

ODTFormer: Efficient Obstacle Detection and Tracking with Stereo Cameras Based on Transformer

Tianye Ding^{1*}, Hongyu Li^{2*}, and Huaizu Jiang¹

Abstract—Obstacle detection and tracking represent a critical component in robot autonomous navigation. In this paper, we propose ODTFormer, a Transformer-based model that addresses both obstacle detection and tracking problems. For the detection task, our approach leverages deformable attention to construct a 3D cost volume, which is decoded progressively in the form of voxel occupancy grids. We further track the obstacles by matching the voxels between consecutive frames. The entire model can be optimized in an end-to-end manner. Through extensive experiments on DrivingStereo and KITTI benchmarks, our model achieves state-of-the-art performance in the obstacle detection task. We also report comparable accuracy to state-of-the-art obstacle tracking models while requiring only a fraction of their computation cost, typically ten-fold to twenty-fold less. Our code is available on <https://github.com/neu-vi/ODTFormer>.

I. INTRODUCTION

Obstacle detection and tracking represents a safety-critical challenge across various domains, including robot autonomous navigation [1–5] and self-driving vehicles [6–9]. For instance, a service robot needs to detect people and pillars surrounding it, track their motions (if any), or even predict their future trajectories to avoid collision. Accurate obstacle detection and tracking are crucial components of autonomous navigation systems, particularly in state-based frameworks, to ensure collision-free navigation [10–13]. Recent research efforts focus on using low-cost visual sensors for obstacle perception to improve affordability [10, 14–17] compared with costly ones (*e.g.*, LiDAR). In this paper, we concentrate on a specific line of research employing stereo cameras, which offer higher 3D perception accuracy, extended sensing range, and enhanced agility for robots compared to monocular-based systems [18, 19].

Previous stereo-based obstacle detection approaches have mostly relied on depth estimation modules [10, 14, 17]. These approaches involve the estimation of a depth map followed by its transformation into a point cloud or a voxel grid. However, this two-stage approach often necessitates tradeoffs between speed and accuracy. Li et al. [14] addressed this tradeoff by introducing StereoVoxelNet, which directly estimates voxel grids from stereo images by constructing a pixel-wise cost volume and employing a 2D-3D encoder-decoder network structure for efficient and accurate voxel

This research is supported by the National Science Foundation under Award Number IIS-2310254.

*Equal contribution.

¹Northeastern University, Boston, MA, 02115. {ding.tian, h.jiang}@northeastern.edu

²Brown University, Providence, RI, 02912. hli230@cs.brown.edu

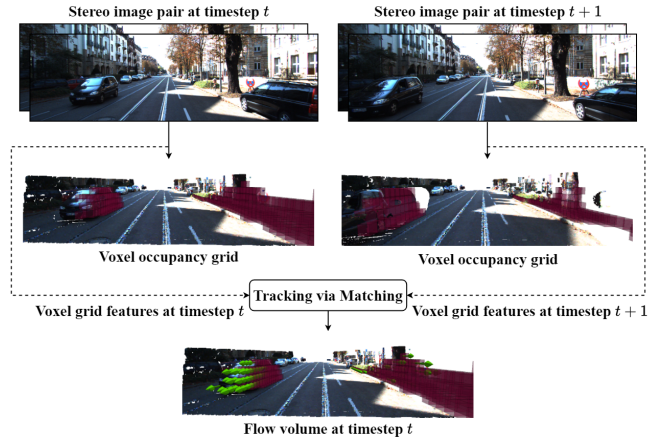


Fig. 1: We propose ODTFormer for joint obstacle detection and tracking using stereo cameras. We first detect obstacles in the form of occupancy grids at each time step and match them across two consecutive frames for tracking. We can see here that our model can successfully detect all obstacles and accurately track them. The obstacle detection results are shown as red cubes, and the tracking results are marked as green arrows. Longer arrows indicate large motion magnitude.

grid reconstruction. However, a notable limitation of StereoVoxelNet is its *implicit* incorporation of camera parameters within the neural network parameters. As a result, the model can only work with specific camera parameters and image resolution, limiting its applicability across datasets with varying resolutions and camera parameters.

Moreover, the assumption of static obstacles presents challenges in dynamic environments, such as social navigation, where obstacles may be randomly moving pedestrians. Numerous studies have addressed the visual dynamic obstacle tracking problem, employing traditional methods such as the Kalman filter to estimate the velocity [11, 13], associating point cloud centroids [10, 13] or bounding boxes [12] from consecutive frames using distance thresholds. However, these techniques often require meticulous parameter tuning, leading to a lack of robustness. Stereo-based scene flow estimation addresses simultaneous 3D structure and motion estimation [20], which could be adapted for obstacle tracking. Nonetheless, these methods estimate a dense motion vector map and depth map for each pixel, resulting in high computational costs and runtime [21, 22].

In this paper, we present ODTFormer, a Transformer-based [23] model, to address the aforementioned obstacle

detection and tracking challenges. Unlike existing stereo-based models, which construct a pixel-wise cost volume for disparity matching [14, 24–27], we propose to use deformable cross-attention [28] from 3D voxel queries to 2D stereo image features to compute the matching cost. Compared with the pixel-wise cost volume used in [14], we directly construct it in the 3D space, conforming better to the scene geometry. More importantly, our approach disentangles dataset specifics from the model design and thus shows better generalization than [14]. The cost volume is then processed by a U-Net decoder to produce occupancy voxel grids progressively for efficiency purposes.

To account for environmental dynamics, we introduce a novel obstacle tracker to capture the motions of the scene by matching similar voxels between two consecutive frames. We integrate physical constraints into the voxel tracking by setting a volumetric bound for each voxel when searching for its corresponding voxel in the next frame, which improves both accuracy and efficiency. On the one hand, compared with traditional approaches, our model has the learning capacity that leverages the deep voxel feature representations to track obstacles accurately. Moreover, both detection and tracking modules can be jointly optimized, leading to high accuracy. On the other hand, compared with scene flow methods, our approach tracks the sparse voxels instead of densely over pixels, leading to lightweight computation with higher efficiency. Our entire model can detect and track obstacles efficiently, resulting in as much as 20 times fewer computations compared to the current state-of-the-art in scene flow estimation [22].

