

# No Blind Spots: Full-Surround Multi-Object Tracking for Autonomous Vehicles using Cameras & LiDARs

Akshay Rangesh<sup>†</sup>, *Member, IEEE*, and Mohan M. Trivedi<sup>†</sup>, *Fellow, IEEE*

**Abstract**—Online multi-object tracking (MOT) is extremely important for high-level spatial reasoning and path planning for autonomous and highly-automated vehicles. In this paper, we present a modular framework for tracking multiple objects (vehicles), capable of accepting object proposals from different sensor modalities (vision and range) and a variable number of sensors, to produce continuous object tracks. This work is a generalization of the MDP framework for MOT proposed in [1], with some key extensions - First, we track objects across multiple cameras and across different sensor modalities. This is done by fusing object proposals across sensors accurately and efficiently. Second, the objects of interest (targets) are tracked directly in the *real world*. This is a departure from traditional techniques where objects are simply tracked in the image plane. Doing so allows the tracks to be readily used by an autonomous agent for navigation and related tasks.

To verify the effectiveness of our approach, we test it on real world highway data collected from a heavily sensorized testbed capable of capturing full-surround information. We demonstrate that our framework is well-suited to track objects through entire maneuvers around the ego-vehicle, some of which take more than a few minutes to complete. We also leverage the modularity of our approach by comparing the effects of including/excluding different sensors, changing the total number of sensors, and the quality of object proposals on the final tracking result.

**Index Terms**—Multi-object tracking (MOT), panoramic surround behavior analysis, highly autonomous vehicles, computer vision, sensor fusion, collision avoidance, path planning.

## I. INTRODUCTION

**T**RACKING for autonomous vehicles involves accurately identifying and localizing dynamic objects in the environment surrounding the vehicle. Tracking of surround vehicles is essential for many tasks crucial to truly autonomous driving, such as *obstacle avoidance*, *path planning*, and *intent recognition*. To be useful for such high-level reasoning, the generated tracks should be accurate, long and robust to sensor noise. In this study, we propose a full-surround MOT framework to create such desirable tracks.

Traditional MOT techniques for autonomous vehicles can roughly be categorized into 3 groups based on the sensory inputs they use - 1) dense point clouds from range sensors, 2) vision sensors, and 3) a fusion of range and vision sensors. Studies like [2]–[5] make use of dense point clouds created by 3D LiDARs like the Velodyne HDL-64E. Such sensors,

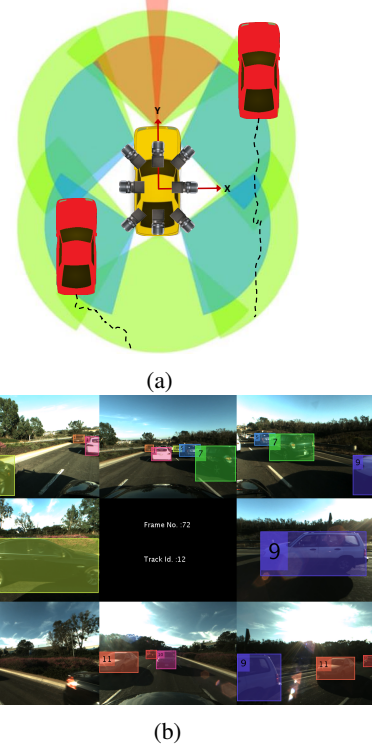


Fig. 1: (a) Illustration of online MOT for autonomous vehicles. The surrounding vehicles (in red) are tracked in a right-handed coordinate system centered on the ego-vehicle (center). The ego-vehicle has full-surround coverage from vision and range sensors, and must fuse proposals from each of them to generate continuous tracks (dotted lines) in the real world.

(b) An example of images captured from a full-surround camera array mounted on our testbed, along with color coded vehicle annotations.

although bulky and expensive, are capable of capturing finer details of the surroundings owing to its high vertical resolution. Trackers can therefore create suitable mid-level representations like 2.5D grids, voxels etc. that retain unique statistics of the volume they enclose, and group such units together to form coherent objects that can be tracked. It must be noted however, that these approaches are reliant on having dense point representations of the scene, and would not scale well to LiDAR sensors that have much fewer scan layers. On the other hand, studies such as [6]–[9] make use of stereo vision alone to perform tracking. The pipeline usually involves

<sup>†</sup>The authors are with the Laboratory for Intelligent & Safe Automobiles, UC San Diego.

<sup>‡</sup>[Video results](#)

<sup>‡</sup>[Dataset download](#)

estimating the disparity image and optionally creating a 3D point cloud, followed by similar mid-level representations like stixels, voxels etc. which are then tracked from frame to frame. These sensors are limited by the quality of disparity estimates and the field of view (FoV) of the stereo pair. Unlike 3D LiDAR based systems, they are unable to track objects in full-surround. There are other single camera approaches to surround vehicle behavior analysis [10], [11], but they too are limited in their FoVs and localization capabilities. Finally, there are fusion based approaches like [12]–[14], that make use of LiDARs, stereo pairs, monocular cameras, and Radars in a variety of configurations. These techniques either perform *early* or *late* fusion based on their sensor setup and algorithmic needs. However, none of them seem to offer full-surround solutions for vision sensors, and are ultimately limited to fusion only in the FoV of the vision sensors.

In this study, we take a different approach to full-surround MOT and try to overcome some of the limitations in previous approaches. Specifically, we propose a framework to perform full-surround MOT using calibrated camera arrays, with varying degrees of overlapping FoVs, and with an option to include low resolution range sensors for accurate localization of objects in 3D. We term this the M<sup>3</sup>OT framework, which stands for multi-perspective, multi-modal, multi-object tracking framework. To train and validate the M<sup>3</sup>OT framework, we make use of naturalistic driving data collected from our testbed (illustrated in Figure 1) that has full-surround coverage from vision and range modalities.

Since we use vision as our primary perception modality, we leverage recent approaches from the 2D MOT community which studies tracking of multiple objects in the 2D image plane. Recent progress in 2D MOT has focused on the *tracking-by-detection* strategy, where object detections from a category detector are linked to form trajectories of the targets. To perform tracking-by-detection *online* (i.e. in a causal fashion), the major challenge is to correctly associate noisy object detections in the current video frame with previously tracked objects. The basis for any data association algorithm is a similarity function between object detections and targets. To handle ambiguities in association, it is useful to combine different cues in computing the similarity, and learn an association based on these cues. Many recent 2D MOT methods such as [15]–[19] use some form of learning (online or offline) to accomplish data association. Similar to these studies, we formulate the online multi-object tracking problem using Markov Decision Processes (MDPs) proposed in [1], where the lifetime of an object is modeled with a MDP (see Figure 5), and multiple MDPs are assembled for multi-object tracking. In this method, learning a similarity function for data association is equivalent to learning a policy for the MDP. The policy learning is approached in a reinforcement learning fashion which benefits from advantages of both offline-learning and online-learning in data association. The M<sup>3</sup>OT framework is also capable to naturally handle the birth/death and appearance/disappearance of targets by treating them as state transitions in the MDP, and also benefits from the strengths of online learning approaches in single object tracking ([20]–[23]).

Our main contributions in this work can be summarized as follows - 1) We extend and improve the MDP formulation originally proposed for 2D MOT [1], and modify it to track objects in 3D (real world). 2) We make the M<sup>3</sup>OT framework capable of tracking objects across multiple vision sensors in calibrated camera arrays by carrying out efficient and accurate fusion of object proposals. 3) The M<sup>3</sup>OT framework is made highly *modular*, capable of working with any number of cameras, with varying degrees of overlapping FoVs, and with the option to include range sensors for improved localization and fusion in 3D. The above contributions and the wider scope and applicability of this work in comparison to traditional 2D MOT approaches are highlighted in Figure 2. Finally, we carry out experiments using naturalistic driving data collected on highways using full-surround sensory modalities, and validate the accuracy, robustness and modularity of our framework.

## II. RELATED RESEARCH

We highlight some representative works in 2D and 3D MOT below. We also summarize the some key aspects of related 3D MOT studies in Table I. For a recent survey on detection and tracking for intelligent vehicles, we refer the reader to [25].

**2D MOT:** Recent research in MOT has focused on tracking-by-detection, where the main challenge is data association for linking object detections to targets. Majority of batch (*offline*) methods ([18], [26]–[31]) formulate MOT as a global optimization problem in a graph-based representation, while *online* methods solve the data association problem either probabilistically [32]–[34] or deterministically (e.g., Hungarian algorithm [35] in [15], [16] or greedy association [36]). A core component in any data association algorithm is a similarity function between objects. Both batch methods [17], [18] and online methods [15], [16], [19] have explored the idea of learning to track, where the goal is to learn a similarity function for data association from training data. In this work, we extend and improve the MDP framework for 2D MOT proposed in [1], which is an online method that uses reinforcement learning to learn associations.

