely. The tracklet similarity score is calculated by taking the mean of the similarity score output by the CNN for all detection pairs in the tracklet pair.

The RNN was trained separately for each of these configurations (except for the ‘CNN only’ setting), and applied to a test set created from ground truth data from the MOT Benchmark 2015. Each sample in the test set contains a pair of tracklets that may belong to the same person (positive sample) or to different persons (negative sample), and are separated by a frame gap. The length of these tracklets is either 1, 5 or 10, and there are 15000-25000 samples for each of these lengths with roughly equal positive and negative samples. In Fig. 7, the mean prediction error (computed in the same way as in Sec. 4.2) is plotted against frame gap for the four network configurations and three track lengths.

For tracklet length = 1, Fig. 7a shows that using only the CNN similarity score results in the lowest error for all frame gaps. This is because there is almost no useful information in the spatio-temporal features if each tracklet only contains a single detection, as evident from the high error for the ‘ST’ configuration. In fact, the ‘ST’ error is approximately 0.5 regardless of the frame gap, meaning that the network’s prediction is no better than a random guess. Furthermore, an RNN is only able to extract meaningful information if it has observed a longer sequence of inputs.

For tracklet length = 5 (Fig. 7b), the trend changes. The spatio-temporal properties now offer useful cues to the network, as evident from the substantially lower error for the ‘ST’ plot for small frame gaps. Also note that the ‘V’ configuration outperforms the ‘CNN only’ configuration for all frame gaps, even though both only utilize visual appearance features. This shows that the ability of the RNN to

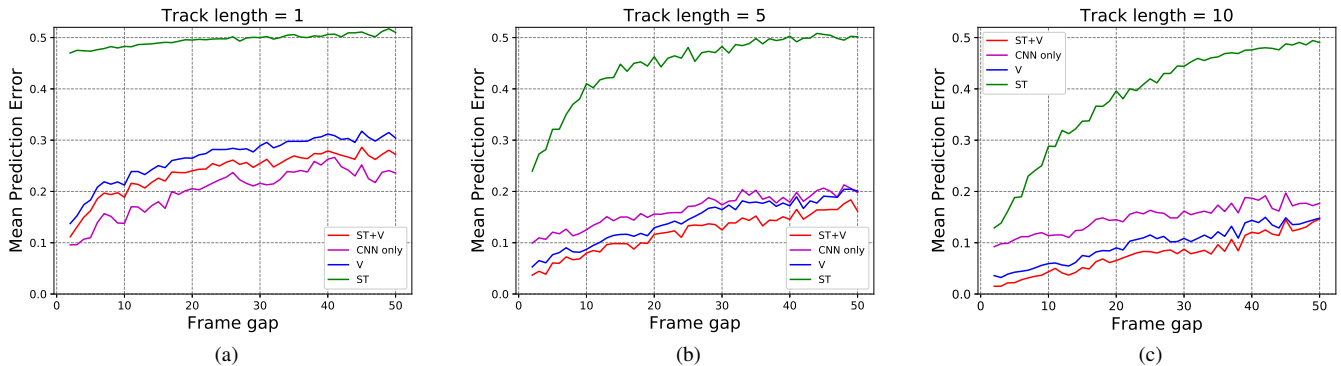


Figure 7: Mean prediction error of the RNN for various frame gaps and network inputs.

learn sequential patterns from the input yields improved performance when a longer sequence is provided. Lastly, the baseline 'ST+V' configuration emerges as the best performer, lending credibility to our claim that an end-to-end network is able to effectively learn and aggregate different types of features. For tracklet length = 10 (Fig. 7c), the same general trend continues; the error for the 'ST' configuration further reduces since spatio-temporal features become more informative, and the performance of the 'CNN only' configuration deteriorates further compared to the 'V' and 'ST+V' configurations. We conjecture that this occurs because the RNN is able to better reason about the sequential pattern of the provided features when the input sequence length is larger.

4.4. Training

For the CNN, a training set with 48954 detection pairs was created from ground truth data in the MOT17 training datasets. Since the ground truth detections bounding boxes are exact, we use the detection boxes output by various detectors that overlap significantly with the ground truth boxes to create training samples. This helps to make the CNN more robust to inaccurate detection boxes. To reduce over-fitting, we employ runtime image augmentation by introducing random brightness offsets and horizontal flips. Dropout [49] with keep probability of 0.8 is also applied to the final feature vector $\psi(\cdot)$ during training. A learning rate of 0.002 with a decay factor of 0.94 after every 7000 iterations was used for optimization using Stochastic Gradient Descent (SGD) with a momentum factor of 0.9.