To validate the effectiveness of our approach, we conduct comprehensive evaluations on DrivingStereo [29] and KITTI [30] datasets. In the obstacle detection task, we demonstrate a significant performance improvement compared to prior works [14, 24–27]. Additionally, our approach exhibits greater generalization across various camera parameters and resolutions than [14]. In the obstacle tracking task, we compare scene flow-based approaches and show that our approach can achieve comparable results to state-of-the-art methods while needing much lighter computation cost. Detailed ablation studies are conducted to validate different design choices. In summary, our contributions are:

- A novel 3D cost volume construction method based on the deformable cross-attention [28], which better conforms to the scene geometry and generalizes well to different camera parameters and image resolutions.
- A novel obstacle tracking method by matching the voxels across two consecutive frames, which can be jointly optimized with the detection module in an end-to-end manner.
- Experimental results show that our approach achieves better or comparable detection and tracking accuracy to state-of-the-art methods while being efficient and running fast at the same time (20fps on the KITTI resolution of 370×1224 on an RTX A5000 GPU without careful postprocessing, *e.g.*, quantization, pruning).

II. RELATED WORKS

In this section, we introduce two related areas to our work. We first introduce the general background of stereo vision, including depth and scene flow estimation. We then survey the obstacle perception works from the perspectives of detection and tracking.

A. Stereo Vision

Stereo vision is a classical computer vision problem that aims to extract 3D information from a pair of RGB cameras. Typically, the transformation between left and right cameras is considered rigid and known. We briefly introduce two relevant applications of stereo vision: depth estimation and scene flow estimation.

Depth estimation involves deriving a depth map from one image of the stereo pair, typically the left one, associating each pixel with a depth value. This is performed by disparity matching between the stereo pair [33]. Recent approaches employ end-to-end deep neural networks and cost volume techniques for matching [14, 24–27].

Scene flow estimation aims to understand the dynamic scenes by extracting the 3D information (disparity) and the motion simultaneously from sequences of stereo images [21, 22, 34]. Unlike the traditional optical flow problem, scene flow estimation provides a denser correspondence between consecutive frames. Stereo-based methods typically separate the 3D flow problem into depth and flow estimation. Some works propose simultaneous estimation of 2D optical flow and disparity maps [21], while others [22] leverage off-the-shelf depth estimators.

Despite their impressive performance, these methods often come with significant computational costs and are, therefore, unsuitable for real-time robotic systems due to the need for dense 2D estimation. Our approach leverages the sparsity of these tasks and voxel grids, a common representation in robotics, significantly reducing computational costs while maintaining accuracy comparable to depth and scene flow methods.

B. Obstacle Perception

Obstacle perception is a fundamental aspect of robotics, playing a crucial role in enabling autonomous robots to navigate among environments safely and effectively. In this paper, we consider two aspects of obstacle perception: obstacle detection and obstacle tracking.

In the obstacle detection task, we detect obstacles in the form of a voxel occupancy grid. Conventionally, stereo-based obstacle detection relies on filtering and voxelizing point clouds obtained by depth estimation module [10, 35] but suffers from high computation complexity. StereoVoxelNet [14] directly infers occupancy grids from stereo pairs but struggles with varying resolutions and camera parameters.

Obstacle tracking aims at estimating the motion of obstacles across time. While most of the previous works are based on LiDAR point cloud [36, 37] or camera-LiDAR fusion [38], visual end-to-end methods, especially stereo-based and occupancy-based ones, remain under-explored. Existing

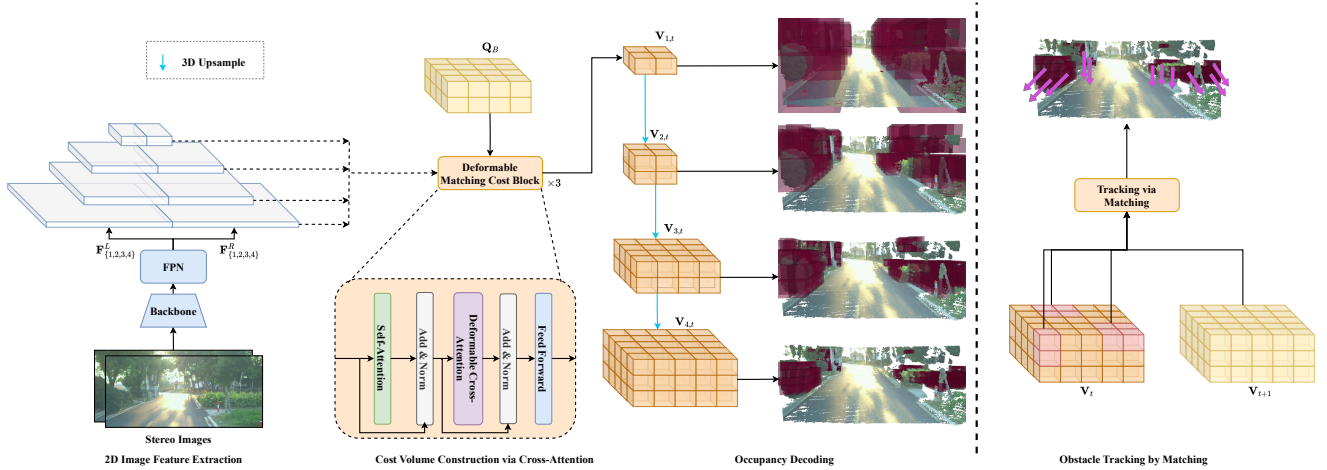


Fig. 2: **Illustration of the overall architecture design.** **Left:** For obstacle detection, we first extract multi-scale 2D feature maps [31, 32] for each of the stereo images. We then encode the voxels in the ROI to cross-attend to the image features to compute the matching cost through our novel cost volume construction method. Such a cost volume is directly constructed in the 3D space, which conforms better to the scene geometry, disentangles dataset specifics from model design, and thus generalizes well. It is then progressively decoded into occupancy voxel grids. **Right:** For obstacle tracking, we cast it as a matching problem by finding the correspondences of voxels across two consecutive frames, where we incorporate physical constraints to improve both the accuracy and efficiency. Both the detection and tracking modules can be jointly optimized in an end-to-end manner and run efficiently.

visual approaches rely on object detectors or Kalman filters [10–12], limiting generalizability and requiring meticulous parameter tuning. We introduce voxel flow, which is the 3D grid-to-grid voxel-level motion of occupied voxels while achieving near real-time obstacle tracking.