**3D MOT for autonomous vehicles:** In [6], the authors use a stereo rig to first calculate the disparity using Semi-Global Matching (SGM), followed by height based segmentation and free-space calculation to create a mid-level representation using *stixels* that encode the height within a cell. Each stixel is then represented by a 6D state vector, which is tracked using the Extended Kalman Filter (EKF). In [7], the authors use a *voxel* based representation instead, and cluster neighboring voxels based on color to create objects that are then tracked using a greedy association model. The authors in [8] use a grid based representation of the scene, where cells are grouped to create objects, each of which is represented by a set of control points on the object surface. This creates a high dimensional state-space representation, which is accounted for by a Rao-Blackwellized particle filter. More recently, the authors of [9] propose to carry out semantic segmentation on the disparity image, which is then used to generate generic object proposals by creating a scale-space representation of the density, followed by multi-scale clustering. The proposed

TABLE I: Relevant research in 3D MOT for intelligent vehicles.

Study	Sensors used					Tracker Type			Experimental Analysis	
	Monocular camera	Stereo pair	Full-surround camera array	LiDAR	Radar	Single object tracker	Multi-object tracker	Online (causal)	Dataset	Evaluation metrics
Choi et al. [2]	-	-	-	✓	-	-	✓	✓	Proposed	Distance and velocity errors
Broggi et al. [7]	-	✓	-	-	-	-	✓	✓	Proposed	True positives, false positives, false negatives
Song et al. [3]	-	-	-	✓	-	✓	✗	✓	KITTI	Position error, intersection ratio
Cho et al. [12]	✓	-	-	✓	✓	-	✓	✓	Proposed	Correctly tracked, falsely tracked, true and false positive rates
Asvadi et al. [4]	-	-	-	✓	-	-	✓	✓	KITTI	-
Vatavu et al. [8]	-	✓	-	-	-	-	✓	✓	Proposed	-
Asvadi et al. [5]	-	-	-	✓	-	-	✓	✓	KITTI	Number of missed and false obstacles
Asvadi et al. [13]	✓	-	-	✓	-	✓	✗	✓	KITTI	Average center location errors in 2D and 3D, orientation errors
Ošep et al. [9]	-	✓	-	-	-	-	✓	✓	KITTI	Class accuracy, GOP recall, tracking recall
Allodi et al. [14]	-	✓	-	✓	-	-	✓	✓	KITTI	MOTP, IDS, Frag, average localization error
Dueholm et al. [24]	-	-	✓	-	-	-	✓	✗	VIVA Surround	MOTA, MOTP, IDS, Frag, Precision, Recall
<b>This work<sup>1</sup> (M<sup>3</sup>OT)</b>	-	-	✓	✓	-	-	✓	✓	Proposed	MOTA, MOTP, MT, ML, IDS for a variety of sensor configurations

<sup>1</sup> this framework can work with or without LiDAR sensors, and with any subset of camera sensors.

clusters are then tracked using a Quadratic Pseudo-Boolean Optimization (QPBO) framework. The work in [24] use a camera setup similar to ours, but the authors propose an offline framework for tracking, hence limiting their use to surveillance related applications.

Alternatively, there are approaches that make use of dense point clouds generated by LiDARs as opposed to creating point clouds from a disparity image. In [2], the authors first carry out ground classification based on the variance in the radius of each scan layer, followed by a 2.5D occupancy grid representation of the scene. The grid is then segmented, and regions of interest (RoIs) identified within, each of which is tracked by a standard Kalman Filter (KF). Data association is achieved by simple global nearest neighbor. Similar to this, the authors in [4] use a 2.5D occupancy grid based representation, but augment this with an occupancy grid *map* which accumulates past grids to create a coherent global map of occupancy by accounting for ego-motion. Using these two representations, a 2.5D *motion* grid is created by comparing the map with the latest occupancy grid, which isolates and identifies dynamic objects in the scene. Although the work in [5] follows the same general idea, the authors propose a piecewise ground plane estimation scheme capable of handling non-planar surfaces. In a departure from grid based methods, the authors in [3] project the 3D point cloud onto a virtual image plane, creating an object appearance model based on 4 image-based cues for each template of the desired target. A particle filtering framework is implemented, where the particle with least reconstruction error with respect to the stored template is chosen to update the tracker. Background filtering and occlusion detection are implemented to improve performance.

Finally, we list recent methods that rely on fusion of different sensor modalities to function. In [12], the authors propose an EKF based fusion scheme, where measurements from each modality are fed sequentially. Vision is used to classify the object category, which is then used to choose

appropriate motion and observation models for the object. Once again, observations are associated based on a global nearest neighbor policy. This is somewhat similar to the work in [13], where given an initial 3D detection box, the authors propose to project 3D points within the box to the image plane and calculate the convex hull of projected points. In this case, a KF is used to perform both fusion and tracking, where fusion is carried out by projecting the 2D hull to a sparse 3D cloud, and using both 3D cues to perform the update. In contrast, the authors of [14] propose using the Hungarian algorithm (for bipartite matching) for both data association and fusion of object proposals from different sensors. The scores for association are obtained from Adaboost classifiers trained on high-level features. The objects are then tracked using an Unscented Kalman Filter (UKF).

### III. FUSION OF OBJECT PROPOSALS

In this study, we make use of full-surround camera arrays comprising of sensors with varying FoVs. The M<sup>3</sup>OT framework, however, is capable of working with any type and number of cameras, as long as they are calibrated. In addition to this, we also propose a variant of the framework for cases where LiDAR point clouds are available. To effectively utilize all available sensors, we propose an *early fusion* of object proposals obtained from each of them. At the very start of each time step during tracking, we identify and fuse all proposals belonging to the same object. These proposals are then utilized by the M<sup>3</sup>OT framework to carry out tracking. It must be noted that this usage of “early fusion” is in contrast to the traditional usage of the term to refer to fusion of raw sensor data to provide a merged representation.

a) *Projection & Back-projection*: It is essential to have a way of associating measurements from different sensors to track objects across different camera views, and to carry out efficient fusion across sensor modalities. This is achieved by defining a set of *projection* mappings, one from each sensor’s

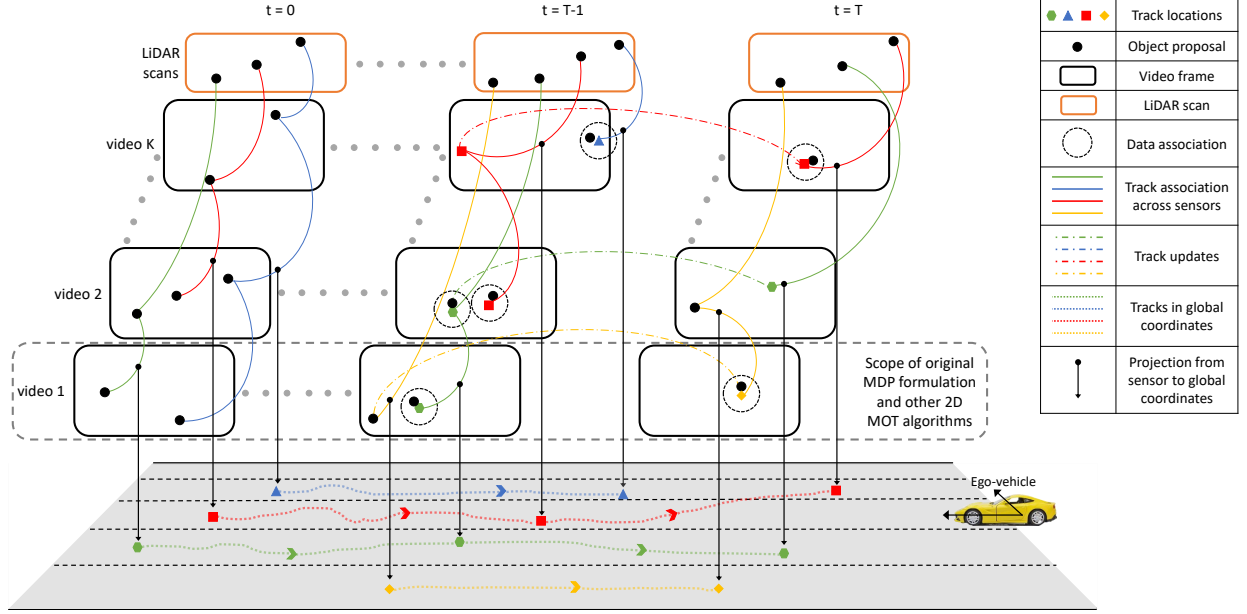


Fig. 2: Illustration of proposed M<sup>3</sup>OT framework and its scope in comparison to traditional 2D MOT algorithms. Data is associated not only within video frames, but also across other videos and sensors. The algorithm produces tracks in each individual sensor coordinate frame, and in the desired global coordinate frame using cross-modal fusion in an online manner.

unique coordinate system to the global coordinate system, and a set of *back-projection* mappings that take measurements in the global coordinate system to individual coordinate systems. In our case, the global coordinate system is centered at the mid-point of the rear axle of the ego-vehicle. The axes form a right-handed coordinate system as shown in Figure 1.