For the RNN, we created 22640 tracklet pairs with roughly equal numbers of positive and negative samples from the MOT17 training datasets. The samples have tracklet lengths ranging from 1 to 20, and the tracklets are separated by frame gaps ranging from 0 to 4 times the tracklet length. Positive samples are created by splitting a known ground truth track at various points. Negative samples are

created in three ways: (1) a portion of another person's track is extracted such that the bounding box center coordinates of the detections of this track are closest to that of the original person's detections in the same frame. This improves the network's performance in cases where the spatial coordinate information is ambiguous. Around 50% of negative samples are created in this manner. (2) Of the remaining, 25% are created by dividing the image plane into four equally sized quadrants, and sampling a portion of another person's track such that the detection centers lie in the same quadrant as that of the original person. (3) The final 25% are created similarly, but by sampling from a track whose detections lie in any other quadrant (i.e., these are easy negatives). All three sample creation techniques are detailed in App. A.

When training the RNN, the weights of the CNN are frozen due to memory constraints. The RNN is trained with a learning rate of 0.002 with a decay factor of 0.95 after every 2000 iterations, and optimized using SGD with a momentum factor of 0.9. Dropout [49] with keep probability 0.5 is also applied to prevent over-fitting.

4.5. Clustering Scheme

The tracklets are iteratively merged to form longer tracklets using CKLJ, as explained in Sec. 3.2. For the first three iterations, the maximum permitted frame gap between vertices is restricted to 1, 2 and 4, respectively. Thereafter, the frame gap is allowed to be at most four times the tracklet length. When no more tracklets can be merged, the maximum allowed frame gap restriction is further relaxed to at most six times the tracklet length. Once no more tracklets can be merged under this setting, the clustering process is said to be complete. Note that this scheme and the associated parameters have been chosen ad-hoc, with the aim of balancing fast convergence and gradual relaxation of the frame gap restriction in a manner that is applicable to both stationary and moving camera videos. We also remark that

Tracker	MOTA(%) \uparrow	MOTP(%) \uparrow	FAF \downarrow	MT(%) \uparrow	ML(%) \downarrow	FP \downarrow	FN \downarrow	ID Sw. \downarrow	Frag \downarrow
LMP [52]	48.8	79.0	1.1	18.2	40.1	6654	86245	481	595
GCRA [39]	48.2	77.5	0.9	12.9	41.1	5104	88586	821	1117
FWT [23]	47.8	75.5	1.5	19.1	38.2	8886	85487	852	1534
MOTDT [38]	47.6	74.8	1.6	15.2	38.3	9253	85431	792	1858
NLLMPa [36]	47.6	78.5	1.0	17.0	40.4	5844	89093	629	768
AMIR [45]	47.2	75.8	0.5	14.0	41.6	2681	92856	774	1675
MCjoint [29]	47.1	76.3	1.1	20.4	46.9	6703	89368	370	598
NOMT [10]	46.4	76.6	1.6	18.3	41.4	9753	87565	359	504
JMC [51]	46.3	75.7	1.1	15.5	39.7	6373	90914	657	1114
HDTR (Ours)	53.6	80.8	0.8	21.2	37.0	4714	79353	618	833

Table 2: Tracking results on the MOT16 test dataset with public detections. \uparrow and \downarrow represent higher is better and lower is better, respectively. The values in bold and blue represent the best and second best performances, respectively.

Tracker	MOTA(%) \uparrow	MOTP(%) \uparrow	FAF \downarrow	MT(%) \uparrow	ML(%) \downarrow	FP \downarrow	FN \downarrow	ID Sw. \downarrow	Frag \downarrow
FWT [23]	51.4	77.0	1.4	21.4	35.2	24101	247921	2648	4279
jCC [28]	51.2	75.9	1.5	20.9	37.0	25937	247822	1802	2984
MOTDT17 [38]	50.9	76.6	1.4	17.5	35.7	24069	250768	2474	5317
MHT_DAM [30]	50.7	77.5	1.3	20.8	36.9	22875	252889	2314	2865
EDMT17 [9]	50.0	77.3	1.8	21.6	36.3	32279	247297	2264	3260
HAM_SADF17 [59]	48.3	77.2	1.2	17.1	41.7	20967	269038	1871	3020
DMAN [62]	48.2	75.7	1.5	19.3	38.3	26218	263608	2194	5378
PHD_GSDL17 [18]	48.0	77.2	1.3	17.1	35.6	23199	265954	3998	8886
MHT_bLSTM [31]	47.5	77.5	1.5	18.2	41.7	25981	268042	2069	3124
HDTR (Ours)	54.1	80.2	1.0	23.3	34.8	18002	238818	1895	2693

Table 3: Tracking results on the MOT17 test dataset.

changing the parameters within a reasonable range does not effect our framework’s performance significantly. Details of the edge creation method employed, and a quantitative analysis of the convergence and computational time of this clustering scheme are presented in App. B, C and D, respectively.