In parallel to our work, there exist seminal occupancy detection approaches tailored for autonomous driving [7, 8]. These approaches commonly operate under the assumption of having access to images from the surrounding views. In contrast, our approach utilizes only a pair of stereo images and emphasizes computational efficiency.

III. PROPOSED METHOD

In this paper, we present a novel approach to tackle both obstacle detection and tracking problems. For the obstacle detection task, our model estimates a voxel occupancy grid in front of the stereo cameras, defining this area as our region of interest (ROI). Subsequently, our model efficiently tracks the motion of each detected obstacle by matching voxels across two consecutive frames.

A. Overview

Obstacle detection. Following StereoVoxelNet [14], we formulate obstacle detection as occupancy detection based on stereo images. As shown in Fig. 2 (a), our detection model consists of three modules: i) 2D image feature extraction, ii) cost volume construction via cross-attention, and iii) occupancy decoding.

1) *2D image feature extraction:* Given the rectified stereo images \mathbf{I}^L and \mathbf{I}^R , we employ a lightweight EfficientNet-B0 [31] backbone enhanced by the FPN [32] to extract multi-scale image features, denoted as \mathbf{F}_i^L and \mathbf{F}_i^R , respectively.

Here the subscript $i \in \{1, 2, 3, 4\}$ indicates the level (scale) of the feature maps unless explicitly mentioned otherwise. With the input image resolution denoted as $H \times W$, the multi-scale feature maps reduce the resolutions to $\{\frac{1}{4}, \frac{1}{8}, \frac{1}{16}, \frac{1}{32}\}$, with the feature channel dimension of D for all the levels.

2) *Cost volume construction via cross-attention:* A crucial difference of our proposed approach from [14] is our novel cost volume construction module design, where we employ cross-attention from 3D voxel queries to 2D stereo image features. Compared with the pixel-wise cost volume used in [14], we directly construct it in the 3D space, conforming better to the scene geometry. More importantly, our approach disentangles dataset specifics from the model design and thus shows better generalization than [14]. We will introduce this module in detail in Sec. III-B.

3) *Occupancy decoding:* Similar to [14], as the final step in the obstacle detection module, we estimate the voxel occupancy in a coarse-to-fine-grained manner. We utilize a 3D U-Net decoder to upsample the cost volume gradually. Specifically, in the time step t , our model outputs four-level coarse-to-fine voxel occupancy volumes $\mathbf{V}_{i,t}$ with the constant ROI size of $s_x \times s_y \times s_z$ meters,

$$\mathbf{V}_{i,t} = \{0, 1\}^{2^{i-1}(n_x \times n_y \times n_z)}, i \in \{1, 2, 3, 4\}, \quad (1)$$

where i denotes the level number. n_x , n_y , and n_z denote the number of voxels in the x , y and z axes, respectively. Each voxel is classified as either empty (denoted by 0) or occupied (denoted by 1) and has a side length of l_i which is decreased by half when the level number increases to keep grid metric size $s_{\{x,y,z\}}$ constant, $l_{i+1} = \frac{1}{2}l_i$.

Obstacle tracking by matching. After getting obstacle detection results in the form of occupied voxels, we’d like to track each of them for a better understanding of scene dynamics. To this end, we find correspondences of occupied voxels across two consecutive time steps via *global matching*, as shown in Fig. 2 (the right part). Such tracking information provides critical cues of dynamic objects for visual navigation. We introduce this part in detail in Sec. III-C.

Finally, we describe how the entire model is trained in Sec. III-D.

B. Cost Volume Construction via Cross-Attention

Cost volume plays a critical role for stereo-based obstacle detection methods [14, 24–27]. Conventional cost volume construction involves computing *pixel-wise* matching costs across a predefined range of disparity/depth levels [14]. However, this approach encounters significant drawbacks. First of all, the dimensions of the cost volume do not directly correlate with the resolution of the voxel grid. To solve this issue, the cost volume is compressed into a single latent vector and subsequently resized into the desired ROI voxel volume in [14]. This process inevitably leads to the loss of geometry information, which is critical for accurate obstacle detection. Furthermore, the hallucination of the latent vector to a 3D voxel volume assumes a fixed ROI dimension, prohibiting its generalization to a different setting¹. In addition, the camera intrinsics and extrinsics are directly embedded into the construction of the cost volume. Such entanglement of the dataset specifics and model design restricts its applicability across different domains and setups, making generalization challenging, especially in scenarios with varying camera parameters and image resolutions.

Instead of constructing a cost volume at the pixel level, we propose a novel approach to build it in the 3D space directly via cross-attention from voxel queries to image features. In a nutshell, we first encode each voxel in the ROI, which will then attend to image features via deformable attention [28]. Compared with the conventional cost volume [14], our approach better conforms to the scene geometry and disentangles dataset specifics and model design, allowing better generalization. We explain different components in detail as follows. For brevity, we omit the time step index t in this section.

Encoding of voxels in the ROI. We partition the ROI into a set of voxels \mathbf{Q}_B with a dimension of $n_x \times n_y \times n_z$ to balance the computation burden and spatial resolution of the cost volume. Each voxel in \mathbf{Q}_B is encoded as the Fourier positional encoding [23] of its normalized centroid coordinates (between 0 and 1), which is fed into a multi-layer perceptron (MLP) to get the encoding \mathbf{q} . It is further enhanced by the pixel-aligned features \mathbf{f} , $\mathbf{q} = \mathbf{q} + \mathbf{f}$ as in [39]. \mathbf{f} is defined as

$$\mathbf{f} = \frac{1}{2}(\mathbf{F}_4^L(\mathcal{P}(p, \theta^L)) + \mathbf{F}_4^R(\mathcal{P}(p, \theta^R))), \quad (2)$$

¹It is in analogy to the existence of the fully connected layers in AlexNet. As a result, the model only accepts an input image of a fixed dimension.

where p is the voxel’s centroid. $\mathcal{P}(p, \theta^L)$ and $\mathcal{P}(p, \theta^R)$ denote the projections of the voxel centroid onto the left and right images, respectively, according to their projection matrices θ^L and θ^R . Bilinear interpolation is used to sample features from the image feature maps \mathbf{F}_4^L and \mathbf{F}_4^R .