The LiDAR sensors output a 3D point cloud in a common coordinate system at every instant. This coordinate frame may either be centered about a single LiDAR sensor, or elsewhere depending on the configuration. In this case, the projection and back-projection mappings are simple 3D coordinate transformations:

$$P_{range \rightarrow G}(\mathbf{x}^{range}) = \mathbf{R}_{range} \cdot \mathbf{x}^{range} + \mathbf{t}_{range}\mathbf{1}, \quad (1)$$

and,

$$P_{range \leftarrow G}(\mathbf{x}^G) = \mathbf{R}_{range}^T \cdot \mathbf{x}^G - \mathbf{R}_{range}^T \mathbf{t}_{range}, \quad (2)$$

where  $P_{range \rightarrow G}(\cdot)$  and  $P_{range \leftarrow G}(\cdot)$  are the projection and back-projection mappings from the LiDAR (range) coordinate system to the global (G) coordinate system and vice-versa. The vectors  $\mathbf{x}^{range}$  and  $\mathbf{x}^G$  are the corresponding coordinates in the LiDAR and global coordinate frames. The  $3 \times 3$  orthonormal rotation matrix  $\mathbf{R}_{range}$  and translation vector  $\mathbf{t}_{range}$  are obtained through calibration.

Similarly, the back-projection mappings for each camera  $k \in \{1, \dots, K\}$  can be defined as:

$$P_{cam_k \leftarrow G}(\mathbf{x}^G) = (u, v)^T, \quad (3)$$

$$\text{s.t.} \begin{bmatrix} u \\ v \\ 1 \end{bmatrix} = \mathbf{C}_k \cdot [\mathbf{R}_k | \mathbf{t}_k] \cdot [(\mathbf{x}^G)^T | 1]^T, \quad (4)$$

where the set of camera calibration matrices  $\{\mathbf{C}_k\}_{k=1}^K$  are obtained after the intrinsic calibration of cameras, the set of

tuples  $\{(\mathbf{R}_k, \mathbf{t}_k)\}_{k=1}^K$  obtained after extrinsic calibration, and  $(u, v)^T$  denotes the pixel coordinates after back-projection.

Unlike the back-projection mappings, the projection mappings for camera sensors are not well defined. In fact, the mappings are one-to-many due to the depth ambiguity of single camera images. To find a good estimate of the projection, we use two different approaches. In case of a vision-only system, we use the *inverse perspective mapping* (IPM) approach:

$$P_{cam_k \rightarrow G}(\mathbf{x}^k) = (x, y)^T, \quad (5)$$

$$\text{s.t.} \begin{bmatrix} x \\ y \\ 1 \end{bmatrix} = \mathbf{H}_k \cdot [(\mathbf{x}^k)^T | 1]^T, \quad (6)$$

where  $\{\mathbf{H}_k\}_{k=1}^K$  are the set of homographies obtained after IPM calibration. Since we are only concerned with lateral and longitudinal displacements of vehicles in the global coordinate system, we only require the  $(x, y)^T$  coordinates, and set the altitude coordinate to a fixed number.

*b) Sensor Measurements & Object Proposals:* As we adopt a tracking-by-detection approach, each sensor is used to produce object proposals to track and associate. In case of vision sensors, a vehicle detector is run on each individual camera's image to obtain multiple detections  $d$ , each of which is defined by a bounding box in the corresponding image. Let  $(u, v)$  denote the top left corner and  $(w, h)$  denote the width and height of a detection  $d$  respectively.

In case of a vision-only system, the corresponding location of  $d$  in the global coordinate system is obtained using the mapping  $P_{cam_k \rightarrow G}((u + \frac{w}{2}, v + h)^T)$ , where  $k$  denotes the camera from which the proposal was generated. This procedure is illustrated in Figure 3a. Alternatively, this could be replaced by a purely vision based 3D detector like the one



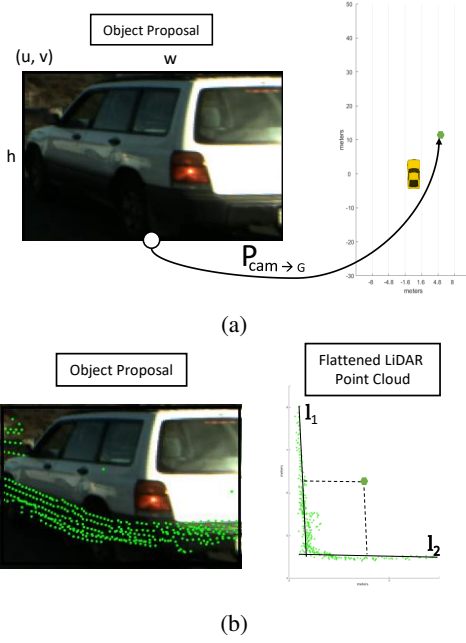


Fig. 3: Projection of object proposals using: (a) **IPM**: The bottom center of the bounding boxes are projected into the global coordinate frame (right), (b) **LiDAR point clouds**: LiDAR points that fall within a detection window are flattened and lines are fit to identify the vehicle center (right).

proposed in [37], where both the global pose and 2D bounding box of an object are obtained from the same algorithm.

In cases where LiDAR sensors are available, an alternative is considered (shown in Figure 3b). First, the back-projected LiDAR points that fall within a detection box  $d$  are identified using a look-up table with pixel coordinates as the key, and the corresponding global coordinates as the value. These points are then flattened by ignoring the altitude component of their global coordinates. Next, a line  $l_1$  is fitted to these points using RANSAC with a small inlier ratio (0.3). This line aligns with the dominant side of the detected vehicle. The other side of the vehicle corresponding to line  $l_2$  is then identified by removing all inliers from the previous step and repeating a RANSAC line fit with a slightly higher inlier ratio (0.4). Finally, the intersection of lines  $l_1$  and  $l_2$  along with the vehicle dimensions yield the global coordinates of the vehicle center. The vehicle dimensions are calculated based on the extent of the LiDAR points along a given side of the vehicle, and stored for each track separately for later use.

Depending on the type of LiDAR sensors used, object proposals along with their dimensions in the real world can be obtained. However, we decide not to make use of LiDAR proposals, but rather use vision-only proposals with high recall by trading off some of the precision. This was seen to provide sufficient proposals to track all surrounding vehicles, at the expense of more false positives which the tracker is capable of handling.

c) *Early Fusion of Proposals*: Since we operate with camera arrays with overlapping FoVs, the same vehicle may be detected in two adjacent views. It is important to identify



Fig. 4: Fusion of object proposals using LiDAR point clouds: Points common to both detections are drawn in green, and the rest are drawn red.

and *fuse* such proposals to track objects across camera views. Once again, we propose two different approaches to carry out this fusion. For vision-only systems, the fusion of proposals is carried out in 4 steps: i) Project proposals from all cameras to the global coordinate system using proposed mappings, ii) Sort all proposals in descending order based on their confidence scores (obtained from the vehicle detector), iii) Starting with the highest scoring proposal, find the subset of proposals whose euclidean distance in the global coordinate system falls within a predefined threshold. These proposals are considered to belong to the same object and removed from the original set of proposals. In practice, we use a threshold of  $1m$  for grouping proposals. iv) The projection of each proposal within this subset is set to the mean of projections of all proposals within the subset. This process is repeated for the remaining proposals until no proposals remain.

Alternatively, for a system consisting of LiDAR sensors, we project the 3D point cloud onto each individual camera image. Next, for each pair of proposals, we make a decision as to whether or not they belong to the same object. This is done by considering the back-projected LiDAR points that fall within the bounding box of each proposal (see Figure 4). Let  $\mathcal{P}_1$  and  $\mathcal{P}_2$  denote the index set of LiDAR points falling within each bounding box. Then, two proposals are said to belong to the same object if:

$$\max\left(\frac{|\mathcal{P}_1 \cap \mathcal{P}_2|}{|\mathcal{P}_1|}, \frac{|\mathcal{P}_1 \cap \mathcal{P}_2|}{|\mathcal{P}_2|}\right) \geq 0.8, \quad (7)$$

where  $|\mathcal{P}_1|$  and  $|\mathcal{P}_2|$  denote the cardinalities of sets  $\mathcal{P}_1$  and  $\mathcal{P}_2$  respectively. It should be noted that after fusion is completed, the union of LiDAR point sets that are back-projected into fused proposals can be used to obtain better projections.

#### IV. M<sup>3</sup>OT FRAMEWORK

Once we have a set of fused object proposals, we feed them into the MDP as illustrated in Figure 5. Although the MDP framework introduced in [1] forms a crucial building block of the proposed M<sup>3</sup>OT framework, we believe that extending this to account for multiple sensors and modalities is a non-trivial endeavor (see Figure 2). This section highlights and explains the modifications we propose to achieve these objectives.

