# FANTrack: 3D Multi-Object Tracking with Feature Association Network

Erkan Baser[1], Venkateshwaran Balasubramanian[2] [*], Prarthana Bhattacharyya[3] [*], Krzysztof Czarnecki[3]

*Abstract*— We propose a data-driven approach to online multi-object tracking (MOT) that uses a convolutional neural network (CNN) for data association in a tracking-by-detection framework. The problem of multi-target tracking aims to assign noisy detections to a-priori unknown and time-varying number of tracked objects across a sequence of frames. A majority of the existing solutions focus on either tediously designing cost functions or formulating the task of data association as a complex optimization problem that can be solved effectively. Instead, we exploit the power of deep learning to formulate the data association problem as inference in a CNN. To this end, we propose to learn a similarity function that combines cues from both image and spatial features of objects. Our solution learns to perform global assignments in 3D purely from data, handles noisy detections and a varying number of targets, and is easy to train. We evaluate our approach on the challenging KITTI dataset and show competitive results. Our code is available at **https://git.uwaterloo.ca/wise-lab/fantrack**.

## I. INTRODUCTION

Multi-object tracking (MOT) is a critical problem in computer vision and has received great attention due to its widespread use in applications such as autonomous driving, robot navigation, and activity recognition. It is the problem of finding the optimal set of trajectories of objects of interest over a sequence of consecutive frames. Most of the successful computer vision approaches to MOT have focused on the *tracking-by-detection* principle [1], [2]. This paradigm allows the problem to be divided into two steps. First, an object detector is used to identify the potential locations of objects in the form of bounding boxes, and then a discrete combinatorial problem is solved to link these noisy detections over time to form trajectories. Despite decades of research, the status quo of tracking is far from reaching human accuracy. Current challenges to the problem include a varying and a-priori unknown number of targets; incorrect and missing detections; changing appearances of targets due
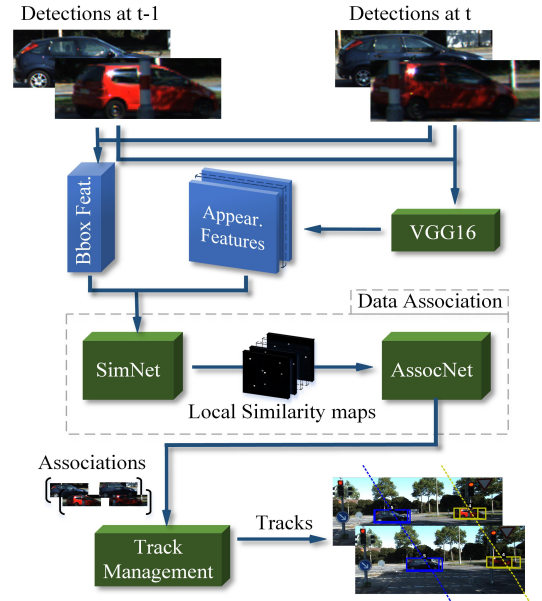


Fig. 1. Overall architecture of the proposed approach

to sensor motion, illumination, and angle of view; frequent occlusions, and abrupt changes in motion.

The linking step called *data association* is arguably the most difficult component of MOT. Traditional *batch* methods usually formulate MOT as a *global* optimization problem, with the assumption that detections from all future frames are available, and solve it by mapping it to a graph based min-cost flow algorithm [3], [4]. *Online* Markovian formulations of MOT on the other hand often employ greedy or bipartite graph matching methods like the Hungarian algorithm to solve the assignment problem [5]–[7]. *Online* approaches are well suited to real-time applications such as tracking road-traffic participants. The success of the final associations is also dependent on the similarity functions used to match the targets and detections. Traditionally cost functions have been handcrafted with representations based on color histograms, bounding box position, and linear motion models [8], [9], but have failed to generalize across tasks and for complex tracking scenarios. Recently, deep neural network architectures have shown superior performance in many vision based tasks. Milan et al. proposed the first end-to-end formulation for MOT, using a recurrent neural network (RNN) to solve the assignment problem for each target independently based on Euclidean cost [10]. However, the use of convolutional neural networks (CNNs), which are easier to train than RNNs, in order to solve the association problem while also learning

[1] Erkan Baser was affiliated to Waterloo Intelligent Systems Engineering Lab, University of Waterloo, 200 University Ave W, Waterloo, ON N2L 3G1. erkanbaser@gmail.com

[2] Venkateshwaran Balasubramanian is with the David R. Cheriton School of Computer Science, University of Waterloo, 200 University Ave W, Waterloo, ON N2L 3G1. venk.b@outlook.com

[3] Prarthana Bhattacharyya and Krzysztof Czarnecki are with the Department of Electrical and Computer Engineering, University of Waterloo, 200 University Ave W, Waterloo, ON N2L 3G1.

{p6bhatta,k2czarne}@uwaterloo.ca

[*] denotes equal contribution

the cost function has not yet been investigated.

In this paper, we propose an *online* MOT formulation that casts the assignment problem as inference in a CNN. We present a two-step learning based approach (see Fig. 1). The first step learns a similarity function that takes advantage of both visual and 3D bounding box data to yield robust matching costs. The second step trains a CNN to predict discrete target assignments from the computed pairwise similarities. The benefit of our proposal is that it is easy to train, takes care of a varying number of targets and noisy detections, and provides a simple way to consider all the targets while making associations. We empirically demonstrate on the KITTI tracking dataset [11] that: (i) Our approach can solve the multi-target association problem by performing inference using CNNs. (ii) It can integrate image based appearance and 3D bounding box features to get a discriminative as well as generalized feature representation, thereby learning a robust cost function for association. (iii) We show competitive qualitative and quantitative 3D tracking results compared to the state of the art.