4.6. Multi-object Tracking Benchmark Results

To assess our framework’s performance, we applied it to test datasets from the MOT16 and MOT17 challenge. The MOT16 test dataset contains 7 video sequences captured in different imaging conditions with varying camera motions and camera angles. The MOT17 challenge contains the same video sequences, but offers detections from 3 different person detectors. Both challenges use the CLEAR MOT performance metrics [8] to rank tracker performance, which include Multi-Object Tracking Accuracy (MOTA) and Precision (MOTP), average False Alarms per Frame (FAF), ratio of Mostly Tracked (MT) and Mostly Lost (ML) targets, False Positives (FP), False Negatives (FN), ID switches (ID Sw.) and trajectory fragmentations (Frag.).

As evident from Tables 2 and 3, our tracker outperforms all other published works on both MOT16 and MOT17 challenges in terms of overall accuracy (MOTA) by an impressive margin of 4.8% and 2.7%, respectively. For the other metrics, it is mostly either ranked first or second. Specifically, our approach more reliably matches tracklets across occlusions, which is evident from our high MT and low

ML scores, and also from the lower false negative count. Moreover, compared to [52] where tracking is performed in one step graph optimization by clustering detections, our approach achieves better results by hierarchically clustering tracklets. Lastly, we recognize ID switches as an area of possible improvement. These switches occur more frequently when there is significant camera motion, which makes spatio-temporal cues less reliable, thus causing the RNN performance to deteriorate. The detailed per video sequence results and annotated videos are available online² and in App. E.

5. Conclusion

We proposed a multi-object tracking framework that hierarchically merges tracklets to effectively resolve lengthy occlusions. Tracklet clustering is formulated as a constrained minimum cost multicut problem and solved using the Constrained Kernighan Lin with Joins Algorithm. To compute similarity metrics between tracklets, a novel deep network was employed that learns and jointly reasons about spatio-temporal and visual appearance features. The framework’s design choices were justified by performing various experiments, and finally, its effectiveness was demonstrated by showing its state-of-the-art performance on the MOT Challenge.

²https://motchallenge.net/tracker/HDTR_16,
https://motchallenge.net/tracker/HDTR_17