##### A. Markov Decision Process

As detailed in [1], we model the lifetime of a target with a Markov Decision Process (MDP). The MDP consists of the tuple  $(\mathcal{S}, \mathcal{A}, T(\cdot, \cdot), R(\cdot, \cdot))$ , where:

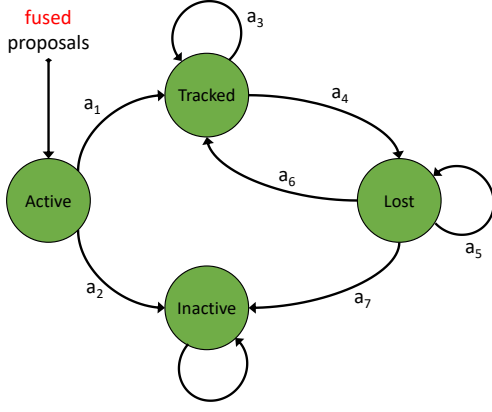


Fig. 5: The Markov Decision Process (MDP) framework proposed in [1]. In this work, we retain the structure of the MDP, and modify the actions, rewards and inputs to enable multi-sensory tracking.

- States  $s \in \mathcal{S}$  encode the status of the target.
- Actions  $a \in \mathcal{A}$  define the actions that can be taken.
- The state transition function  $T : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$  dictates how the target transitions from one state to another, given an action.
- The real-valued reward function  $R : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$  assigns the immediate reward received after executing action  $a$  in state  $s$ .

In this study, we retain the states  $\mathcal{S}$ , actions  $\mathcal{A}$  and the state transition function  $T(\cdot, \cdot)$  from the MDP framework for 2D MOT [1], while changing only the reward function  $R(\cdot, \cdot)$ .

**States:** The state space is partitioned into four subspaces, i.e.,  $\mathcal{S} = \mathcal{S}_{Active} \cup \mathcal{S}_{Tracked} \cup \mathcal{S}_{Lost} \cup \mathcal{S}_{Inactive}$ , where each subspace contains infinite number of states which encode the information of the target depending on the feature representation, such as appearance, location, size and history of the target. Figure 5 illustrates the transitions between the four subspaces. *Active* is the initial state for any target. Whenever an object is detected by the object detector, it enters an *Active* state. An active target can transition to *Tracked* or *Inactive*. Ideally, a true positive from the object detector should transition to a *Tracked* state, while a false alarm should enter an *Inactive* state. A tracked target can stay tracked, or transition to a *Lost* state if the target is not visible due to some reason, e.g. occlusion, or disappearance from sensor range. Likewise, a lost target can stay *Lost*, or go back to a *Tracked* state if it appears again, or transition to an *Inactive* state if it has been lost for a sufficiently long time. Finally, *Inactive* is the terminal state for any target, i.e., an inactive target stays inactive forever.

**Actions and Transition Function:** Seven possible transitions are designed between the state subspaces, which correspond to seven actions in our target MDP. Figure 5 illustrates these transitions and actions. In the MDP, all the actions are deterministic, i.e., given the current state and an action, we specify a new state for the target. For example, executing action  $a_6$  on a *Lost* target would transfer the target into a *Tracked* state, i.e.,  $T(s_{Lost}, a_6) = s_{Tracked}$ .

**Reward Function:** As in the original study [1], we learn the reward function from training data, i.e., an inverse reinforcement learning problem, where we use ground truth trajectories of the targets as supervision.

## B. Policy

a) *Policy in Active States:* In an *Active* state  $s$ , the MDP makes the decision between transferring an object proposal into a *Tracked* or *Inactive* state based on whether the detection is true or noisy. To do this, we train a set of binary Support Vector Machines (SVM) offline, one for each camera view, to classify a detection belonging to that view into *Tracked* or *Inactive* states using a normalized 5D feature vector  $\phi_{Active}(s)$ , i.e., 2D image plane coordinates, width, height and score of the detection, where training examples are collected from training video sequences.

This is equivalent to learning the reward function:

$$R_{Active}(s, a) = y(a)((\mathbf{w}_{Active}^k)^T \cdot \phi_{Active}(s) + b_{Active}^k), \quad (8)$$

for an object proposal belonging to camera  $k \in \{1, \dots, K\}$ .  $(\mathbf{w}_{Active}^k, b_{Active}^k)$  defines the learned weights and bias of the SVM for camera  $k$ ,  $y(a) = +1$  if action  $a = a_1$ , and  $y(a) = -1$  if  $a = a_2$  (see Figure 5). Training a separate SVM for each camera view allows weights to be learned based on object dimensions and locations in that particular view, and thus works better than training a single SVM for all views. Since a single object can result in multiple proposals, we initialize a tracker for that object if any of the fused proposals result in a positive reward. Note that a false positive from the object detector can still be misclassified and transferred to a *Tracked* state, which we then leave to be handled by the MDP in *Tracked* and *Lost* states.

b) *Policy in Tracked States:* In a *Tracked* state, the MDP needs to decide whether to keep tracking the target or to transfer it to a *Lost* state. As long as the target is visible, we should keep tracking it. Else, it should be marked “lost”. We build an appearance model for the target online and use it to track the target. If the appearance model is able to successfully track the target in the next video frame, the MDP leaves the target in a *Tracked* state. Otherwise, the target is transferred to a *Lost* state.

**Template Representation:** The appearance of the target is simply represented by a template that is an image patch of the target in a video frame. Whenever an object detection is transferred to a *Tracked* state, we initialize the target template with the detection bounding box. If the target is initialized with multiple fused proposals, then each detection is stored as a template. We make note of detections obtained from different camera views, and use these to model the appearance of the target in that view. This is crucial to track objects across camera views under varying perspective changes. When the target is being tracked, the MDP collects its templates in the tracked frames to represent the history of the target, which will be used in the *Lost* state for decision making.

**Template Tracking:** Tracking of templates is carried out by performing dense optical flow as described in [1]. The stability of the tracking is measured using the median of the

Forward-Backward (FB) errors [23] of all sampled points:  $e_{medFB} = \text{median}(e(\mathbf{u}_i)_{i=1}^n)$ , where  $\mathbf{u}_i$  denotes each sampled point, and  $n$  is the total number of points. If  $e_{medFB}$  is larger than some threshold, the tracking is considered to be unstable. Moreover, after filtering out unstable matches whose FB error is larger than the threshold, a new bounding box of the target is predicted using the remaining matches by measuring scale change between points before and after. This process is carried out for all camera views in which a target template has been initialized and tracking is in progress.

Similar to the original MDP framework, we use the optical flow information in addition to the object proposals history to prevent drifting of the tracker. To do this, we compute the bounding box overlap between the target box for  $l$  past frames, and the corresponding detections in each of those frames. Then we compute the mean bounding box overlap for the past  $L$  tracked frames  $o_{mean}$  as another cue to make the decision. Once again, this process is repeated for each camera view the target is being tracked in. In addition to the above features, we also *gate* the target track. This involves introducing a check to see if the current global position of the tracked target falls within a window (gate) of it's last known global position. This forbids the target track from latching onto objects that appear close on the image plane, yet are much farther away in the global coordinate frame. We denote the last known global position and the currently tracked global position of the target as  $\mathbf{x}^G(t-1)$  and  $\hat{\mathbf{x}}^G(t)$  respectively.

Finally, we define the reward function in a Tracked state  $s$  using the feature set  $\phi_{Tracked}(s) = (\{e_{medFB}^{k'}\}_{k'=1}^{K'}, \{o_{mean}^{k'}\}_{k'=1}^{K'}, \mathbf{x}^G(t-1), \hat{\mathbf{x}}^G(t))$  as:

$$R_{Tracked}(s, a) = \begin{cases} y(a), & \text{if } \exists k' \in \{1, \dots, K'\} \text{ s.t.} \\ & (e_{medFB}^{k'} < e_0) \wedge (o_{mean}^{k'} > o_0) \\ & \wedge (|\mathbf{x}^G(t-1) - \hat{\mathbf{x}}^G(t)| \leq t_{gate}), \\ -y(a), & \text{otherwise,} \end{cases} \quad (9)$$

where  $e_0$  and  $o_0$  are fixed thresholds,  $y(a) = +1$  if action  $a = a_3$ , and  $y(a) = -1$  if  $a = a_4$  (see Figure 5).  $k'$  above indexes camera views in which the target is currently being tracked and  $t_{gate}$  denotes the gating threshold. So the MDP keeps the target in a Tracked state if  $e_{medFB}$  is smaller and  $o_{mean}$  is larger than their respective thresholds for any one of  $K'$  camera views in addition to satisfying the gating check. Otherwise, the target is transferred to a Lost state.