## II. LITERATURE REVIEW

### A. Data Association in MOT

Classical approaches solve the data association problem by considering multiple hypotheses for an assignment (MHT) [12], or by jointly considering all possible assignment hypotheses (JPDA) [13]. These formulations prove to be very computationally intensive, however.

Many recent works process sequences in *batch* mode, using a graph-based representation with detections as nodes and possible assignments as edges. The optimization is then cast as a linear program solved to (near) global optimality with relaxation, min-cost or shortest path algorithms [14]–[16]. More complex optimization schemes include MCMC [17] and discrete-continuous settings [18]. However, *global* optimization formulations are unsuited to real-time applications like autonomous navigation.

*Online* methods estimate the current state using the information only from the past frames and the current one. Commonly used state-estimators include the Kalman filter [19] for linear motion and particle filters [20] for multimodal posteriors. The two-frame association problem is often solved using a greedy or Hungarian algorithm [6]. Approaches based on local associations tend to be susceptible to track fragmentation and noisy detections, however.

*Deep learning* has achieved state-of-the-art results in perception tasks like image classification, segmentation, and single object tracking. Milan et al. proposed the first fully end-to-end multi-object tracking method based on deep learning. The method predicts the assignment of each target, one at a time, using an RNN [10]. In contrast, our approach feeds all detections and their learned similarity scores at once into a CNN to predict the assignments. Our model is easier to optimize than an RNN, handles noisy detections and a varying number of targets, and considers all targets at once when performing assignments.

### B. Measuring Similarity

Tracking algorithms have used distance functions such as Euclidean [21] and Mahalanobis distance [22] as matching costs for data association. Other similarity measures include color-based appearance features [23], SIFT-like features [24], and linear and non-linear motion models and their various weighted combinations [25]. These tediously hand-crafted features fail to generalize across complex scenarios and backgrounds, however.

Recent works explore learning pairwise costs using deep structured SVM [3], CNNs [26], and RNNs [27]. For CNNs, similarity learning often exploits Siamese networks. Leal-Taixe et al. [28] and Frossard et al. [29] use them to learn descriptors for matching with multi-modal inputs. While we also use Siamese networks to learn *generalized* and *discriminative* features from 3D object configurations and visual information conditioned on similarity, we adapt our objective function to use the cosine-similarity metric with hard-mining which has a positive impact on convergence.

## III. OUR APPROACH

Our proposed framework is based on tracking by detection paradigm. Our problem setup assumes at any time instant $t$ we have $N$ number of targets, $M$ number of detections and track labels for every $i^{th}$ track. We use AVOD [30] as our 3D object detector since it achieves state-of-the-art results on KITTI and is open-source, but in principle, any other 3D object detector could be used. The motivation for building FANTrack is to leverage the power of Siamese networks to model the similarities between targets and detections, CNNs to solve the data association problem in MOT, and an online track management module to update, initialize and prune tracks. We describe these modules in the following sections.

### A. Similarity Network

Figure 2 gives an overview of our proposed similarity network *SimNet*. The network has two *input* pairs with each pair corresponding to target and detection data, and consists of 3D bounding box parameters ($1 \times 7$ dimensional vector) and image convolutional features ($7 \times 7 \times 320$ dimensional vector). The *output* from *SimNet* is a set of $N_{max}$ number of maps for each existing target corresponding to a *local* $5m \times 5m$ region around the target and of $0.5m$ resolution. These output maps contain the similarity scores in each target's local neighbourhood with respect to all detections at a particular time step. The *SimNet* output is subsequently used for data association.

*Functionally*, *SimNet* computes a similarity score for every detection and target pair. It has two branches: a *bounding box branch* and an *appearance branch*, each of which uses a trainable Siamese network to learn object representations conditioned on whether two objects are similar or not. The outputs of these branches are vector representations of targets and detections. Their respective contribution towards the final similarity score computation is weighted using the *importance branch*. Finally, cosine-similarities of each target-detection vector representation are computed and the
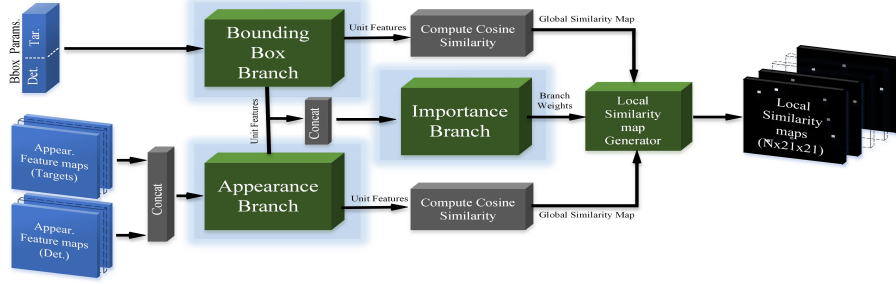
Fig. 2. Architecture of the proposed Siamese network for similarity learning. The branches highlighted in blue have trainable parameters.

scalars are mapped to their corresponding positions on the above-mentioned set of local maps.

We describe our formulation of *SimNet* in the remainder of this sub-section.