References

- [1] J. K. Aggarwal and M. S. Ryoo. Human activity analysis: A review. *ACM Computing Surveys (CSUR)*, 43(3):16, 2011. 1
- [2] A. Alahi, V. Ramanathan, and L. Fei-Fei. Socially-aware large-scale crowd forecasting. In *2014 IEEE Conference on Computer Vision and Pattern Recognition*, pages 2211–2218, 2014. 1
- [3] D. G. Aviv. Abnormality detection and surveillance system, 1997. US Patent 5,666,157. 1
- [4] M. Babae, D. T. Dinh, and G. Rigoll. A deep convolutional neural network for video sequence background subtraction. *Pattern Recognition*, 76:635–649, 2018. 2
- [5] M. Babae, Y. You, and G. Rigoll. Pixel level tracking of multiple targets in crowded environments. In *European Conference on Computer Vision*, pages 692–708. Springer, 2016. 2
- [6] M. Babae, Y. You, and G. Rigoll. Combined segmentation, reconstruction, and tracking of multiple targets in multi-view video sequences. *Computer Vision and Image Understanding*, 154:166–181, 2017. 2
- [7] N. Bansal, A. Blum, and S. Chawla. Correlation clustering. *Machine learning*, 56(1-3):89–113, 2004. 5
- [8] K. Bernardin and R. Stiefelhagen. Evaluating multiple object tracking performance: the clear mot metrics. *Journal on Image and Video Processing*, 2008:1, 2008. 8, 12
- [9] J. Chen, H. Sheng, Y. Zhang, and Z. Xiong. Enhancing detection model for multiple hypothesis tracking. In *Conf. on Computer Vision and Pattern Recognition Workshops*, pages 2143–2152, 2017. 8
- [10] W. Choi. Near-online multi-target tracking with aggregated local flow descriptor. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3029–3037, 2015. 8
- [11] S. Chopra and M. R. Rao. The partition problem. *Mathematical Programming*, 59(1-3):87–115, 1993. 4
- [12] Q. Chu, W. Ouyang, H. Li, X. Wang, B. Liu, and N. Yu. Online multi-object tracking using cnn-based single object tracker with spatial-temporal attention mechanism. In *2017 IEEE International Conference on Computer Vision (ICCV) (Oct 2017)*, pages 4846–4855, 2017. 2
- [13] A. Dehghan, S. Modiri Assari, and M. Shah. Gmmcp tracker: Globally optimal generalized maximum multi clique problem for multiple object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4091–4099, 2015. 2
- [14] E. D. Demaine, D. Emanuel, A. Fiat, and N. Immerlica. Correlation clustering in general weighted graphs. *Theoretical Computer Science*, 361(2-3):172–187, 2006. 5
- [15] M. M. Deza and M. Laurent. Geometry of cuts and metrics, volume 15 of algorithms and combinatorics, 1997. 4
- [16] K. Fang, Y. Xiang, X. Li, and S. Savarese. Recurrent autoregressive networks for online multi-object tracking. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 466–475. IEEE, 2018. 2
- [17] P. F. Felzenszwalb, R. B. Girshick, D. McAllester, and D. Ramanan. Object detection with discriminatively trained part-based models. *IEEE transactions on pattern analysis and machine intelligence*, 32(9):1627–1645, 2010. 12
- [18] Z. Fu, P. Feng, F. Angelini, J. Chambers, and S. Naqvi. Particle phd filter based multiple human tracking using online group-structured dictionary learning. *IEEE Access*, pages 14764–14778, 2018. 8
- [19] A. Geiger, M. Lauer, C. Wojek, C. Stiller, and R. Urtasun. 3d traffic scene understanding from movable platforms. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 36(5):1012–1025, 2014. 2
- [20] B. Harwood, G. Carneiro, I. Reid, and T. Drummond. Smart mining for deep metric learning. *space*, 9(13):22, 2017. 6
- [21] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *Computer Vision (ICCV), 2017 IEEE International Conference on*, pages 2980–2988. IEEE, 2017. 3
- [22] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 3
- [23] R. Henschel, L. Leal-Taixé, D. Cremers, and B. Rosenhahn. Fusion of head and full-body detectors for multi-object tracking. In *Computer Vision and Pattern Recognition Workshops (CVPRW)*, 2018. 1, 8
- [24] E. Hoffer and N. Ailon. Deep metric learning using triplet network. In *International Workshop on Similarity-Based Pattern Recognition*, pages 84–92. Springer, 2015. 5
- [25] M. Hofmann, D. Wolf, and G. Rigoll. Hypergraphs for joint multi-view reconstruction and multi-object tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3650–3657, 2013. 2
- [26] B. W. Kernighan and S. Lin. An efficient heuristic procedure for partitioning graphs. *The Bell system technical journal*, 49(2):291–307, 1970. 5
- [27] M. Keuper, E. Levinkov, N. Bonneel, G. Lavoué, T. Brox, and B. Andres. Efficient decomposition of image and mesh graphs by lifted multicuts. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1751–1759, 2015. 2, 5
- [28] M. Keuper, S. Tang, B. Andres, T. Brox, and B. Schiele. Motion segmentation & multiple object tracking by correlation co-clustering. *IEEE transactions on pattern analysis and machine intelligence*, 2018. 8
- [29] M. Keuper, S. Tang, Y. Zhongjie, B. Andres, T. Brox, and B. Schiele. A multi-cut formulation for joint segmentation and tracking of multiple objects. *arXiv preprint arXiv:1607.06317*, 2016. 8
- [30] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg. Multiple hypothesis tracking revisited. In *International Conference on Computer Vision*, pages 4696–4704, 2015. 8
- [31] C. Kim, F. Li, and J. M. Rehg. Multi-object tracking with neural gating using bilinear lstm. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 200–215, 2018. 8
- [32] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012. 2