**Template Updating:** The appearance model of the target needs to be regularly updated in order to accommodate appearance changes. As in the original work, we adopt a “lazy” updating rule and resort to the object detector in preventing tracking drift. This is done so that we don’t accumulate tracking errors, but rather rely on data association to handle appearance changes and continue tracking. In addition to this, templates are initialized in views where the target is yet to be tracked by using proposals that are fused with detections corresponding to the tracked location in an adjacent camera view. This helps track objects that move across adjacent camera views, by creating target templates in the new view as soon as they are made available.

```

input : Set of multi-video sequences  $\mathcal{V} = \{(v_i^1, \dots, v_i^K)\}_{i=1}^N$ ,
        ground truth trajectories  $\mathcal{T}_i = \{t_{ij}\}_{j=1}^{N_i}$ , object proposals
         $\mathcal{D}_i = \{d_{im}\}_{m=1}^{M_i}$  and their corresponding projections for
        each multi-video sequence  $(v_i^1, \dots, v_i^K)$ 

output: Binary classifiers  $\{(\mathbf{w}_{Lost}^k, b_{Lost}^k)\}_{k=1}^K$  for data
        association

1 repeat
2   foreach multi-video sequence  $(v_i^1, \dots, v_i^K)$  in  $\mathcal{V}$  do
3     foreach target  $t_{ij}$  do
4       Initialize the MDP in an Active state;
5        $l \leftarrow$  index of the first frame in which  $t_{ij}$  is correctly
        detected;
6       Transfer the MDP to a Tracked state and initialize
        the target template for each camera view in which
        target is observed;
7       while  $l \leq$  index of last frame of  $t_{ij}$  do
8         Fuse object proposals as described in III;
9         Follow the current policy and choose an action
         $a$ ;
10        Compute the action  $a_{gt}$  according to the
        ground truth;
11        if current state is lost and  $a \neq a_{gt}$  then
12          foreach camera view  $k$  in which the target
            has been seen do
13            Decide the label  $y_{m_k}^k$  of the pair
             $(t_{ij}^k, d_{im_k}^k)$ ;
14             $S^k \leftarrow S^k \cup \{(\phi(t_{ij}^k, d_{im_k}^k), y_{m_k}^k)\}$ ;
15             $(\mathbf{w}_{Lost}^k, b_{Lost}^k) \leftarrow$  solution of Eq.11
            on  $S^k$ ;
16          end
17          break;
18        else
19          Execute action  $a$ ;
20           $l \leftarrow l + 1$ ;
21        end
22        if  $l >$  index of last frame of  $t_{ij}$  then
23          Mark target  $t_{ij}$  as successfully tracked;
24        end
25      end
26    end
27  end
28 until all targets are successfully tracked;

```

**Algorithm 1:** Reinforcement learning of the binary classifier for data association.

*c) Policy in Lost States:* In a Lost state, the MDP needs to decide whether to keep the target in a Lost state, or transition it to a Tracked state, or mark it as Inactive. We simply mark a lost target as Inactive and terminate the tracking if the target has been lost for more than  $L_{Lost}$  frames. The more challenging task is to make the decision between tracking the target and keeping it as lost. This is treated as a data association problem where, in order to transfer a lost target into a Tracked state, the target needs to be associated with an object proposal, else, the target retains its Lost state.

**Data Association:** Let  $t$  denote a lost target, and  $d$  be an object detection. The goal of data association is to predict the label  $y \in \{+1, -1\}$  of the pair  $(t, d)$  indicating that the target is linked ( $y = +1$ ) or not linked ( $y = -1$ ) to the detection. Assuming that the detection  $d$  belongs to camera view  $k$ , this binary classification is performed using the real-valued linear function  $f^k(t, d) = (\mathbf{w}_{Lost}^k)^T \cdot \phi_{Lost}(t, d) + b_{Lost}^k$ , where  $(\mathbf{w}_{Lost}^k, b_{Lost}^k)$  are the parameters that control the function (for camera view  $k$ ), and  $\phi_{Lost}(t, d)$  is the feature vector which captures the similarity between the target and the detection.

TABLE II: Features used for data association [1]. We introduce two new features (highlighted in bold) based on the global coordinate positions of targets and detections.

Type	Notation	Feature Description
FB error	$\phi_1, \dots, \phi_5$	Mean of the median forward-backward errors from the entire, left half, right half, upper half and lower half of the templates obtained from optical flow
NCC	$\phi_6$	Mean of the median Normalized Correlation Coefficients (NCC) between image patches around the matched points in optical flow
	$\phi_7$	Mean of the median Normalized Correlation Coefficients (NCC) between image patches around the matched points obtained from optical flow
Height ratio	$\phi_8$	Mean of ratios of the bounding box height of the detection to that of the predicted bounding boxes obtained from optical flow
	$\phi_9$	Ratio of the bounding box height of the target to that of the detection
Overlap	$\phi_{10}$	Mean of the bounding box overlaps between the detection and the predicted bounding boxes from optical flow
Score	$\phi_{11}$	Normalized detection score
Distance	$\phi_{12}$	Euclidean distance between the centers of the target and the detection after motion prediction of the target with a linear velocity model
	$\phi_{13}$	<b>Lateral offset between last known global coordinate position of the target and that of the detection</b>
	$\phi_{14}$	<b>Longitudinal offset between last known global coordinate position of the target and that of the detection</b>

The decision rule is given by  $y = +1$  if  $f^k(t, d) \geq 0$ , else  $y = -1$ . Consequently, the reward function for data association in a lost state  $s$  given the feature set  $\{\phi_{Lost}(t, d_j)\}_{m=1}^M$  is defined as

$$R_{Lost}(s, a) = y(a) \left( \max_{m=1}^M ((\mathbf{w}_{Lost}^{k_m})^T \cdot \phi_{Lost}(t, d_m) + b_{Lost}^{k_m}) \right), \quad (10)$$

where  $y(a) = +1$  if action  $a = a_6$ ,  $y(a) = -1$  if  $a = a_5$  (see Figure 5), and  $m$  indexes  $M$  potential detections for association. Potential detections for association with a target are simply obtained by applying a gating function around the last known location of the target in the global coordinate system. Note that based on which camera view each detection  $d_m$  originates from, the appropriate weights  $(\mathbf{w}_{Lost}^{k_m}, b_{Lost}^{k_m})$  associated with that view are used. As a result, the task of policy learning in the Lost state reduces to learning the set of parameters  $\{(\mathbf{w}_{Lost}^k, b_{Lost}^k)\}_{k=1}^K$  for the decision functions  $\{f^k(t, d)\}_{k=1}^K$ .

**Reinforcement Learning:** We train the binary classifiers described above using the reinforcement learning paradigm. Let  $\mathcal{V} = \{(v_i^1, \dots, v_i^K)\}_{i=1}^N$  denote a set of multi-video sequences for training, where  $N$  is the number of sequences and  $K$  is the total number of camera views. Suppose there are  $N_i$  ground truth targets  $\mathcal{T}_i = \{t_{ij}\}_{j=1}^{N_i}$  in the  $i^{th}$  multi-video sequence  $(v_i^1, \dots, v_i^K)$ . Our goal is to train the MDP to successfully track all these targets across all camera views they appear in. We start training with initial weights  $(\mathbf{w}_0^k, b_0^k)$  and an empty training set  $\mathcal{S}_0^k = \emptyset$  for the binary classifier corresponding to each camera view  $k$ . Note that when the weights of the binary classifiers are specified, we have a complete policy for the MDP to follow. So the training

algorithm loops over all the multi-video sequences and all the targets, follows the current policy of the MDP to track the targets. The binary classifier or the policy is updated only when the MDP makes a mistake in data association. In this case, the MDP takes a different action than what is indicated by the ground truth trajectory. Suppose the MDP is tracking the  $j^{th}$  target  $t_{ij}$  in the video  $v_i^k$ , and on the  $l^{th}$  frame of the video, the MDP is in a lost state. Consider the two types of mistakes that can happen: i) The MDP associates the target  $t_{ij}(l)$  to an object detection  $d_m^k$  which disagrees with the ground truth, i.e., the target is incorrectly associated to a detection. Then  $\phi(t_{ij}(l), d_m^k)$  is added to the training set  $\mathcal{S}^k$  of the binary classifier for camera  $k$  as a negative example. ii) The MDP decides to not associate the target to any detection, but the target is visible and correctly detected by a detection  $d_m^k$  based on the ground truth, i.e., the MDP missed the correct association. Then  $\phi(t_{ij}(l), d_m^k)$  is added to the training set as a positive example. After the training set has been augmented, we update the binary classifier by re-training it on the new training set. Specifically, given the current training set  $\mathcal{S}^k = \{(\phi(t_m^k, d_m^k), y_m^k)\}_{m=1}^M$ , we solve the following soft-margin optimization problem to obtain a max-margin classifier for data association in camera view  $k$ :

$$\begin{aligned} \min_{\mathbf{w}_{Lost}^k, b_{Lost}^k, \xi} \quad & \frac{1}{2} \|\mathbf{w}_{Lost}^k\|^2 + C \sum_{m=1}^M \xi_m \\ \text{s.t.} \quad & y_m^k ((\mathbf{w}_{Lost}^k)^T \cdot \phi(t_m^k, d_m^k) + b_{Lost}^k) \geq 1 - \xi_m, \xi_m \geq 0, \forall m, \end{aligned} \quad (11)$$

where  $\xi_m, m = 1, \dots, M$  are the slack variables, and  $C$  is a regularization parameter. Once the classifier has been updated, we obtain a new policy which is used in the next iteration of the training process. Note that based on which view the data association is carried out in, the weights of the classifier in that view are updated in each iteration. We keep iterating and updating the policy until all the targets are successfully tracked. Algorithm 1 summarizes the policy learning algorithm.

**Feature Representation:** We retain the same feature representation described in [1], but add two features based on the lateral and longitudinal displacements of the last known target location and the object proposal location in the global coordinate system. This leverages 3D information that is otherwise unavailable in 2D MOT. Table II summarizes our feature representation.