*1) Bounding Box Branch:* The bounding box branch outputs a discriminative, robust vector representation for 3D object configurations of targets and detections, conditioned on whether they are similar or not. We train a Siamese network with *stacked* input pairs of target and detection 3D bounding boxes for this purpose. The 3D bounding boxes are defined by their centroids $(x, y, z)$, axis-aligned dimensions $(l, w, h)$, and rotation around the $z$-axis $(\theta_z)$ in the ego-car's IMU/GPS coordinates. To prevent learning variations induced due to ego-motion, detection centroids are converted to coordinates at a common time-step using GPS data.

*Architecture:* The input to this branch is a $(N + M) \times 1 \times 7$ tensor where the third dimension consists of the 7 bounding box parameters defined above. The inputs are fed to a convolutional layer with 256 $1 \times 1$ filters to capture complex interactions across the 7 channels by pooling the parameters of each targets and detections independently [31]. These object-independent features are then fed into two fully-connected layers with 512 neurons, with dropout regularization. We apply L2 normalization on the output features, and henceforth refer to the result as *unit* features. Finally, the unit features of dimensions $(N + M) \times 512$ are sliced along the first dimension into target features and detection features using their respective counts (see Fig 3). These are used to compute the bounding box cosine similarities as described in subsection *A.4*. We use batch normalization and leaky-ReLU across all layers.

*2) Appearance Branch:* The appearance branch outputs a robust and invariant vector representation for 2D visual cues of targets and detections conditioned on whether they belong to similar or dissimilar objects. We train another stacked Siamese network for this purpose. As its input, we concatenate convolutional features of targets and detections obtained from AVOD's ([30]) image feature extractors. Specifically, we use the second layer's convolutional features and the interpolated fourth layer's convolutional features. This is because the low-level features are local and more discriminative whereas high-level features are abstract, and are more invariant to appearance changes [32], [33].
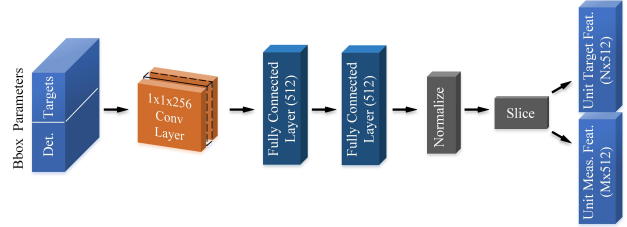


Fig. 3. Detailed architecture of the bounding box branch. Inputs are the concatenated bounding boxes of targets and detections. Outputs are sliced unit feature vectors.

*Architecture:* The input to this branch is a $(N + M) \times 7 \times 7 \times 320$ convolutional feature. The architecture of the branch is shown in Fig. 4. First, we apply 256 $3 \times 3$ convolutions to obtain promising features for similarity learning by preserving the spatial size of the input. Before flattening the feature maps for the fully-connected layers with 512 neurons, the Global Average Pooling (GAP) [34] layer extracts one abstract feature from each feature map. Similar to the bounding box branch, L2 normalization yields a vector of dimension $(N + M) \times 512$. As in the case of the bounding box branch, the $(N + M) \times 512$ features are sliced along the *first* dimension to obtain appearance features of detections and targets to compute the appearance cosine similarities.

*3) Importance Branch:* The aim of this branch is to determine the *relative importance* of the bounding box and appearance features in the computation of the final cosine similarity score (see Fig. 5).

*Architecture:* The inputs to this branch are the unit features from the other two branches. First, the vector representation of an object obtained from the appearance and bounding box branches is concatenated to form a single vector (dimension 1024). Then, a fully-connected layer with two neurons, ReLU activation, and a softmax layer that computes two scalars indicating importance weights (probabilities) of the two branches, for each target and detection. The importance weights obtained ($\omega_{bbox}$ and $\omega_{appear}$) are normalized to sum up to unity.

*4) Similarity Maps:* A similarity map (see Fig. 6) is computed for every target (for $N$ targets we have $N$ similarity maps) and contains its cosine similarity scores with all the detections within a 2D region of interest ($[-40, 40] \times [0, 80]$
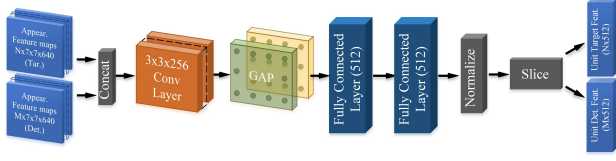
Fig. 4. Detailed architecture of the appearance branch. Inputs are concatenated appearance feature maps of targets and detections. Outputs are sliced unit feature vectors.

m) in the ego-car's IMU/GPS coordinates. This map is referred to as the *global* similarity map. To compute the similarity scores, we perform the following:

*i*) Each global map is split into grids with $0.5$ m resolution.
*ii*) The detection appearance and bounding box features are positioned into the appropriate grid locations based on their location of detection.
*iii*) The target appearance and bounding box features are used as kernels to compute the similarity scores by the convolution with strides equal to $1$. The computed scores correspond to cosine similarities as the features are normalized to unit vectors by the network branches.

Local similarity maps are then obtained from the global similarity maps for each target by cropping them around the target's local $5\text{m} \times 5\text{m}$ region corresponding to $10 \times 10$ cells. *SimNet* thus finally outputs $N \times 21 \times 21$ local similarity maps to be used for data association.