- [33] L. Lan, X. Wang, S. Zhang, D. Tao, W. Gao, and T. S. Huang. Interacting tracklets for multi-object tracking. *IEEE Transactions on Image Processing*, 2018. 2
- [34] L. Leal-Taixé, C. Canton-Ferrer, and K. Schindler. Learning by tracking: Siamese cnn for robust target association. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 33–40, 2016. 1, 2
- [35] L. Leal-Taixé, A. Milan, I. Reid, S. Roth, and K. Schindler. Motchallenge 2015: Towards a benchmark for multi-target tracking. *arXiv preprint arXiv:1504.01942*, 2015. 2
- [36] E. Levinkov, J. Uhrig, S. Tang, M. Omran, E. Insafutdinov, A. Kirillov, C. Rother, T. Brox, B. Schiele, and B. Andres. Joint graph decomposition & node labeling: Problem, algorithms, applications. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 2, 8
- [37] H. Li, Y. Li, and F. Porikli. Deeptrack: Learning discriminative feature representations online for robust visual tracking. *IEEE Transactions on Image Processing*, 25(4):1834–1848, 2016. 2
- [38] C. Long, A. Haizhou, Z. Zijie, and S. Chong. Real-time multiple people tracking with deeply learned candidate selection and person re-identification. In *ICME*, 2018. 2, 8
- [39] C. Ma, C. Yang, F. Yang, Y. Zhuang, Z. Zhang, H. Jia, and X. Xie. Trajectory factory: Tracklet cleaving and reconnection by deep siamese bi-gru for multiple object tracking. *arXiv preprint arXiv:1804.04555*, 2018. 2, 8
- [40] 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. 12
- [41] A. Milan, S. H. Rezatofighi, A. Dick, I. Reid, and K. Schindler. Online multi-target tracking using recurrent neural networks. In *AAAI*, 2016. 2
- [42] W. Ouyang and X. Wang. Joint deep learning for pedestrian detection. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2056–2063, 2013. 2
- [43] H. Pirsiavash, D. Ramanan, and C. C. Fowlkes. Globally-optimal greedy algorithms for tracking a variable number of objects. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1201–1208. IEEE, 2011. 2
- [44] 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*, pages 91–99, 2015. 12
- [45] A. Sadeghian, A. Alahi, and S. Savarese. Tracking the untrackable: Learning to track multiple cues with long-term dependencies. *arXiv preprint arXiv:1701.01909*, 4(5):6, 2017. 2, 8
- [46] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 5
- [47] J. Son, M. Baek, M. Cho, and B. Han. Multi-object tracking with quadruplet convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5620–5629, 2017. 1, 2
- [48] J. Son, M. Baek, M. Cho, and B. Han. Multi-object tracking with quadruplet convolutional neural networks. In *Computer Vision and Pattern Recognition*, pages 5620–5629, 2017. 2
- [49] N. Srivastava, G. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929–1958, 2014. 7
- [50] S. Tang, B. Andres, M. Andriluka, and B. Schiele. Subgraph decomposition for multi-target tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5033–5041, 2015. 2
- [51] S. Tang, B. Andres, M. Andriluka, and B. Schiele. Multi-person tracking by multicut and deep matching. In *Proceedings of the European Conference on Computer Vision*, pages 100–111. Springer, 2016. 8
- [52] S. Tang, M. Andriluka, B. Andres, and B. Schiele. Multiple people tracking by lifted multicut and person reidentification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3539–3548, 2017. 1, 2, 8
- [53] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu. Learning fine-grained image similarity with deep ranking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1386–1393, 2014. 5, 6
- [54] N. Wang, S. Li, A. Gupta, and D.-Y. Yeung. Transferring rich feature hierarchies for robust visual tracking. *arXiv preprint arXiv:1501.04587*, 2015. 2
- [55] L. Wen, Z. Lei, S. Lyu, S. Z. Li, and M.-H. Yang. Exploiting hierarchical dense structures on hypergraphs for multi-object tracking. *IEEE transactions on pattern analysis and machine intelligence*, 38(10):1983–1996, 2016. 2
- [56] N. Wojke, A. Bewley, and D. Paulus. Simple online and real-time tracking with a deep association metric. *arXiv preprint arXiv:1703.07402*, 2017. 2
- [57] B. Yang and R. Nevatia. An online learned crf model for multi-target tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2034–2041. IEEE, 2012. 2
- [58] F. Yang, W. Choi, and Y. Lin. Exploit all the layers: Fast and accurate cnn object detector with scale dependent pooling and cascaded rejection classifiers. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2129–2137, 2016. 12
- [59] Y.-c. Yoon, A. Boragule, K. Yoon, and M. Jeon. Online multi-object tracking with historical appearance matching and scene adaptive detection filtering. *arXiv preprint arXiv:1805.10916*, 2018. 8
- [60] L. Zhang, Y. Li, and R. Nevatia. Global data association for multi-object tracking using network flows. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–8. IEEE, 2008. 2
- [61] X. Zhou, P. Jiang, Z. Wei, H. Dong, and F. Wang. Online multi-object tracking with structural invariance constraint. In *IEEE International Conference on British Machine Vision Conference (BMVC)*, 2018. 2
- [62] J. Zhu, H. Yang, N. Liu, M. Kim, W. Zhang, and M.-H. Yang. Online multi-object tracking with dual matching attention

Appendices

A. RNN Training Sample Generation

An abstract example of the RNN sample generation techniques employed in the framework is given in Fig. 8, where samples are being created which contain two tracklets, each with two detections, and a frame gap of 1 between them. Circles denote detections, and circles of the same color belong to the same person ID. The first tracklet in the sample is T_1 . Since a frame gap is required, the red detection in frame 3 is skipped, and tracklet T_2 is chosen as the second tracklet to create a positive sample.