Deformable cross-attention for matching cost. We leverage the deformable matching cost (DMC) block to update voxel queries with calculated matching costs for each query. Our DMC block is built on top of the transformer decoder block [23] with deformable attention [28].

Specifically, for the voxel centered at p , we first sample points on multi-scale feature maps $(\mathbf{F}_i^L, \mathbf{F}_i^R)$ (similar to the pixel-aligned features \mathbf{f} but on multi-scale features instead of just the last level) around the projected points of the voxel centroid p with the offset δ , respectively. δ is generated by a learnable 3D offset sampler g [28] based on the voxel’s encoding \mathbf{q} , $\delta = g(\mathbf{q})$, where g is an MLP plus a sigmoid output layer (scaled by the 3D voxel size of the coarsest output, see experiment settings in Sec. IV-A). We then directly concatenate the sampled features along the channel dimension and employ a scale-specific MLP to compute the cross-view matching cost \mathbf{c}_i :

$$\mathbf{c}_i(p, \mathbf{q}, \mathbf{F}_i) = \text{MLP}_i(\mathbf{F}_i^L(\mathcal{P}(p + \delta, \theta^L)) \oplus \mathbf{F}_i^R(\mathcal{P}(p + \delta, \theta^R))), \quad (3)$$

where \oplus denotes the feature concatenation. MLP_i has two linear layers with batch normalization and ReLU activation in between. Since we project the same 3D sample point onto multi-scale feature maps, the resultant multi-scale matching costs for this 3D point can be aggregated by concatenating along the channel dimension and employing an MLP to compute the voxel matching cost \mathbf{c} :

$$\mathbf{c}(p, \mathbf{q}, \mathbf{F}_{\{1,2,3,4\}}) = \text{MLP}(\mathbf{c}_1(p, \mathbf{q}, \mathbf{F}_1) \oplus \mathbf{c}_2(p, \mathbf{q}, \mathbf{F}_2) \oplus \mathbf{c}_3(p, \mathbf{q}, \mathbf{F}_3) \oplus \mathbf{c}_4(p, \mathbf{q}, \mathbf{F}_4)), \quad (4)$$

where MLP has the same layer configuration as MLP_i . Instead of using a single offset, multiple ones are usually generated to enhance the computing of the matching cost, and different matching costs for different samples are aggregated as follows.

$$\mathbf{C}_B(p, \mathbf{q}, \mathbf{F}) = \sum_{s=1}^{N_s} \mathbf{A}_s \mathbf{W}_s \mathbf{c}(p, \delta_s, \mathbf{F}_{\{1,2,3,4\}}), \quad (5)$$

where N_s is the total number of sampling points. δ_s denotes the s -th learned offset. $\mathbf{A}_s \in [0, 1]$ is the learnable weights for the cost generation and $\mathbf{W}_s \in \mathbb{R}^{D \times D}$ are learnable parameters. We can get a cost volume $\mathbf{C}_B \in \mathbb{R}^{n_x \times n_y \times n_z \times D}$ by computing the matching cost for all voxels in \mathbf{Q}_B .

Deformable cross-attention has been used in previous works [6–9] to aggregate multi-view and multi-scale features. But instead of simply averaging them as in existing work, we leverage the geometric constraints between two stereo images and use the deformable attention to construct a cost volume in 3D, leading to better accuracy, as evidenced by our ablation studies.

The cost volume \mathbf{C} will be processed to output occupancy voxel grid $\mathbf{V}_{i,t}$ progressively using a U-Net decoder [14].

C. Voxel Tracking by Matching

Robots usually navigate in a dynamic environment by independently moving objects around (*e.g.*, pedestrians). Enabling a robot to track an object is essential to ensure safe navigation (*e.g.*, avoiding bumping into a person).

Given the estimated occupancy grid \mathbf{V}_t and \mathbf{V}_{t+1} at the time step t and $t + 1$, respectively², our goal is to find the 3D motion vector of each voxel in \mathbf{V}_t by finding its correspondence in \mathbf{V}_{t+1} . Apparently, we do not need to worry about unoccupied voxels. For the occupied ones, we compute their cosine similarities $\mathbf{S}_t \in [-1, 1]^{N_t \times N_{t+1}}$ by comparing their feature representations in the U-Net decoder right before the output layer. N_t and N_{t+1} denote the number of occupied voxels in \mathbf{V}_t and \mathbf{V}_{t+1} , respectively. Optionally, we can match all voxels in \mathbf{V}_{t+1} regardless of their occupancy status. We will ablate this design choice in the experiment section. The matching distribution of the occupied voxels can then be computed as

$$\mathbf{P}_t = \text{softmax}(\tau \mathbf{S}_t), \quad (6)$$

where τ is a learnable logit scale. The softmax is applied in a row-wise manner. Motion vectors \mathbf{M}_t of the occupied voxels can be computed as

$$\mathbf{M}_t = \mathbf{P}_t \mathbf{G}_{t+1} - \mathbf{G}_t, \quad (7)$$

where $\mathbf{G}_t \in \mathbb{Z}^{N_t \times 3}$ denotes the 3D coordinates of each voxel’s centroid in \mathbf{V}_t .

In practice, the motions of the occupied voxels are bounded by physical constraints. Instead of naively comparing all possible pairs of occupied voxels, we set a volumetric bound surrounding each occupied voxel in \mathbf{V}_t and only measure the similarities for voxels contained in this boundary when computing \mathbf{S}_t . For the voxels outside of the boundary, the similarity score is set to be $-\infty$. Specifically, assuming that a stereo image sequence is captured at least f frames per second (FPS) and the obstacles have a smaller relative velocity of than v meters per second to us, the maximum displacement of any trackable voxel across two frames is $d = \frac{v}{f}$. With the voxel at the center of the volumetric region, the volume bound’s dimensionality can be calculated as $(2\lceil \frac{d}{l_4} \rceil + 1) \times 3 \times (2\lceil \frac{d}{l_4} \rceil + 1)$ for x , y and z axes, where l_4 is the voxel side length in the last level of the U-Net decoder. The maximum displacements along the y -axis (pointing downward the ground plane) in both directions have been set to 1 to ignore drastic vertical movements within the driving scenario. As we will see in the experiment section, it not only eases the computation burden but also leads to better obstacle tracking accuracy.