*5) Loss Function:* To learn the trainable parameters of the appearance branch, bounding box branch, and importance branch, we use the *weighted cosine distance* given by:

$$L(\Theta_1) = \frac{1}{N^+} \sum_{i=1}^{N} w_{skew}^{(i)} \times w_{cost}^{(i)} \times \left( 1 - y^{(i)} \times \hat{y}^{(i)}(\Theta_1) \right) \quad (1)$$

where $\Theta_1$ is the network parameters, $N^+$ is the number of examples with nonzero weights, $y^{(i)}$ denotes the ground truth value of the $i^{th}$ example, i.e., $y^{(i)} \in \{-1, 1\}$. $\hat{y}^{(i)}$ is the estimated cosine similarity score computed using the cosine similarities from the two branches and their normalized importance weights as follows:

$$\hat{y}^{(i)}(\Theta_1) = \omega_{bbox}(\Theta_1)^{(i)} \times \hat{y}_{bbox}^{(i)}(\Theta_1) + \omega_{appear}^{(i)}(\Theta_1) \times \hat{y}_{appear}^{(i)}(\Theta_1) \quad (2)$$

$w_{skew}^{(i)}$ is the weight used to remove the imbalance of negative examples in the training dataset. $w_{cost}^{(i)}$ scales the loss function according to how easy or hard it is to distinguish between each pair of examples so that the training can revolve around a sparse set of the selected hard examples [35].

*6) Creating training examples for SimNet:* In this section, we describe creating positive and negative pairs of examples by augmentation from the KITTI training set to train the similarity network. A new bounding box proposal is a positive pair if its intersection over union (IoU) with its ground truth on images exceeds $0.8$. The selected IoU thresholds should
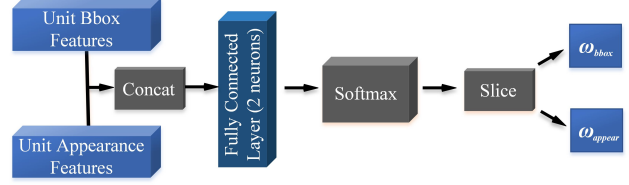


Fig. 5. Detailed architecture of the importance branch. The inputs are the unit bounding box and appearance features of both targets and detections. Outputs are branch weights for bounding box and appearance branches. These weights can be further sliced for targets and detections separately.
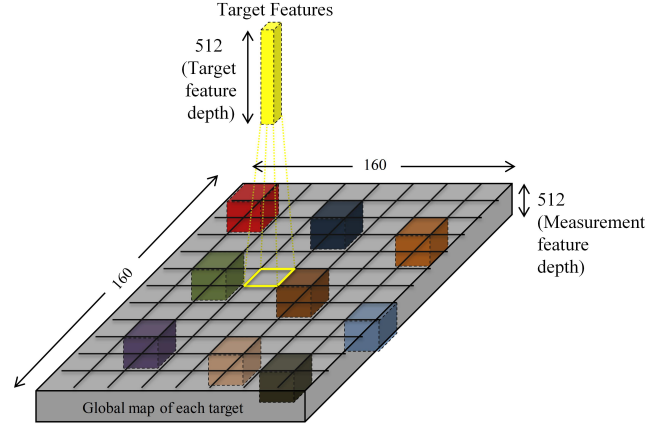


Fig. 6. Construction of global similarity map, one for each target. The target feature vector (yellow solid) is convolved with those of the detections to compute the similarity scores. Locations that do not include a detection feature vector are filled by zero vectors.

be at least greater than those used for the final detections in AVOD $(0.65)$ [30]. In addition, the diversity among bounding box proposals for each object is maintained by rejecting a new proposal whose IoUs with existing ones are greater than $0.95$.

*B. Data Association Network*

Fig. 7 introduces our proposed data association network, referred to as *AssocNet*. The purpose of this network is to associate targets to the detections. The *input* to the network is the set of local similarity maps of dimension $N \times 21 \times 21$ obtained from *SimNet*, containing cosine similarity scores for probable target-detection pairs. The *output* from the network is the target-to-detection association probabilities for each existing target.

We first describe how our framework handles noisy detections and a varying number of targets. In order to deal with varying number of inputs, we create $N_{max} - N$ extra channels with dummy maps, where $N_{max}$ denotes the maximum number of targets that can be tracked. The dummy maps are a matrix of zeros - a reasonable representation since zero inputs don't have any impact on the output of the convolutional layers. We deal with missed detections by introducing an extra cell for each of the $N_{max}$ targets to account for *spurious detections* and concatenate it to their map of logits as described in the architecture below.
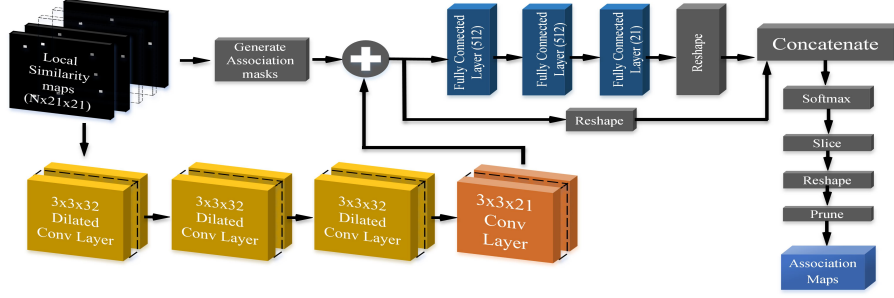
Fig. 7. The architecture of the proposed association network. The inputs are local similarity maps from the proposed Siamese network. The outputs are the association maps which provide target-to-detection association and detection probabilities.

*Architecture:* The main building blocks of the *AssocNet* are convolutional, dilated convolutional (d-Conv), and fully-connected layers, with batch-normalization and leaky-ReLU activation used in all the layers. We take advantage of the increased receptive field of dilated convolutions [36] to compensate for the sparsity of local similarity maps.