### C. Multi-Object Tracking with MDPs

After learning the policy/reward of the MDP, we apply it to the multi-object tracking problem. We dedicate one MDP for each target, and the MDP follows the learned policy to track the object. Given a new input video frame, targets in tracked states are processed first to determine whether they should stay as tracked or transfer to lost states. Then we compute pairwise similarity between lost targets and object detections which are not covered by the tracked targets, where non-maximum suppression based on bounding box overlap is employed to suppress covered detections, and the similarity score is



```

input : A multi-video sequence  $(v^1, \dots, v^K)$ , corresponding object
proposals  $\mathcal{D} = \{d_m\}_{m=1}^M$  and their projections, learned binary
classifier weights  $\{(\mathbf{w}_{Active}^k, b_{Active}^k)\}_{k=1}^K$  and
 $\{(\mathbf{w}_{Lost}^k, b_{Lost}^k)\}_{k=1}^K$ 
output: Trajectories of targets  $\mathcal{T} = \{t_j\}_{j=1}^N$  in the sequence
1 foreach frame  $l$  in  $(v^1, \dots, v^K)$  do
2   Fuse object proposals as described in III;
3   /* process targets in tracked states */
4   foreach tracked target  $t_j$  in  $\mathcal{T}$  do
5     Follow the policy, move the MDP of  $t_j$  to the next state;
6   end
7   /* process targets in lost states */
8   foreach lost target  $t_j$  in  $\mathcal{T}$  do
9     foreach proposal  $d_m$  not covered by any tracked target do
10      Compute
11       $f^{k_m}(t_j, d_m) = (\mathbf{w}^{k_m})_{Lost}^T \cdot \phi(t_j, d_m) + b_{Lost}^{k_m}$ ;
12    end
13  end
14  Data association with Hungarian algorithm for the lost targets;
15  Initialize target templates for uninitialized camera views using
  matched (fused) proposals;
16  foreach lost target  $t_j$  in  $\mathcal{T}$  do
17    Follow the assignment, move the MDP of  $t_j$  to the next state;
18  end
19  /* initialize new targets */
20  foreach proposal  $d_m$  not covered by any tracked target in  $\mathcal{T}$  do
21    Initialize a MDP for a new target  $t$  with proposal  $d_m$ ;
22    if action  $a_1$  is taken following the policy then
23      Transfer  $t$  to the tracked state and initialize the target
      template for each camera view in which target is
      observed;
24       $\mathcal{T} \leftarrow \mathcal{T} \cup \{t\}$ ;
25    else
26      Transfer  $t$  to the Inactive state;
27    end
28  end

```

**Algorithm 2:** Multi-object tracking with MDPs.

computed by the binary classifier for data association. After that, the similarity scores are used in the Hungarian algorithm [35] to obtain the assignment between detections and lost targets. According to the assignment, lost targets which are linked to some object detections are transferred to tracked states. Otherwise, they stay as lost. Finally, we initialize a MDP for each object detection which is not covered by any tracked target. Algorithm 2 describes our 3D MOT algorithm using MDPs in detail. Note that, tracked targets have higher priority than lost targets in tracking, and detections covered by tracked targets are suppressed to reduce ambiguities in data association.

## V. EXPERIMENTAL ANALYSIS

**Testbed:** Since we propose full-surround MOT using vision sensors, we use a testbed comprising of 8 outside looking RGB cameras (seen in Figure 1). This setup ensures full surround coverage of the scene around the vehicle, while retaining a sufficient overlap between adjacent camera views. Frames captured from these cameras along with annotated surround vehicles are shown in Figure 1. In addition to full vision coverage, the testbed has full-surround Radar and LiDAR FoVs. Despite the final goal of this study being full-surround MOT, we additionally consider cases where only a subset of the vision sensors are used to illustrate the modularity of the

approach. More details on the sensors, their synchronization and calibration, and the testbed can be found in [43].

**Dataset:** To train and test our 3D MOT system, we collect a set of four sequences, each 3–4 minutes long, comprising of multi-camera videos and LiDAR point clouds using our testbed described above. The sequences are chosen much longer than traditional MOT sequences so that long range maneuvers of surround vehicles can be tracked. This is very crucial to autonomous driving. We also annotate all vehicles in the 8 camera videos for each sequence with their bounding box, as well as track IDs. It should be noted that each unique vehicle in the scene is assigned the same ID in all camera views. With these sequences set up, we use one sequence for training our tracker, and reserve the rest for testing. All our results are reported on the entire test set.

**Evaluation Metrics:** We use multiple metrics to evaluate the multiple object tracking performance as suggested by the MOT Benchmark [38]. Specifically, we use the 3D MOT metrics described in [39]. These include Multiple Object Tracking Accuracy (MOTA), Multiple Object Tracking Precision (MOTP), Mostly Track targets (MT, percentage of ground truth objects whose trajectories are covered by the tracking output for at least 80%), Mostly Lost targets (ML, percentage of ground truth objects whose trajectories are covered by the tracking output less than 20%), and the total number of ID Switches (IDS). In addition to listing the metrics in Table III, we also draw arrows next to each of them indicating if a high ( $\uparrow$ ) or low ( $\downarrow$ ) value is desirable. Finally, we provide top-down visualizations of the tracking results in a global coordinate system centered on the ego-vehicle for qualitative evaluation.

### A. Experimenting with Number of Cameras

As our approach to tracking is designed to be extremely modular, we test our tracker with different camera configurations. We experiment with 2, 3, 4, 6 and 8 cameras respectively. Top-down visualizations of the generated tracks for a test sequence are depicted in Figure 6. The ground truth tracks are provided for visual comparison. As can be seen, the tracker provides consistent results in its FoV irrespective of the camera configuration used, even if the cameras have no overlap between them.

The quantitative results on the test set for each camera configuration are listed in Table III. It must be noted that the tracker for each configuration is scored only based on the ground truth tracks visible in that camera configuration. The tracker is seen to score very well on each metric, irrespective of the number of cameras used. This illustrates the robustness of the M<sup>3</sup>OT framework. More importantly, it is seen that our tracker performs exceptionally well in the MT and ML metrics, especially in camera configurations with overlapping FoVs. Even though our test sequences are about 3 minutes long in duration, the tracker mostly tracks more than 70% of the targets, while mostly losing only a few. This demonstrates that our M<sup>3</sup>OT framework is capable of long-term target tracking.

### B. Effect of Projection Scheme

Figure 7 depicts the tracking results for a test sequence using the two projection schemes proposed. It is obvious that

TABLE III: Quantitative results showing ablative analysis of our proposed tracker.

Criteria for Comparison	Tracker Variant	Sensor Configuration		MOT Metrics [38], [39]				
		# of Cameras	Range Sensors	MOTA ( $\uparrow$ )	MOTP ( $\downarrow$ )	MT ( $\uparrow$ )	ML ( $\downarrow$ )	IDS ( $\downarrow$ )
<b>Number of Cameras Used</b> (Section V-A)	-	2	✓	73.38	0.03	71.36%	16.13%	16
	-	3	✓	77.26	0.03	77.34%	14.49%	38
	-	4	✓	72.81	0.05	72.48%	20.76%	49
	-	4 $\uparrow$	✓	74.18	0.05	74.10%	18.18%	45
	-	6	✓	79.06	0.04	79.66%	11.93%	51
	-	8	✓	75.10	0.04	70.37%	14.07%	59
<b>Projection Scheme</b> (Section V-B)	Point cloud based projection	8	✓	<b>75.10</b>	<b>0.04</b>	<b>70.37%</b>	<b>14.07%</b>	<b>59</b>
	IPM projection	8	✓(for fusion)	47.45	0.39	53.70%	19.26%	152
<b>Fusion Scheme</b> (Section V-C)	Point cloud based fusion	8	✓	<b>75.10</b>	<b>0.04</b>	<b>70.37%</b>	14.07%	<b>59</b>
	Distance based fusion	8	✓(for projection)	72.20	<b>0.04</b>	68.23%	<b>12.23%</b>	65
<b>Sensor Modality</b> (Sections V-B,V-C)	Cameras+LiDAR	8	✓	<b>75.10</b>	<b>0.04</b>	<b>70.37%</b>	<b>14.07%</b>	<b>59</b>
	Cameras	8	✗	40.98	0.40	50.00%	27.40%	171
<b>Vehicle Detector</b> (Section V-D)	RefineNet [40]	8	✓	<b>75.10</b>	<b>0.04</b>	<b>70.37%</b>	<b>14.07%</b>	<b>59</b>
	RetinaNet [41]	8	✓	73.89	0.05	68.37%	17.07%	72
	SubCat [42]	8	✓	69.93	<b>0.04</b>	66.67%	22.22%	81
<b>Global Position based Features</b> (Section V-E)	with $\{\phi_{13}, \phi_{14}\}$	8	✓	<b>75.10</b>	<b>0.04</b>	<b>70.37%</b>	<b>14.07%</b>	<b>59</b>
	without $\{\phi_{13}, \phi_{14}\}$	8	✓	71.32	0.05	64.81%	17.78%	88