²They are the output of the last layer in the U-Net decoder, $\mathbf{V}_{4,t}$ and $\mathbf{V}_{4,t+1}$. For brevity, we simply use \mathbf{V}_t and \mathbf{V}_{t+1} in this section when the context is clear.

TABLE I: **Quantitative obstacle detection results on the DrivingStereo testing set.** The best result is bold and the second-best result is underlined. We use an RTX A5000 GPU for measuring the inference speed.

Method	Level	CD (meters) ↓		IoU (%) ↑		MACs ↓	Params ↓	FPS ↑
		15.0m	30.0m	15.0m	30.0m			
2D-MSNet [24]	4	9.72	21.01	18.14	9.38	91.74G	<u>2.23M</u>	6.7
3D-MSNet [24]	4	5.96	13.14	21.11	11.96	414.59G	1.77M	3.6
ACVNet [27]	4	8.43	18.82	18.32	13.59	801.33G	7.11M	4.4
Lac-GwcNet [26]	4	7.19	<u>12.90</u>	29.84	17.25	777.19G	9.37M	4.8
CFNet [25]	4	9.79	18.65	21.21	10.35	456.14G	22.24M	4.9
StereoVoxelNet [14]	1	4.54	5.12	80.07	71.94			
	2	3.07	4.34	70.08	56.39			
	3	2.75	6.87	51.80	37.14	<u>16.03G</u>	4.58M	21.8
	4	3.15	17.20	38.27	21.92			
StereoVoxelNet w/ EfficientNet	1	4.36	5.17	80.26	71.93			
	2	2.81	4.28	70.55	56.32	14.83G	10.46M	20.2
	3	2.49	5.95	52.59	37.83			
	4	2.90	13.56	38.69	23.04			
ODTFormer (Ours)	1	3.09	4.13	85.25	76.75			
	2	1.92	3.22	76.10	62.69	25.05G	6.14M	<u>21.7</u>
	3	1.64	3.82	57.94	43.79			
	4	1.87	8.72	41.62	26.62			

D. Model Training

We optimize the obstacle detection model using the weighted sum of the Intersection of Union (IoU) loss \mathcal{L}_D over four U-Net decoder levels:

$$\mathcal{L}_D = \sum_{j=1}^4 w_j (1 - \text{IoU}(\mathbf{V}_{j,t}, \bar{\mathbf{V}}_{j,t})), \quad (8)$$

where $\bar{\mathbf{V}}_{j,t}$ denotes the ground-truth voxel grid at level j . w_j is the weight assigned to the U-Net decoder level, and we set w_1 to w_4 as [0.30, 0.27, 0.23, 0.20] empirically.

The tracker is trained using the endpoint error (EPE) loss, which is the average of element-wise L2 distances between all estimated and ground-truth 3D motion vectors. We note, however, that optimizing the loss for detected occupied voxels leads to degenerated performance because the detection results may be wrong. Instead, we define the loss for *all* voxels regardless of their detection results. The loss is defined as

$$\mathcal{L}_T = \frac{1}{N'_t} \sum_{k=1}^{N'_t} \|\mathbf{M}'_t[k] - \mathbf{M}_t[k]\|_2, \quad (9)$$

where \mathbf{M}' contains the motion vectors for all the voxels. For the k -th voxel, if it is detected as occupied, $\mathbf{M}'[k] = \mathbf{M}[k]$. Otherwise, $\mathbf{M}'[k] = \mathbf{0}$. $\bar{\mathbf{M}}$ denotes the ground-truth motion vectors.

The ODTFormer is trained progressively. We first train the detection module to obtain reliable detection results. We then jointly train the detection and tracking modules together, where the feature representations of the voxels are refined to encode the spatio-temporal cues effectively.

IV. EXPERIMENTS

In this section, we evaluate the performance of ODTFormer against various baselines on both obstacle detection (Sec. IV-A) and obstacle tracking (Sec. IV-B) tasks. In the end, we perform ablation studies on our design choices for both tasks (Sec. IV-C).

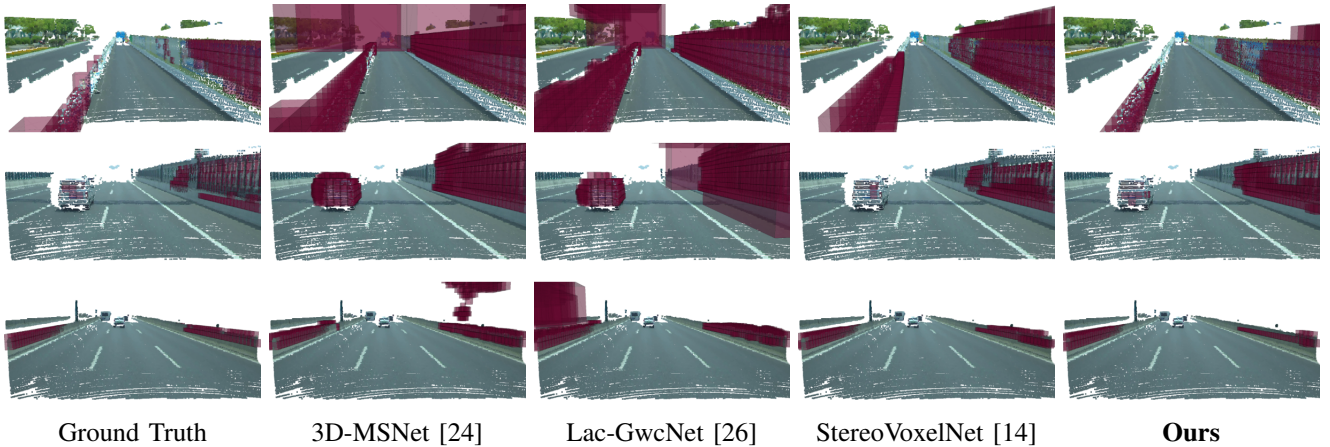


Fig. 3: Visual results of obstacle detection on the DrivingStereo dataset.