We now discuss the flow of information through *AssocNet*. The network processes the input using three dilated convolutional (d-Conv) layers with dilation factors of $2, 4$, and $6$ respectively. The neighbouring fields have slightly overlapping fields of view due to increased dilation size [37]. The convolutional layer enables interactions between these neighbouring units which effectively results in considering all the detections simultaneously while making assignments. Thus to aggregate information, we employ a $3 \times 3$ convolutional layer at the end to compute the maps of logits (the vector of non-normalized predictions).

*AssocNet* is to be trained to predict assignment probabilities between a target and its probable detections. Since the locations of probable detections are known in each local similarity map, there is no need to train *AssocNet* to predict assignment probabilities of other locations as zero. To implement this idea, we generate association masks for each local similarity map. In the association masks, cells of probable detections are set to zero, while the other cells are set to a minimum negative number. Then the association masks are added to the map of logits obtained from the convolutional layer with $3 \times 3 \times 21$ filters (see Fig.7). This maintains the values of the logits computed for probable detections, but makes other logits insignificant for further computation.

After masking the maps of logits, *Assocnet* is split into two branches. One branch consisting of fully-connected layers predicts the $N_{max}$ logit values of spurious detections. The other branch reshapes the logit map into 1D vectors to concatenate logits of spurious detections with those of probable detections. This results in a $N_{max} \times (21 \times 21 + 1)$ tensor. The softmax then computes the association probabilities for each target, which are is our required output.

The association probabilities are sliced and reshaped in order to obtain 2D association maps. The probabilities computed for spurious detections are missed-detection probabilities. Finally, we get rid of the association maps corresponding

to the $N_{max} - N$ dummy channels.

*1) Loss Function:* Training *AssocNet* can be considered as training a classification problem in which labels are association maps showing the true data association for each existing target. To train the data association network we use a multi-task loss function given by:

$$L(\Theta) = l(\Theta)_{assoc} + l(\Theta)_{reg} \tag{3}$$

where $\Theta$ is the set parameters of the association network, $l(\Theta)_{reg}$ is the regularization loss. $l(\Theta)_{assoc}$ is the binary cross-entropy computed for the association maps as follows:

$$q_{vec} = q_{assoc}^{(t)}(i,j) \times \log \left( \min \left( \hat{q}_{assoc}^{(t)}(i,j;\Theta) + 0.01, 1 \right) \right)$$
$$p_{vec} = p_{assoc}^{(t)}(i,j) \times \log \left( \min \left( \hat{p}_{assoc}^{(t)}(i,j;\Theta) + 0.01, 1 \right) \right)$$
$$l(\Theta)_{assoc} = \sum_{t=1}^{N} \sum_{i,j=1}^{21+1} (-q_{vec}) + (-p_{vec})$$
$$\tag{4}$$

where $q_{assoc}^{(t)}(i,j) = 1 - p_{assoc}^{(t)}(i,j)$ and $0.01$ is the margin used to ignore negligible errors in the predicted probabilities $\hat{p}_{assoc}^{(t)}(i,j;\Theta)$.

### C. Track Management

The track management module takes care of state estimation, initiation, update, and termination of tracks. We use a Kalman filter for motion prediction and state estimation. We initiate, update and prune tracks with a Bayesian estimation model as specified in [38] with a probability of existence $P_e$. Our complete tracking algorithm is described in Algorithm 1.

## IV. EXPERIMENTS

In this section, we describe the dataset, training parameters, and experimental evaluation results for the tracker built using our proposed data association networks.

### A. Dataset

We used the KITTI Tracking benchmark dataset for training and evaluation of our approach. The KITTI Tracking dataset consists of 21 training sequences and 29 test sequences. As the training sequences have different levels

**Algorithm 1:** Tracker Algorithm

---

$\theta_{ex} = 0.40$ ;         `// Existence threshold`
**while** *true* **do**
    Get Detections $m^\tau$ at time $\tau$
    **if** $\tau = 0$ **then**
        **foreach** $m_i^0$ **do**
            Create new track i
        **end**
    **else**
        Perform Kalman Filter Prediction;
        Predict $P_e \; \forall i \in T$;
        $DataAssociation$ for $t^k$ and $m^k$;
        Perform Kalman Filter Update;
        Update $P_e \; \forall i \in T$;
        $\forall (t_i^k, m_j^k)$ Update track $i$ with $m_j^k$
        $\forall (t_i^k, None)$ Propagate predicted $t_i^k$ to
          $\tau = k + 1$
        $\forall (None, m_j^k)$ Create a new track;

        $\forall i \in T$ **if** $P_{e_i}^k < \theta_{ex}$ **then**
           $Prune \quad i$
        **end**
    **end**
**end**

---

of difficulty, occlusion, and clutter, we split the 20% of every training sequence for validation. This way, training and validation datasets are not skewed. For training SimNet, we construct a training dataset from the training sequences by generating positive and negative examples in consecutive frames using ground truth information. Geometric transformations (translation, rotation, and scaling) are applied to the ground-truth bounding box parameters to model partial occlusion and detector noise. This gives a large training set in which the ratio of negatives to positives is approximately $18 : 25$. We trained the object detector using a combined dataset consisting of the KITTI 3D object detection dataset and the 80% split of the KITTI training dataset mentioned earlier after pre-training on a synthetic dataset [44].

### B. Training Parameters

*1) Similarity Network:* SimNet is trained with mini-batches of size 128. Each mini-batch consists of the spatial indices of detections in the global map, the number of targets ($N$), target centroids in $x$-$y$ coordinates, target and detection appearance features, their bounding box parameters, and the labels of each example. To optimize the loss function (1) we used Adam optimizer and exponentially-decaying learning rate [45]. The learning rate is initially set to $1e - 5$ and then decreased every 100 epochs with a base of 0.95.