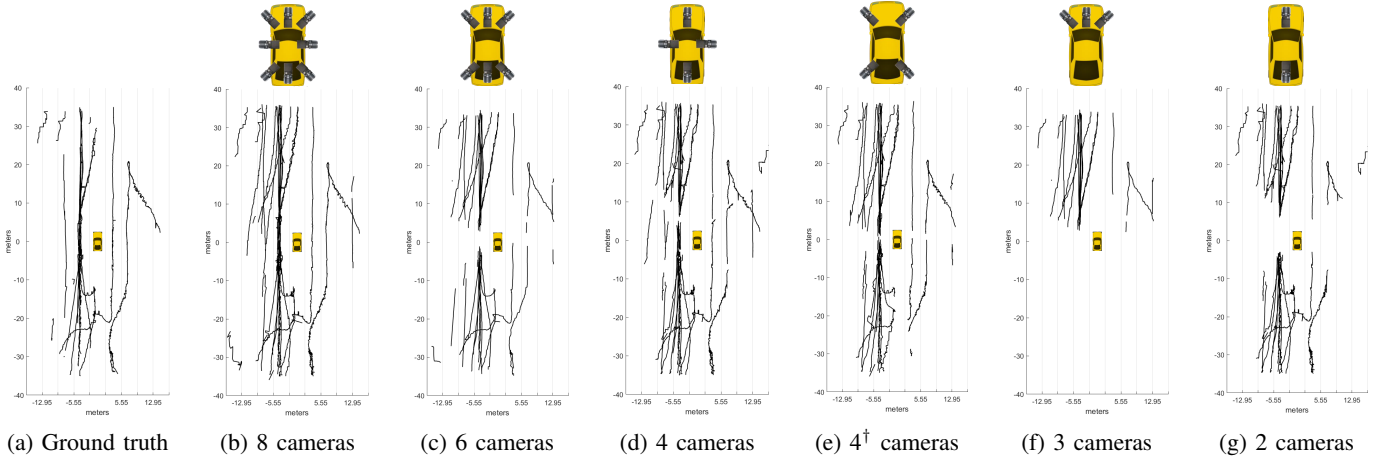


Fig. 6: Tracking results with different number of cameras. The camera configuration used is depicted above each result.

LiDAR based projection results in much better localization in 3D, which leads to more stable tracks and fewer fragments. The IPM based projection scheme is very sensitive to small changes in the input domain, and this leads to considerable errors during gating and data association. This phenomenon is verified by the high MOTP value obtained for IPM based projection as listed in Table III.

### C. Effect of Fusion Scheme

Once again, we see that the LiDAR point cloud based fusion scheme is more reliable in comparison to the distance based approach, albeit this difference is much less noticeable when proposals are projected using LiDAR point clouds. The LiDAR based fusion scheme results in objects being tracked longer (across camera views), and more accurately. The distance based fusion approach on the other hand fails to associate certain proposals, which results in templates not being stored for new camera views, thereby cutting short the track as soon as the target exits the current view. This superiority is reflected

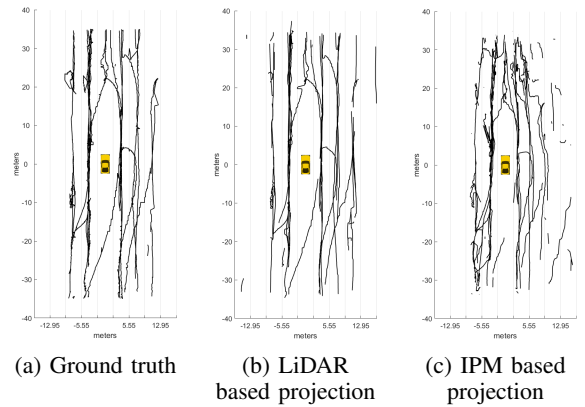


Fig. 7: Tracking results on a test sequence with different projection schemes.

in the quantitative results shown in Table III. The drawbacks of the distance based fusion scheme are exacerbated when using

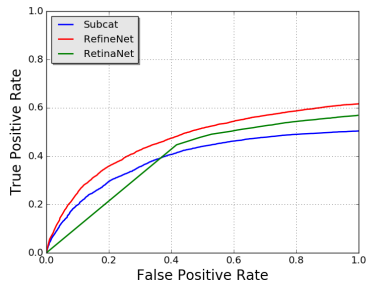


Fig. 8: ROC curves for different vehicle detectors on the 4 test sequences.

IPM to project proposals, reflected by the large drop in MOTA for a purely vision based system. This drop in performance is to be expected in the absence of LiDAR sensors. However, it must be noted that half the targets are still tracked for most of their lifetime, while only a quarter of the targets are mostly lost.

#### D. Effect of using Different Vehicle Detectors

Ideally, a tracking-by-detection approach should be detector agnostic. To observe how the tracking results change for different vehicle detectors, we ran the proposed tracker on vehicle detections obtained from three commonly used object detectors [40]–[42]. All three detectors were trained on the KITTI dataset [44] and have not seen examples from the proposed multi-camera dataset. The ROC curves for the detectors on the proposed dataset are shown in Figure 8. The corresponding tracking results for each detector are listed in Table III. Despite the sub-optimal performance of all three detectors in addition to significant differences in their ROC curves, the tracking results are seen to be relatively unaffected. This indicates that the tracker is less sensitive to errors made by the detector, and consistently manages to correct for it.

#### E. Effect of Global Position based Features

Table III indicates a clear benefit in incorporating features  $\{\phi_{13}, \phi_{14}\}$  for data association in Lost states. These features express how near/far a proposal is from the last known location of a target. This helps the tracker disregard proposals that are unreasonably far away from the latest target location. Introduction of these features leads to an improvement in all metrics and therefore justifies their inclusion.

### VI. CONCLUDING REMARKS

In this work, we have described a full-surround camera and LiDAR based approach to multi-object tracking for autonomous vehicles. To do so, we extend a 2D MOT approach based on the tracking-by-detection framework, and make it capable of tracking objects in the real world. The proposed M<sup>3</sup>OT framework is also made highly modular so that it is capable of working with any camera configuration with varying FoVs, and also with or without LiDAR sensors. An efficient and fast early fusion scheme is adopted to handle

object proposals from different sensors within a calibrated camera array. We conduct extensive testing on naturalistic full-surround vision and LiDAR data collected on highways, and illustrate the effects of different camera setups, fusion schemes and 2D-to-3D projection schemes, both qualitatively and quantitatively. Results obtained on the dataset support the modular nature of our framework, as well as its ability to track objects for a long duration. In addition to this, we believe that the M<sup>3</sup>OT framework can be used to test the utility of any camera setup, and make suitable modifications thereof to ensure optimum coverage from vision and range sensors.

### VII. ACKNOWLEDGMENTS

We would like to thank the Toyota Collaborative Safety Research Center (CSRC) for their generous and continued support. We would also like to thank our colleagues Kevan Yuen, Nachiket Deo, Ishan Gupta and Borhan Vasli at the Laboratory for Intelligent and Safe Automobiles (LISA), UC San Diego for their useful inputs and help in collecting and labeling the dataset.

### REFERENCES

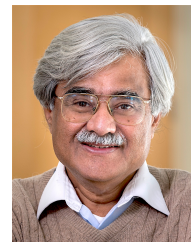
- [1] Y. Xiang, A. Alahi, and S. Savarese, “Learning to track: Online multi-object tracking by decision making,” in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 4705–4713.
- [2] J. Choi, S. Ulbrich, B. Lichte, and M. Maurer, “Multi-target tracking using a 3d-lidar sensor for autonomous vehicles,” in *Intelligent Transportation Systems (ITSC), 2013 16th International IEEE Conference on*. IEEE, 2013, pp. 881–886.
- [3] S. Song, Z. Xiang, and J. Liu, “Object tracking with 3d lidar via multi-task sparse learning,” in *Mechatronics and Automation (ICMA), 2015 IEEE International Conference on*. IEEE, 2015, pp. 2603–2608.
- [4] A. Asvadi, P. Peixoto, and U. Nunes, “Detection and tracking of moving objects using 2.5 d motion grids,” in *Intelligent Transportation Systems (ITSC), 2015 IEEE 18th International Conference on*. IEEE, 2015, pp. 788–793.
- [5] A. Asvadi, C. Premebida, P. Peixoto, and U. Nunes, “3d lidar-based static and moving obstacle detection in driving environments: An approach based on voxels and multi-region ground planes,” *Robotics and Autonomous Systems*, vol. 83, pp. 299–311, 2016.
- [6] D. Pfeiffer and U. Franke, “Efficient representation of traffic scenes by means of dynamic stixels,” in *Intelligent Vehicles Symposium (IV), 2010 IEEE*. IEEE, 2010, pp. 217–224.
- [7] A. Broggi, S. Cattani, M. Patander, M. Sabbatelli, and P. Zani, “A full-3d voxel-based dynamic obstacle detection for urban scenario using stereo vision,” in *Intelligent Transportation Systems (ITSC), 2013 16th International IEEE Conference on*. IEEE, 2013, pp. 71–76.
- [8] A. Vataavu, R. Danescu, and S. Nedevschi, “Stereovision-based multiple object tracking in traffic scenarios using free-form obstacle delimiters and particle filters,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 1, pp. 498–511, 2015.
- [9] A. Osep, A. Hermans, F. Engelmann, D. Klostermann, M. Mathias, and B. Leibe, “Multi-scale object candidates for generic object tracking in street scenes,” in *Robotics and Automation (ICRA), 2016 IEEE International Conference on*. IEEE, 2016, pp. 3180–3187.
- [10] R. K. Satzoda, S. Lee, F. Lu, and M. M. Trivedi, “Vision based front & rear surround understanding using embedded processors,” *IEEE Transactions on Intelligent Vehicles*, 2017.
- [11] S. Sivaraman and M. M. Trivedi, “Dynamic probabilistic drivability maps for lane change and merge driver assistance,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 15, no. 5, pp. 2063–2073, 2014.
- [12] H. Cho, Y.-W. Seo, B. V. Kumar, and R. R. Rajkumar, “A multi-sensor fusion system for moving object detection and tracking in urban driving environments,” in *Robotics and Automation (ICRA), 2014 IEEE International Conference on*. IEEE, 2014, pp. 1836–1843.
- [13] A. Asvadi, P. Girão, P. Peixoto, and U. Nunes, “3d object tracking using rgb and lidar data,” in *Intelligent Transportation Systems (ITSC), 2016 IEEE 19th International Conference on*. IEEE, 2016, pp. 1255–1260.