A. Obstacle Detection

We optimize our voxel occupancy prediction module using the DrivingStereo dataset [29], which contains over 180,000 stereo pairs for the driving scenario with a resolution of 400×880 . During preprocessing, we augment the training images using color jittering.

We set our ROI (defined in Sec. III-A.3) with respect to the left camera as $30m$ ahead, $8m$ to the left, $10m$ to the right, and $3m$ to the top and bottom. We set the voxel sizes to be $3m$, $1.5m$, $0.75m$, and $0.375m$ for each level, and the finest level will have a voxel grid of size $n_x \times n_y \times n_z = 48 \times 16 \times 80$. We precompute and store the precise voxel occupancy grid ground truth by cropping our ROI and removing all voxels that are $1.5m$ under the viewpoint on the y -axis (ground plane).

The network is trained for 20 epochs with a batch size of 16 using AdamW [40] optimizer. We set the initial learning rate as 0.0001 and decrease it each epoch using the Cosine Annealing learning rate scheduler with a minimum learning rate of 1×10^{-8} .

Evaluation metrics. We choose IoU and Chamfer Distance (CD) as our evaluation metrics. IoU measures the intersection between the ground truth and the predicted voxel occupancy grids. CD measures the distance between the point clouds (by presenting the centroids of occupied voxels as point clouds). We report the output accuracies with the $15.0m$ range (half of the ROI depth) and the $30.0m$ range (entire ROI depth) separately, which represents easier and harder-to-identify obstacles.

Besides the accuracy metrics, we also consider computational cost metrics: multiply-accumulate operations (MACs), the number of model parameters (Params), and frame-per-second (FPS).

Baseline methods. We compare obstacle detection performance against two types of seminal works: stereo-based depth estimation and voxel occupancy prediction. For depth estimation models, following the prior works [10, 14], we first convert the estimated depth map into a point cloud and subsequently voxelize it into a voxel occupancy grid.

TABLE II: Quantitative obstacle tracking results on the KITTI 2015 scene flow dataset. RAFT-3D MACs are evaluated including MobileStereoNet. The best result is bold, and the second-best result is underlined.

Method	EPE (meters) ↓	Foreground EPE (meters) ↓	MACs ↓
PWOC-3D [34]	0.043	1.23	N/A
SENSE [41]	0.048	1.18	284.73G
RAFT-3D (2D-MSNet) [22]	<u>0.025</u>	0.91	649.46G
RAFT-3D (3D-MSNet) [22]	0.028	0.87	1295.16G
ODTFormer (Ours)	0.021	1.08	64.47G

We include MobileStereoNet (MSNet, both 2D and 3D version) [24], ACVNet [27], Lac-GwcNet [26] and CFNet [25]. For the occupancy prediction approach, we compare against StereoVoxelNet [14]. To rule out the effect of feature extraction, we also ablate the original feature extraction backbone of StereoVoxelNet with EfficientNet (same as ours).

As shown in Tab. I, ODTFormer significantly outperforms all other approaches, particularly depth-based methods [24–27], in terms of both IoU and CD. Although ODTFormer exhibits higher MACs compared to StereoVoxelNet, we achieve similar FPS. This discrepancy arises because we do not construct cost volumes iteratively based on stereo matching, enabling better parallelism.

B. Obstacle Tracking

We finetune ODTFormer for voxel tracking using the synthetic SceneFlow Driving dataset [42] and the real-world KITTI 2015 dataset [30], which contains over 4,000 stereo image sequences with the resolution of 540×960 and 200 stereo images with the resolution of 370×1224 , respectively. We resize the images and camera parameters to have the same resolution 400×880 as in DrivingStereo, then apply the same data augmentation as obstacle detection and voxelize the scene flow by averaging the scene flow vectors of the point cloud within occupied voxels. We set the FPS and velocity thresholds for the matching boundary as $f = 26$ FPS and $v = 33.3m/s$ (or $120 km/h$). The volumetric region of each voxel during tracking is set to be the surrounding $9 \times 3 \times 9$ voxels along x, y, z axes.

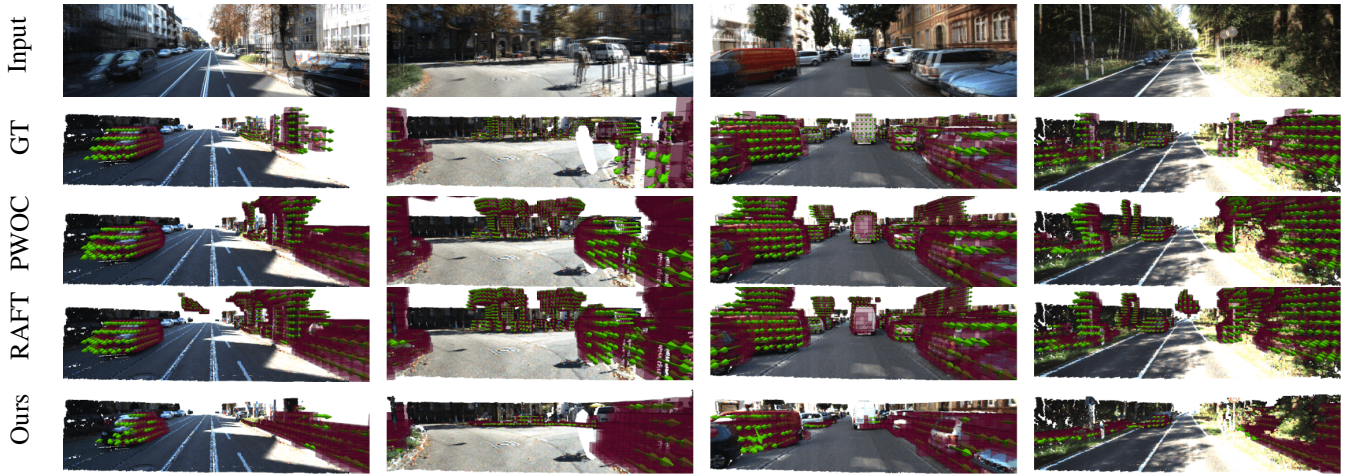


Fig. 4: **Visual results of obstacle tracking on the KITTI 2015 dataset.** The first row shows the stacked images from consecutive frames. The obstacle detection results are shown as red cubes, and the tracking results are marked as green arrows. Longer arrows indicate large motion magnitude.