*2) AssocNet:* To optimize the loss function in (3) we used Adam optimizer and exponentially-decaying learning rate. The learning rate was initially set to $1e-6$ and then decreased every 20 epochs with a base of 0.95.

### C. Evaluation Metrics

We use the popular CLEAR MOT metrics [46] for evaluating our tracker. Multiple Object Tracking Accuracy (MOTA) gives us an estimate of the tracker's overall performance. However, this is dependent on the performance of the object detector. Hence, we also look at tracking specific metrics like Mostly Tracked (MT), Mostly Lost (ML), ID Switches (IDS) and fragmentation (FRAG), which evaluate the efficiency of the tracker in assigning the right IDs with reduced switches or fragmentation in the tracks.

### D. Ablation Study

We do an ablation study to evaluate the components in our approach by comparing them with traditional approaches. Firstly, we study the impact of the similarity network. In Table II, Euclidean and Manhattan denote the baseline distances modeled with the 3D position estimates. Bhattacharyya and ChiSquare metrics are built from the image histograms of the cropped targets and detections to study the image-only configuration. SimNet and AssocNet denote our similarity and Association networks respectively. From Table II, we could infer that conventional similarity approaches were not able to achieve comparable accuracy (MOTA) as the features involved in the computation of the similarity scores were not robust. We also study the impact of our association network by replacing it with a baseline Hungarian approach. Again, we could observe that the baseline approaches like Hungarian couldn't fare better than ours.

### E. Qualitative Evaluation

We perform a qualitative evaluation by running our tracker on the KITTI tracking validation and testing sequences. We analyze different scenarios including occlusions, clutter, parked vehicles and false negatives from the detector. Fig. 8 shows an example from sequence 0 in the test set. Different tracks representing the vehicles are color coded and the track IDs are displayed for reference. The tracker is able to perform well in spite of the clutter due to the closely parked cars. In Fig. 9 we see an example from test sequence 17 in which the false negative by the detector is overcome with the help of the prediction of the tracker. These examples show the robustness of the tracker and its ability to perform better even with an average object detector. There were also some cases where the data association fails and as a result ID switching and fragmentation happen. In Fig. 10 the track 38 was previously assigned to a nearby car but after an occlusion in the detection ID switching happens. This could be due to the low-lit conditions of the two cars.

### F. Benchmark Results

We evaluate our approach on the test sequences on the KITTI evaluation server for the 'Car' class. The results are presented in Table I. Due to the challenging nature of online tracking approach and to do a fair comparison, we only consider published online tracking approaches for our comparison. We achieve competitive results with respect to the state of the art in online tracking with improved

TABLE I

RESULTS ON KITTI TEST SET FOR 'CAR' CLASS

| Method | MOTA ↑ | MOTP ↑ | MT ↑ | ML ↓ | IDS ↓ | FRAG ↓ |
|---|---|---|---|---|---|---|
| MOTBeyondPixels [39] | **84.24** % | **85.73** % | **73.23** % | **2.77** % | 468 | 944 |
| JCSTD [40] | 80.57 % | 81.81 % | 56.77 % | 7.38 % | **61** | 643 |
| 3D-CNN/PMBM [41] | 80.39 % | 81.26 % | 62.77 % | 6.15 % | 121 | 613 |
| extraCK [42] | 79.99 % | 82.46 % | 62.15 % | 5.54 % | 343 | 938 |
| MDP [43] | 76.59 % | 82.10 % | 52.15 % | 13.38 % | 130 | **387** |
| *FANTrack (Ours)* | 77.72 % | 82.32 % | 62.61 % | 8.76 % | 150 | 812 |

TABLE II

ABLATION STUDY ON KITTI VALIDATION SET FOR 'CAR' CLASS

| Method | MOTA ↑ | MOTP ↑ | MT ↑ | PT ↑ | ML ↓ | IDS ↓ | FRAG ↓ |
|---|---|---|---|---|---|---|---|
| Euclidean+*AssocNet* | 56.16 % | 84.84 % | 72.22 % | 18.51 % | 9.25 % | 269 | 320 |
| Manhattan+*AssocNet* | 56.75 % | 84.83 % | **73.14 %** | 17.59 % | 9.25 % | 265 | 319 |
| Bhattacharyya+*AssocNet* | 56.69 % | 84.81 % | 72.22 % | 18.51 % | 9.25 % | 256 | 307 |
| ChiSquare+*AssocNet* | 57.17 % | 84.81 % | **73.14 %** | 18.51 % | **8.33 %** | 262 | 311 |
| *SimNet*+Hungarian | 74.59 % | **84.92 %** | 65.74 % | **23.14 %** | 11.11 % | 26 | 93 |
| *SimNet*+*AssocNet* | **76.52 %** | 84.81 % | **73.14 %** | 17.59 % | 9.25 % | **1** | **54** |