- [14] M. Allodi, A. Broggi, D. Giaquinto, M. Patander, and A. Prioletti, "Machine learning in tracking associations with stereo vision and lidar observations for an autonomous vehicle," in *Intelligent Vehicles Symposium (IV), 2016 IEEE*. IEEE, 2016, pp. 648–653.
- [15] S.-H. Bae and K.-J. Yoon, "Robust online multi-object tracking based on tracklet confidence and online discriminative appearance learning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1218–1225.
- [16] S. Kim, S. Kwak, J. Feyereisl, and B. Han, "Online multi-target tracking by large margin structured learning," in *Asian Conference on Computer Vision*. Springer, 2012, pp. 98–111.
- [17] C.-H. Kuo, C. Huang, and R. Nevatia, "Multi-target tracking by on-line learned discriminative appearance models," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 685–692.
- [18] Y. Li, C. Huang, and R. Nevatia, "Learning to associate: Hybridboosted multi-target tracker for crowded scene," in *Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on*. IEEE, 2009, pp. 2953–2960.
- [19] X. Song, J. Cui, H. Zha, and H. Zhao, "Vision-based multiple interacting targets tracking via on-line supervised learning," in *European Conference on Computer Vision*. Springer, 2008, pp. 642–655.
- [20] B. Babenko, M.-H. Yang, and S. Belongie, "Robust object tracking with online multiple instance learning," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 8, pp. 1619–1632, 2011.
- [21] C. Bao, Y. Wu, H. Ling, and H. Ji, "Real time robust l1 tracker using accelerated proximal gradient approach," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 1830–1837.
- [22] S. Hare, S. Golodetz, A. Saffari, V. Vineet, M.-M. Cheng, S. L. Hicks, and P. H. Torr, "Struck: Structured output tracking with kernels," *IEEE transactions on pattern analysis and machine intelligence*, vol. 38, no. 10, pp. 2096–2109, 2016.
- [23] Z. Kalal, K. Mikolajczyk, and J. Matas, "Tracking-learning-detection," *IEEE transactions on pattern analysis and machine intelligence*, vol. 34, no. 7, pp. 1409–1422, 2012.
- [24] J. V. Dueholm, M. S. Kristoffersen, R. K. Satzoda, T. B. Moeslund, and M. M. Trivedi, "Trajectories and maneuvers of surrounding vehicles with panoramic camera arrays," *IEEE Transactions on Intelligent Vehicles*, vol. 1, no. 2, pp. 203–214, 2016.
- [25] S. Sivaraman and M. M. Trivedi, "Looking at vehicles on the road: A survey of vision-based vehicle detection, tracking, and behavior analysis," *IEEE Transactions on Intelligent Transportation Systems*, vol. 14, no. 4, pp. 1773–1795, 2013.
- [26] J. Berclaz, F. Fleuret, E. Turetken, and P. Fua, "Multiple object tracking using k-shortest paths optimization," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 9, pp. 1806–1819, 2011.
- [27] A. A. Butt and R. T. Collins, "Multi-target tracking by lagrangian relaxation to min-cost network flow," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1846–1853.
- [28] A. Milan, S. Roth, and K. Schindler, "Continuous energy minimization for multitarget tracking," *IEEE transactions on pattern analysis and machine intelligence*, vol. 36, no. 1, pp. 58–72, 2014.
- [29] J. C. Niebles, B. Han, and L. Fei-Fei, "Efficient extraction of human motion volumes by tracking," in *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*. IEEE, 2010, pp. 655–662.
- [30] H. Pirsiavash, D. Ramanan, and C. C. Fowlkes, "Globally-optimal greedy algorithms for tracking a variable number of objects," in *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*. IEEE, 2011, pp. 1201–1208.
- [31] L. Zhang, Y. Li, and R. Nevatia, "Global data association for multi-object tracking using network flows," in *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*. IEEE, 2008, pp. 1–8.
- [32] 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, 2005.
- [33] S. Oh, S. Russell, and S. Sastry, "Markov chain monte carlo data association for multi-target tracking," *IEEE Transactions on Automatic Control*, vol. 54, no. 3, pp. 481–497, 2009.
- [34] K. Okuma, A. Taleghani, N. d. Freitas, J. J. Little, and D. G. Lowe, "A boosted particle filter: Multitarget detection and tracking," *Computer Vision-ECCV 2004*, pp. 28–39, 2004.
- [35] J. Munkres, "Algorithms for the assignment and transportation problems," *Journal of the society for industrial and applied mathematics*, vol. 5, no. 1, pp. 32–38, 1957.
- [36] M. D. Breitenstein, F. Reichlin, B. Leibe, E. Koller-Meier, and L. Van Gool, "Online multiperson tracking-by-detection from a single, uncalibrated camera," *IEEE transactions on pattern analysis and machine intelligence*, vol. 33, no. 9, pp. 1820–1833, 2011.
- [37] A. Rangesh and M. M. Trivedi, "Ground plane polling for 6dof pose estimation of objects on the road," *arXiv preprint arXiv:1811.06666*, 2018.
- [38] K. Bernardin and R. Stiefelhagen, "Evaluating multiple object tracking performance: the clear mot metrics," *EURASIP Journal on Image and Video Processing*, vol. 2008, no. 1, pp. 1–10, 2008.
- [39] A. Milan, K. Schindler, and S. Roth, "Challenges of ground truth evaluation of multi-target tracking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2013, pp. 735–742.
- [40] R. N. Rajaram, E. Ohn-Bar, and M. M. Trivedi, "Refinenet: Refining object detectors for autonomous driving," *IEEE Transactions on Intelligent Vehicles*, vol. 1, no. 4, pp. 358–368, 2016.
- [41] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," *arXiv preprint arXiv:1708.02002*, 2017.
- [42] E. Ohn-Bar and M. M. Trivedi, "Fast and robust object detection using visual subcategories," in *Computer Vision and Pattern Recognition Workshops (CVPRW), 2014 IEEE Conference on*. IEEE, 2014, pp. 179–184.
- [43] A. Rangesh, K. Yuen, R. K. Satzoda, R. N. Rajaram, P. Gunaratne, and M. M. Trivedi, "A multimodal, full-surround vehicular testbed for naturalistic studies and benchmarking: Design, calibration and deployment," *arXiv preprint arXiv:1709.07502*, 2017.
- [44] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision meets robotics: The kitti dataset," *The International Journal of Robotics Research*, vol. 32, no. 11, pp. 1231–1237, 2013.



**Akshay Rangesh** is currently working towards his PhD in electrical engineering from the University of California at San Diego (UCSD), with a focus on intelligent systems, robotics, and control. His research interests span computer vision and machine learning, with a focus on object detection and tracking, human activity recognition, and driver safety systems in general. He is also particularly interested in sensor fusion and multi-modal approaches for real time algorithms.



**Mohan Manubhai Trivedi** is a Distinguished Professor at University of California, San Diego (UCSD) and the founding director of the UCSD LISA: Laboratory for Intelligent and Safe Automobiles, winner of the IEEE ITSS Lead Institution Award (2015). Currently, Trivedi and his team are pursuing research in intelligent vehicles, machine perception, machine learning, human-robot interactivity, driver assistance, active safety systems. Three of his students have received "best dissertation" recognitions. Trivedi is a Fellow of IEEE, ICPR and SPIE. He received the IEEE ITS Society's highest accolade "Outstanding Research Award" in 2013. Trivedi serves frequently as a consultant to industry and government agencies in the USA and abroad.