We randomly split the KITTI training dataset into a custom training set and a validation set using an 8:2 ratio. We first train the network for 20 epochs with a batch size of 8 on the SceneFlow Driving dataset, then finetune on the KITTI dataset for 50 epochs.

Evaluation metrics. We employ EPE and foreground EPE to assess our model’s performance in voxel flow estimation. EPE quantifies the Euclidean distance between the ground-truth flow vector and the estimated flow vector. Foreground EPE is the EPE loss measured upon those voxels classified as occupied within the ground truth so that the measurement considers both occupancy detection and flow estimation accuracies.

Baseline methods. We compare against stereo-based scene flow estimation methods: SENSE [41], PWOC-3D [34], and RAFT-3D [22] with MobileStereoNet (MSNet) [24] for depth estimation. We use their model weights trained on the *entire* KITTI 2015 training set, including our custom validation set, to evaluate all the methods. We first estimate the scene flow given stereo image pairs from consecutive frames, then voxelize it using the same process as we generate the ground truth (Sec. IV-B).

Quantitative results. The quantitative results are evaluated using our custom validation set within the KITTI 2015 training data since our model cannot be directly evaluated through the KITTI test submission. The results are summarized in Table II. To filter outliers in dense scene flow estimation methods, we clipped the results from compared methods into the same bounded tracking boundary of [9, 3, 9]. Although the validation set we use is included in the baseline methods’ training data, with significantly fewer MACs than all dense scene flow methods, ODTFormer still outperforms baselines on both metrics except RAFT-3D. It validates the effectiveness of our proposed approach to track obstacles. Notably, we achieve comparable performance to RAFT-3D using only ten-fold or twenty-fold less MAC

TABLE III: Ablation study for obstacle detection.

Method	Level	CD (meters) ↓		IoU (%) ↑		MACs ↓	Params ↓
		15.0m	30.0m	15.0m	30.0m		
Ours w/o geometric constraints	1	6.17	9.92	77.77	64.01	23.79G	7.52M
	2	4.69	9.34	66.94	47.36		
	3	4.56	12.97	48.58	29.40		
	4	4.98	23.65	34.07	17.10		
Ours w/o 3D volumetric sampling	1	3.73	5.49	82.53	72.92	24.65G	6.22M
	2	2.51	4.53	72.66	57.72		
	3	2.23	5.74	54.18	38.45		
	4	2.42	12.03	38.30	22.57		
Ours	1	3.09	4.13	85.25	76.75	25.05G	6.14M
	2	1.92	3.22	76.10	62.69		
	3	1.64	3.82	57.94	43.79		
	4	1.87	8.72	41.62	26.62		

operations. In terms of inference speed, the tracing module adds negligible burden to the detection module thanks to our design. The entire model runs at 20fps for the KITTI resolution (370×1224) on an RTX A5000 GPU without careful postprocessing such as quantization, pruning, etc.

C. Ablation Studies

Obstacle prediction. To study the effectiveness of our design choices, we conduct ablation studies and report the results in Table III using the DrivingStereo dataset.

We first replace the feature concatenations in Eq.(3) and (4) with cross-view averaging as in existing work [6–9]. In this way, it discards the geometric constraints between two stereo images. As we can see at the top of Table III, it leads to significantly worse detection accuracy.

We also study the effectiveness of the matching cost refinement and simply use C_B as the final cost volume. As shown in the middle of Table III, it shows worse results than our full model with fewer model parameters and a slight decrease in computation burden. Given the fast inference speed of our model, we include it for the best obstacle detection accuracy.

Voxel tracking. The ablation study for our voxel tracker design is reported in Table IV using KITTI 2015. We see that without the bounded tracking design, the tracker directly results in an out-of-memory error when evaluating with an

TABLE IV: Ablation study for obstacle tracking. OOM indicates an out-of-memory error.

Bounded	Match Occupied Voxels Only in the Second Frame	EPE↓ (meters)	Foreground EPE↓ meters
✗	✗	OOM	OOM
✓	✓	0.021	1.092
✓	✗	0.021	1.087

RTX A5000 GPU. Whether to match occupied voxels in the second frame is an open design choice. It leads to a minor performance degradation in the foreground area, partially because the occupancy detection results may not always be correct. We therefore opt to track all voxels for the best tracking accuracy.

V. CONCLUSION

This paper contributed ODTFormer, a Transformer-based model to address obstacle detection and tracking problems. We detect the obstacles in the form of a voxel occupancy grid using a deformable attention mechanism and track them by matching voxels across two consecutive frames. We confirmed the effectiveness of ODTFormer by providing extensive quantitative and qualitative results on DrivingStereo and KITTI datasets. Our results showed that ODTFormer achieves state-of-the-art performance on obstacle detection tasks. For the obstacle tracking task, we showed that ODTFormer can achieve accuracy comparable to state-of-the-art methods. The entire model is efficient and runs fast at 20fps for the KITTI resolution (370×1224) on an RTX A5000 GPU without careful postprocessing (e.g., quantization, pruning). Being able to detect and track obstacles at the same time will empower a set of downstream tasks, for instance, ensuring safe navigation when people move around a robot.