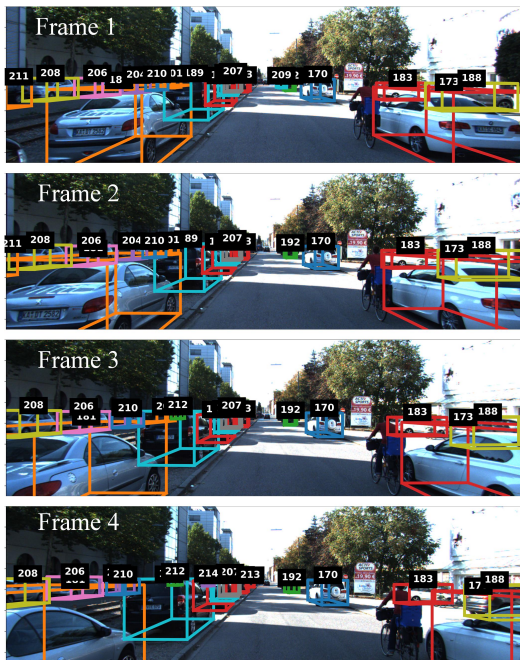(↑ denotes higher values are better. ↓ denotes lower values are better)



Fig. 8. Qualitative Evaluation - An example from video 14 in test set where the tracker performs well in a cluttered scene with parked cars.
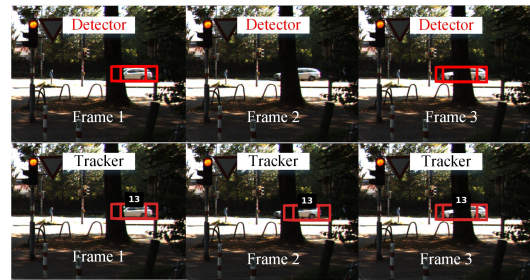


Fig. 9. Qualitative Evaluation - In this example (video 17 in test set) the detection was missed by the detector and reappears in the next frame. But the tracker was able to successfully maintain the track
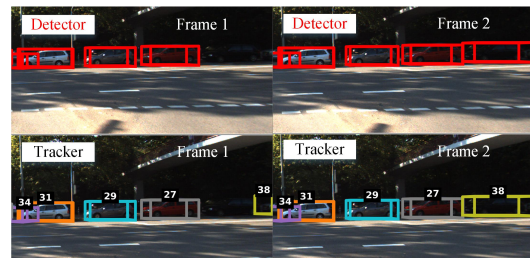


Fig. 10. Qualitative Evaluation - An example from video 15 in test set where ID switching occurs for Track 38 due to low-lit conditions

MOTP which is better than most of the online methods. Our Mostly Tracked and Mostly Lost (MT & ML) values are also competitive which show the effectiveness of our data association approach. Further, our approach gives inferences in 3D and KITTI evaluations are done in 2D, which is not completely representative of our approach. It should also be noted that none of these approaches use deep learning for data association. On the other side, we have used a simple Kalman filter for state estimation and motion prediction which could potentially be improved by better tuning of

parameters or trying out more sophisticated approaches for track management.

After optimizing the convolution operation in subsection *A.4* with selective dot products our tracking algorithm has an average runtime of 0.04s per frame ( 25 Hz) on Nvidia GeForce GTX 1080 Ti and with a single thread on Intel Core i7-7700 CPU @ 3.60GHz.

## V. Conclusions

In this paper, we presented a solution to the problem of data association in 3D online multi-object tracking using deep learning with multi-modal data. We have shown that a learning-based data association framework helps in combining different similarity cues in the data and provides more accurate associations than conventional approaches, which helps in increased overall tracking performance. We demonstrated the effectiveness of the tracker built using this model with a multitude of experiments and evaluations and show competitive results in the KITTI tracking benchmark. In the future, we plan to integrate this solution with an object detection framework more tightly and perform end-to-end training.

## References

[1] K. Okuma, A. Taleghani, N. de Freitas, J. J. Little, and D. G. Lowe, "A boosted particle filter: Multitarget detection and tracking," in *ECCV*, 2004.

[2] B. Wu and R. Nevatia, "Detection and tracking of multiple, partially occluded humans by bayesian combination of edgelet based part detectors," *International Journal of Computer Vision*, vol. 75, pp. 247–266, 01 2007.

[3] "Tracking multiple targets based on min-cost network flows with detection in rgb-d data," *Int. J. Comput. Sci. Eng.*, vol. 15, no. 3-4, pp. 330–339, Jan. 2017.

[4] J. Berclaz, F. Fleuret, E. Türetken, and P. Fua, "Multiple object tracking using k-shortest paths optimization," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 33, pp. 1806–1819, 2011.

[5] J. Munkres, "Algorithms for the assignment and transportation problems," 1957.

[6] H. Kuhn, "The hungarian method for the assignment problem," *Naval Res. Logist. Quart.*, vol. 2, pp. 83–98, 01 1955.

[7] M. D. Breitenstein, F. Reichlin, B. Leibe, E. Koller-Meier, and L. V. Gool, "Robust tracking-by-detection using a detector confidence particle filter, what," *2009 IEEE 12th International Conference on Computer Vision*, pp. 1515–1522, 2009.

[8] C. Kim, F. Li, A. Ciptadi, and J. M. Rehg, "Multiple hypothesis tracking revisited," *2015 IEEE International Conference on Computer Vision (ICCV)*, pp. 4696–4704, 2015.

[9] S.-I. Oh and H.-B. Kang, "Multiple objects fusion tracker using a matching network for adaptively represented instance pairs," *Sensors (Basel, Switzerland)*, vol. 17, 04 2017.

[10] A. Milan, S. H. Rezatofighi, A. R. Dick, K. Schindler, and I. D. Reid, "Online multi-target tracking using recurrent neural networks," *CoRR*, vol. abs/1604.03635, 2016.

[11] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The KITTI Dataset," in *IJRR*, 2013, pp. 1229–1235.

[12] D. B. Reid, "an algorithm for tracking multiple targets," vol. 24, 02 1978, pp. 1202 – 1211.