For negative sample creation, there are three possibilities for choosing the second tracklet:

1. Selecting detections from another person’s track such that these detections are spatially close to the detections in T_2 . This is approximately achieved by searching for the detection whose bounding box center is closest to that of the red detection in frame 4, which turns out to be the blue detection. T_4 is therefore chosen as the second tracklet.
2. Selecting detections from another person’s track such that these detections lie in the same quadrant of the image. For this, the frame images are divided into four equally sized quadrants, as shown by the dotted lines. We then search for detections in frame 4 belonging to other persons that lie in the same quadrant as the red detection. This happens to be the detection in green, and therefore T_3 is chosen as the second tracklet. A random selection is made in case there are multiple candidates.
3. Selecting detections from another person’s track such that these detections lie in another quadrant of the image. The procedure for this is the same as above, except that a detection in any of the other quadrants is chosen. For this scenario, tracklet T_5 is a suitable candidate.

B. Edge Creation

Since the proposed framework employs a computationally expensive neural network to compute tracklet similarities, and moreover, involves multiple iterations of graph clustering, creating edges between all combinations of tracklet pairs within the allowable frame gap range results

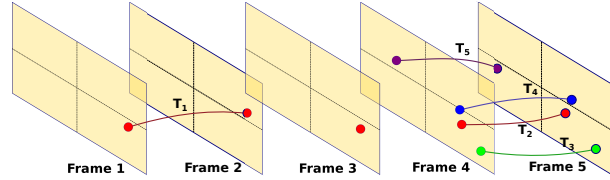


Figure 8: Example of training sample creation for RNN.

in a very dense graph which requires a long time to evaluate. To mitigate this problem, we employ an edge creation method that avoids creating irrelevant edges. The underlying intuition is that it is very unlikely that two tracklets belong to the same person if they are spatially far apart, but temporally close. To formulate this mathematically, we compute statistics for the average, per frame bounding box movement using the provided ground-truth data in the training datasets. A separate set of statistics is computed for static and moving camera videos. Moreover, these statistics are normalized by the image dimensions, framerate, and the bounding box dimensions (larger detections are normally closer to the camera, and can therefore be expected to experience larger movements).

When applying the framework to test datasets, the pre-computed average statistics are first inflated as a safety measure, and then used to define a feasible radius around the detections of a tracklet. Edges are created only with those tracklets whose detections lie within this radius. Naturally, the radius is scaled according to the frame gap between the detections. This process is illustrated in Fig. 9 for the simple case where each tracklet contains a single detection. The picture sequence shows how the acceptable radius increases in size as the frame gap between the detections being considered increases. To clearly show the increasing radius, each image in the sequence occurs 10 frames after its predecessor.

It is worth pointing out that even though these statistics are heuristically computed parameters that may not generalize to all video sequences, this method of edge creation only serves to reduce computational time, and has negligible impact on the accuracy of the framework’s output. This is because the edges discarded in this manner are trivial cases which the RNN can easily detect as being dissimilar.

C. Graph Complexity and Convergence

In Fig. 11, the number of vertices and edges in the graph for each iteration of the clustering process, and for each test dataset in the MOT16 challenge are illustrated. The key observations from these results are:

- The number of vertices is, on average, reduced by approximately 90% after the first iteration, suggesting that a significant part of the clustering is already com-

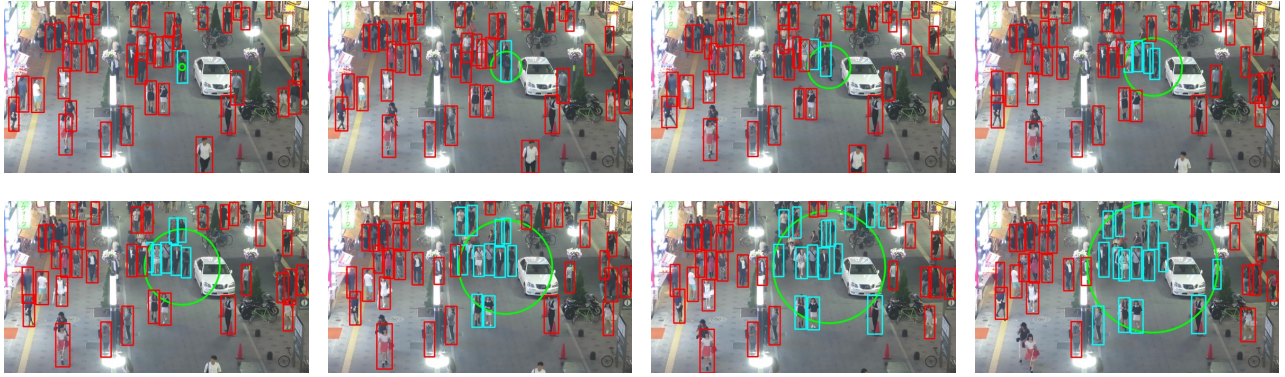


Figure 9: Visualization of the edge creation criteria. Here, we wish to find feasible edges for the detection in cyan in the upper left image. The green circle denotes the acceptable radius. The other detections in cyan are those which are accepted as edge connections, whereas those in red are not. Ordering is from left to right, and then top to bottom.

plete. This behavior is encouraging because it means that the complexity of the graph is greatly reduced after just one iteration.