REFERENCES

- [1] J. Frey *et al.*, “Locomotion Policy Guided Traversability Learning using Volumetric Representations of Complex Environments,” Aug. 2022, arXiv:2203.15854 [cs].
- [2] Y. F. Chen *et al.*, “Socially aware motion planning with deep reinforcement learning,” in *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, Sep. 2017, pp. 1343–1350.
- [3] N. U. Akmandor *et al.*, “Deep Reinforcement Learning based Robot Navigation in Dynamic Environments using Occupancy Values of Motion Primitives,” in *IROS*, 2022.
- [4] L. Zhao *et al.*, “E(2)-equivariant graph planning for navigation,” *IEEE Robotics and Automation Letters*, vol. 9, no. 4, 2024.
- [5] H. Li *et al.*, “Stereonavnet: Learning to navigate using stereo cameras with auxiliary occupancy voxels,” Mar. 2024, arXiv:2403.12039 [cs].
- [6] Z. Li *et al.*, “BEVFormer: Learning Bird’s-Eye-View Representation from Multi-camera Images via Spatiotemporal Transformers,” in *Computer Vision – ECCV 2022*, S. Avidan *et al.*, Eds., 2022, pp. 1–18.
- [7] Y. Li *et al.*, “VoxFormer: Sparse Voxel Transformer for Camera-Based 3D Semantic Scene Completion,” 2023, pp. 9087–9098.
- [8] Y. Wei *et al.*, “SurroundOcc: Multi-camera 3D Occupancy Prediction for Autonomous Driving,” 2023, pp. 21 729–21 740.
- [9] Y. Huang *et al.*, “Tri-Perspective View for Vision-Based 3D Semantic Occupancy Prediction,” 2023, pp. 9223–9232.
- [10] T. Eppenberger *et al.*, “Leveraging Stereo-Camera Data for Real-Time Dynamic Obstacle Detection and Tracking,” in *IROS*, 2020.
- [11] Z. Xu *et al.*, “A real-time dynamic obstacle tracking and mapping system for UAV navigation and collision avoidance with an RGB-D camera,” in *ICRA*, 2023.
- [12] M. Lu *et al.*, “Perception and Avoidance of Multiple Small Fast Moving Objects for Quadrotors With Only Low-Cost RGBD Camera,” *IEEE RA-L*, vol. 7, no. 4, Oct. 2022.

- [13] Y. Wang *et al.*, “Autonomous Flights in Dynamic Environments with Onboard Vision,” in *IROS*, 2021.
- [14] H. Li *et al.*, “Stereovoxelnet: Real-time obstacle detection based on occupancy voxels from a stereo camera using deep neural networks,” in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, 2023.
- [15] S. Kareer *et al.*, “ViNL: Visual Navigation and Locomotion Over Obstacles,” Jan. 2023, arXiv:2210.14791 [cs].
- [16] B. Zhou *et al.*, “RAPTOR: Robust and Perception-aware Trajectory Replanning for Quadrotor Fast Flight,” Jul. 2020, arXiv:2007.03465.
- [17] A. Loquercio *et al.*, “Learning high-speed flight in the wild,” *Science Robotics*, vol. 6, no. 59, p. eabg5810, 2021.
- [18] D. Falanga *et al.*, “How Fast Is Too Fast? The Role of Perception Latency in High-Speed Sense and Avoid,” *IEEE Robotics and Automation Letters*, vol. 4, no. 2, pp. 1884–1891, Apr. 2019.
- [19] Y. Wang *et al.*, “Pseudo-LiDAR From Visual Depth Estimation: Bridging the Gap in 3D Object Detection for Autonomous Driving,” 2019, pp. 8445–8453.
- [20] S. Vedula *et al.*, “Three-dimensional scene flow,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 27, no. 3, pp. 475–480, Mar. 2005.
- [21] H. Jiang *et al.*, “SENSE: a Shared Encoder Network for Scene-flow Estimation,” Oct. 2019, arXiv:1910.12361 [cs].
- [22] Z. Teed and J. Deng, “RAFT-3D: Scene Flow using Rigid-Motion Embeddings,” Apr. 2021, arXiv:2012.00726 [cs].
- [23] A. Vaswani *et al.*, “Attention is All you Need,” in *Advances in Neural Information Processing Systems*, vol. 30, 2017.
- [24] F. Shamsafar *et al.*, “MobileStereoNet: Towards Lightweight Deep Networks for Stereo Matching,” 2022, pp. 2417–2426.
- [25] Z. Shen *et al.*, “CFNet: Cascade and Fused Cost Volume for Robust Stereo Matching,” 2021, pp. 13 906–13 915.
- [26] B. Liu *et al.*, “Local Similarity Pattern and Cost Self-Reassembling for Deep Stereo Matching Networks,” *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 2, pp. 1647–1655, Jun. 2022, number: 2.
- [27] G. Xu *et al.*, “Attention Concatenation Volume for Accurate and Efficient Stereo Matching,” 2022, pp. 12 981–12 990.
- [28] X. Zhu *et al.*, “Deformable detr: Deformable transformers for end-to-end object detection,” *arXiv preprint arXiv:2010.04159*, 2020.
- [29] G. Yang *et al.*, “Drivingstereo: A large-scale dataset for stereo matching in autonomous driving scenarios,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [30] M. Menze and A. Geiger, “Object scene flow for autonomous vehicles,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 3061–3070.
- [31] M. Tan and Q. Le, “EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks,” in *Proceedings of the 36th International Conference on Machine Learning*. PMLR, May 2019, pp. 6105–6114.
- [32] T.-Y. Lin *et al.*, “Feature pyramid networks for object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.
- [33] H. Hirschmuller, “Accurate and efficient stereo processing by semi-global matching and mutual information,” in *CVPR*, vol. 2, 2005.
- [34] R. Saxena *et al.*, “Pwoc-3d: Deep occlusion-aware end-to-end scene flow estimation,” in *2019 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2019, pp. 324–331.
- [35] I. Shepel *et al.*, “Occupancy grid generation with dynamic obstacle segmentation in stereo images,” *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 14 779–14 789, 2021.
- [36] H. Wu *et al.*, “Casa: A cascade attention network for 3-d object detection from lidar point clouds,” *IEEE Transactions on Geoscience and Remote Sensing*, vol. 60, pp. 1–11, 2022.
- [37] —, “Virtual sparse convolution for multimodal 3d object detection,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 21 653–21 662.
- [38] X. Wang *et al.*, “You only need two detectors to achieve multi-modal 3d multi-object tracking,” *arXiv preprint arXiv:2304.08709*, 2023.
- [39] Y. Xie *et al.*, “Pixel-Aligned Recurrent Queries for Multi-View 3D Object Detection,” in *ICCV*, 2023.
- [40] I. Loshchilov and F. Hutter, “Decoupled Weight Decay Regularization,” Sep. 2018.
- [41] H. Jiang *et al.*, “Sense: A shared encoder network for scene-flow estimation,” in *ICCV*, 2019.
- [42] N. Mayer *et al.*, “A large dataset to train convolutional networks for disparity, optical flow, and scene flow estimation,” in *CVPR*, 2016.