[13] T. E. Fortmann, Y. bar shalom, and M. Scheffe, "Multi-target tracking using joint probabilistic data association," vol. 2, 01 1981, pp. 807 – 812.

[14] A. A. Butt and R. T. Collins, "Multi-target tracking by lagrangian relaxation to min-cost network flow," *2013 IEEE Conference on Computer Vision and Pattern Recognition*, pp. 1846–1853, 2013.

[15] L. Zhang, Y. Li, and R. Nevatia, "Global data association for multi-object tracking using network flows," 06 2008.

[16] Z. Xi, H. Liu, H. Liu, and B. Yang, "Multiple object tracking using the shortest path faster association algorithm," *TheScientificWorldJournal*, vol. 2014, p. 481719, 08 2014.

[17] Z. Khan, T. Balch, and F. Dellaert, "Mcmc-based particle filtering for tracking a variable number of interacting targets," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 11, pp. 1805–1819, Nov 2005.

[18] A. Andriyenko, K. Schindler, and S. Roth, "Discrete-continuous optimization for multi-target tracking," in *CVPR*, 2012.

[19] R. E. Kálmán, "A new approach to linear filtering and prediction problems," 2000.

[20] N. J. Gordon, D. Salmond, and A. F. M. Smith, "Novel approach to nonlinear / non-gaussian bayesian state estimation," 2004.

[21] F. Y. Shih and C. C. Pu, "A skeletonization algorithm by maxima tracking on euclidean distance transform," *Pattern Recognition*, vol. 28, pp. 331–341, 1995.

[22] U. Franke, C. Rabe, H. Badino, and S. Gehrig, "6d-vision: Fusion of stereo and motion for robust environment perception," vol. 3663, 08 2005, pp. 216–223.

[23] A. Roshan Zamir, A. Dehghan, and M. Shah, "Gmcp-tracker: Global multi-object tracking using generalized minimum clique graphs," *ECCV*, vol. 7573, 01 2012.

[24] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *Int. J. Comput. Vision*, vol. 60, no. 2, pp. 91–110, Nov. 2004.

[25] W. Choi, "Near-online multi-target tracking with aggregated local flow descriptor," 04 2015.

[26] J. Son, M. Baek, M. Cho, and B. Han, "Multi-object tracking with quadruplet convolutional neural networks," *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 3786–3795, 2017.

[27] A. Sadeghian, A. Alahi, and S. Savarese, "Tracking the untrackable: Learning to track multiple cues with long-term dependencies," *CoRR*, vol. abs/1701.01909, 2017.

[28] L. Leal-Taixé, C. Canton-Ferrer, and K. Schindler, "Learning by tracking: Siamese CNN for robust target association," *CoRR*, vol. abs/1604.07866, 2016.

[29] D. Frossard and R. Urtasun, "End-to-end learning of multi-sensor 3d tracking by detection," *CoRR*, vol. abs/1806.11534, 2018.

[30] J. Ku, M. Mozifian, J. Lee, A. Harakeh, and S. L. Waslander, "Joint 3d proposal generation and object detection from view aggregation," arxiv:1712.02294, 2017.

[31] M. Lin, Q. Chen, and S. Yan, "Network in network," arxiv:1312.4400, 2013.

[32] H. Lee, R. Grosse, R. Ranganath, and A. Y. Ng, "Convolutional deep belief networks for scalable unsupervised learning of hierarchical representations," in *ICML*, 2009, pp. 609–616.

[33] M. D. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," in *ECCV*, 2014, pp. 818–833.

[34] M. Lin, Q. Chen, and S. Yan, "Network in network," *arXiv preprint arXiv:1312.4400*, 2013.

[35] T.-Y. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *ICCV*, 2017, pp. 2999–3007.

[36] F. Yu and V. Koltun, "Multi-scale context aggregation by dilated convolutions," arxiv:1511.07122, 2015.

[37] H. Ryuhei, F. Aito, N. Keisuke, I. Tomoyuki, and H. Shuhei, "Effective use of dilated convolutions for segmenting small object instances in remote sensing imagery," in *WACV*, 2018, pp. 1442–1450.

[38] A. O. Pak, J. Correa, M. Adams, D. Clark, E. Delande, J. Houssineau, and J. Franco, "Joint target detection and tracking filter for chilbolton advanced meteorological radar data processing," 2016.

[39] S. Sharma, J. A. Ansari, J. K. Murthy, and K. M. Krishna, "Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking," *CoRR*, vol. abs/1802.09298, 2018.

[40] J. Xiao, H. Cheng, H. Sawhney, and F. Han, "Vehicle detection and tracking in wide field-of-view aerial video," 2010.

[41] S. Scheidegger, J. Benjaminsson, E. Rosenberg, A. Krishnan, and K. Granstrom, "Mono-camera 3d multi-object tracking using deep learning detections and pmbm filtering," *arXiv preprint arXiv:1802.09975*, 2018.

[42] G. Gündüz and T. Acarman, "A lightweight online multiple object vehicle tracking method," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 427–432.

[43] Y. Xiang, A. Alahi, and S. Savarese, "Learning to track: Online multi-object tracking by decision making," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 4705–4713.

[44] B. Hurl, K. Czarnecki, and S. Waslander, "Precise synthetic image and lidar (presil) dataset for autonomous vehicle perception," *arXiv preprint arXiv:1905.00160v1*, 2019.

[45] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," arxiv:1412.6980, 2014.

[46] K. Bernardin and R. Stiefelhagen, "Evaluating multiple object tracking performance: the clear mot metrics," *Journal on Image and Video Processing*, vol. 2008, p. 1, 2008.