- The number of edges and vertices are roughly of the same order. This shows the effectiveness of the edge creation scheme described earlier.
- The number of iterations required for convergence under the currently employed clustering scheme is 8-12. This is despite the fact that the test dataset videos were captured in varying environments, and have different numbers of detections and frames, and different detection densities (average number of detections per frame).
- Recall that the criteria for determining the maximum allowable frame gap between tracklets was relaxed in two steps: the first relaxation occurs in the fourth iteration, when the maximum allowed frame gap is increased from 4, to 4 times the tracklet length. This is reflected by the spike in the number of edges created in the fourth iteration. The second relaxation comes when the algorithm initially converges, after which edges are allowed to span 6 times the tracklet length. Here, it is again observable that the number of edges increases. Moreover, the number of vertices usually decreases after this relaxation, even though the algorithm had converged under the previous criteria. The iteration number of both relaxations is marked in magenta colored arrows on the graphs.

D. Timing Analysis

A common concern with any framework that employs graph optimization is its scalability. Fortunately, we observed that the total computational time (including inference and clustering) required to process each dataset is

strongly correlated with the number of edges created in the first iteration. In Fig. 10, a scatter plot is drawn of the processing time required for each MOT17 test dataset against the number of edges created in the first iteration. The line of best fit between all the points is also plotted in magenta. It can be seen that despite the varying nature of the datasets, there is a strong linear relation between the two parameters. This trend offers strong empirical support for the scalability of our framework.

All results were obtained on a desktop system with an Intel Xeon E5-1620 CPU running at 3.5GHz with 16GB RAM, and an Nvidia GTX TITAN X GPU.

E. MOT Challenge Results

The performance metrics for each video sequence in the MOT16 [40] test dataset are provided in Table 4. In Fig. 12, screenshots of annotated video sequences of the MOT16 test dataset are given. Three screenshots from each of the seven videos are given in each row of the tiled figure. The results provided in the two tables, as well as the full video sequences are available online for both MOT16³ and MOT17 challenges⁴.

As mentioned in the main text, the MOT17 challenge contains the same video sequences as MOT16, but with a different ground truth, and with three sets of detections which are produced by different publicly available object detectors: DPM [17], SDP [58] and FRCNN [44]. In Fig. 13, a bar plot shows how the MOT Accuracy score (as defined by MOT Clear metrics [8]) varies for each of the seven video sequences depending on which detector the detections came from.

³https://motchallenge.net/tracker/HDTR_16

⁴https://motchallenge.net/tracker/HDTR_17

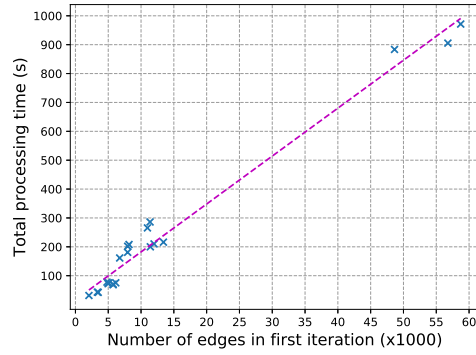


Figure 10: Plot of processing time against number of edges in the first iteration for all MOT17 test datasets.

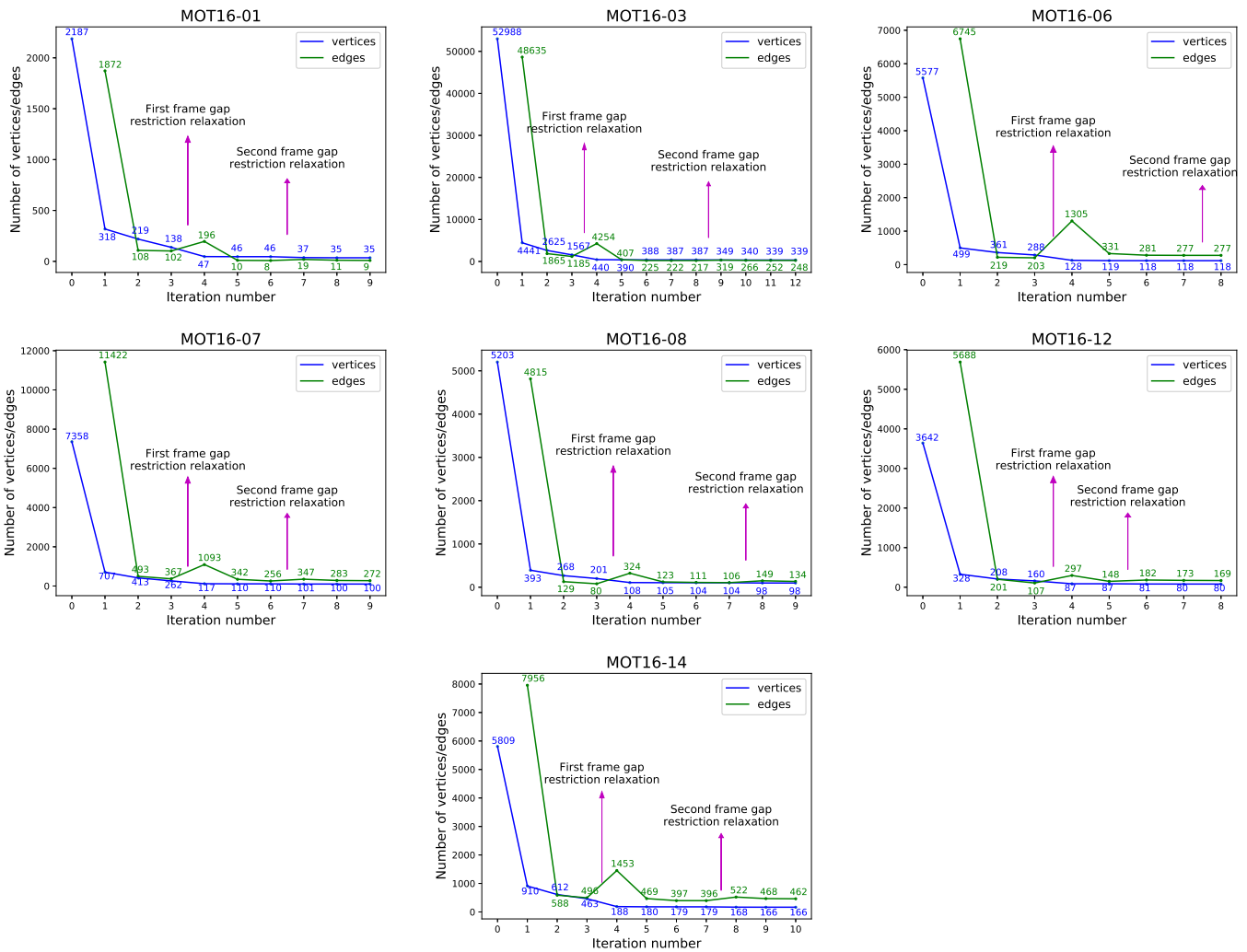


Figure 11: Number of vertices and edges in the graph at each clustering iteration for all MOT16 test datasets.



Figure 12: Screenshots of annotated video sequences from the MOT16 test datasets. Different colors correspond to different person IDs.

Tracker	MOTA(%) \uparrow	IDF1(%) \uparrow	MOTP(%) \uparrow	FAF \downarrow	MT(%) \uparrow	ML(%) \downarrow	FP \downarrow	FN \downarrow	ID Sw. \downarrow	Frag \downarrow
MOT16-01	49.5	38.2	80.6	0.1	30.4	30.4	40	3178	10	15
MOT16-03	62.9	52.1	81.0	0.8	31.1	16.9	1175	37473	185	262
MOT16-06	43.2	35.2	81.2	1.3	27.1	42.1	1565	4898	88	130
MOT16-07	46.9	42.6	79.5	1.6	20.4	24.1	788	7778	105	143
MOT16-08	37.7	37.2	82.1	0.3	17.5	38.1	163	10219	45	73
MOT16-12	42.4	49.9	80.4	0.3	15.1	48.8	247	4482	53	46
MOT16-14	34.0	33.2	78.2	1.0	7.9	47.0	736	11325	132	164

Table 4: Tracking results for each video sequence in the MOT16 test dataset with public detections. \uparrow and \downarrow represent higher is better and lower is better, respectively. The values in bold and blue represent the best and second best performances, respectively.

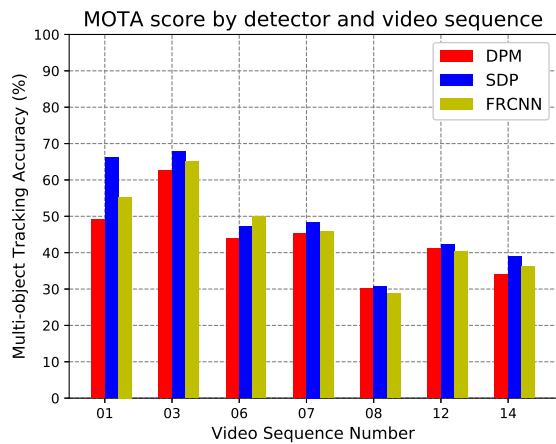


Figure 13: Example of training sample creation for RNN.