# An End-to-End Framework of Road User Detection, Tracking, and Prediction from Monocular Images

Hao Cheng[1], Mengmeng Liu[2], and Lin Chen[3]

*Abstract*—Perception that involves multi-object detection and tracking, and trajectory prediction are two major tasks of autonomous driving. However, they are currently mostly studied separately, which results in most trajectory prediction modules being developed based on ground truth trajectories without taking into account that trajectories extracted from the detection and tracking modules in real-world scenarios are noisy. These noisy trajectories can have a significant impact on the performance of the trajectory predictor and can lead to serious prediction errors. In this paper, we build an end-to-end framework for detection, tracking, and trajectory prediction called ODTP (Online Detection, Tracking and Prediction). It adopts the state-of-the-art online multi-object tracking model, QD-3DT, for perception and trains the trajectory predictor, DCENet++, for perception and trains the trajectory predictor, DCENet++, directly based on the detection results without purely relying on ground truth trajectories. We evaluate the performance of ODTP on the widely used nuScenes dataset for autonomous driving. Extensive experiments show that ODPT achieves high performance end-to-end trajectory prediction. DCENet++, with the enhanced dynamic maps, predicts more accurate trajectories than its base model. It is also more robust when compared with other generative and deterministic trajectory prediction models trained on noisy detection results.

## I. INTRODUCTION

Trajectory prediction plays a crucial role in achieving autonomous driving. It involves observing the behavior of agents like vehicles, pedestrians, and other road users in a few past time steps. This observation information includes perceiving the road users' type and their past trajectories, which is then fed into a trajectory predictor to forecast their potential trajectories in the following time steps. Despite the rapid development of trajectory prediction methods, they are usually developed independently of the perception module, assuming that the ground truth information of the road users' past trajectories is already available. This means that trajectory predictors trained on ground truth data may not reflect the complexities of real-world driving scenarios [1]. Moreover, the input data to the trajectory prediction module is prone to noise because the perception module is imperfect due to long-standing issues such as changes in lighting, scale, background interference, sensor limitations, and multiple occlusions. Therefore, this paper aims to address the task of trajectory prediction by integrating the perception module and developing an end-to-end framework for road user detection, tracking, and prediction.

Object detection results lay as the foundation for multi-object tracking and trajectory prediction. Our focus is primarily on mixed traffic scenarios that comprise not only vehicles and pedestrians but also other types of road users. In this paper, we have selected monocular images obtained from a moving vehicle as the input for the perception module, as camera sensors are low-cost for capturing all the objects in the field of view and straightforward to deploy. The perception module we have employed is the monocular Quasi-Dense 3D Object Tracking (QD-3DT) [2] as the base model, which can effectively associate moving agents over time and estimate their complete 3D bounding box information from a sequence of 2D images captured on a mobile platform.

To achieve end-to-end prediction, the prediction module takes the 2D positional information at discrete time steps from the perception module as input. In order to consider the multimodality nature of agents' movements and their mutual influence during interactions, we use the DCENet model [3] as the base model for the prediction module. DCENet explores the spatial and temporal information captured by dynamic maps and leverages the attention mechanisms [4] and Conditional Variational Autoencoder [5] to predict agents' multimodal trajectories. Instead of relying on ground truth point-wise trajectory data, we extend DCENet to take as input the perceived trajectory data, including the agents' shape and pose information. This approach allows the dynamic maps to more accurately map the position, velocity, and pose information of each agent into the 2D grid cells that are projected from the agents' detected 3D shape. Furthermore, the extended DCENet model is trained based on detection results, which are more robust against detection noise compared to ground truth data. We term this new prediction model DCENet++ to reflect its improved performance.

Figure 1 depicts the end-to-end detection, tracking, and prediction framework overview. It comprises three primary components: monocular image sensor input captured from the ego vehicle, the QD-3DT-based 3D detection and tracking module, and the DCENet++ multimodal prediction module for trajectory prediction. We name our end-to-end framework ODTP (Online Detection, Tracking and Prediction). The **key contributions** of our work are as follows:

- We propose an end-to-end framework ODTP, which automatically detects and tracks various types of road users from monocular images and predicts their multimodal trajectories. The prediction module in ODTP is

[1]Scene Understanding Group, University of Twente, The Netherlands, `h.cheng-2@utwente.nl`, Cheng is funded by MSCA European Postdoctoral Fellowships under the 101062870 - VeVuSafety project.

[2]Institute of Cartography and Geoinformatics, Leibniz University Hannover, Germany, `mengmeng.liu1998@gmail.com`.

[3]VISCODA GmbH, Schneiderberg 32, 30167 Hannover, Germany, `chen@viscoda.com`.

Fig. 1: The overview end-to-end detection, tracking and trajectory prediction framework using camera sensors.

trained directly on the detection and tracking results, rather than on the manually labeled ground truth data.

- We extend the trajectory prediction model DCENet by adding road users' shape and pose information acquired by the detection module, called DCENet++, which achieves better performance than DCENet that only relies on the 2D positional information.

## II. RELATED WORK

### A. Object Detection and Tracking

In this paper, we focus on object detection and tracking based on images. In the paradigm of object detection, the most attention has been given to two-stage detectors represented by R-CNN-based models with region proposal networks [6]–[8] and one-stage detectors represented by YOLO-based models [9, 10]. SORT [11] uses Faster R-CNN [8] for object detection and predicts and updates the motion track using the Kalman filter [12]. The Hungarian dichotomous matching algorithm [13] with Intersection Over Union (IOU) as the matching criterion is then used to match detected and tracked objects. DeepSORT [14] further adds appearance representation of detections using deep neural networks to enhance the tracking performance. Several works, such as [15]–[19], propose the use of 3D information to narrow the search area and make the object's trajectory smoother. A dominant approach to associating data in multi-object tracking (MOT) problems is to utilize different kinds of costs, including trajectory priors, bounding box center locations, optical flow, bounding box overlap, and appearance information or deep appearance features [2, 16, 20]. For example, [16] employs simple complementary costs for data association, which include 2D-3D cost, 3D-3D cost, appearance cost, and shape and pose cost. [20] employs discriminative feature embeddings and the greedy bipartite matching method to match new detections and trajectories. Moreover, DEFT [21] proposes a joint detection and tracking method that relies on the appearance features from a detection backbone for object-to-track association. 3D-Times [22] employs attention mechanisms to learn spatial-temporal information cues for joint 3D detection and tracking from monocular videos. In this paper, we use QD-3DT [2] as our base model for joint 3D detection and tracking. It employs Faster R-CNN for 2D detection basis and Region of Interest (ROI) features for each proposal extracted by a region proposal network to regress the 3D dimensions, 2D

projection center, depth, and orientation. Moreover, different similarity cues, such as the deep representation similarity, the overlap of 3D bounding boxes, and the motion similarity, are utilized for data association. In addition, QD-3DT uses an LSTM-based module instead of using 3D KF for motion refinement.

### B. Trajectory Prediction

Recent years, deep learning methods, such as Generative Adversarial Network (GAN) [23], Conditional Variational AutoEncoder (CVAE) [5, 24], and attention mechanisms [25, 26], have been introduced to the trajectory prediction task. Gupta et al. [27] propose a fusion of LSTM and GAN, using the global pooling module of LSTM as the encoder-decoder generator and a discriminator composed of multiple LSTMs. For machine navigation, Altan et al. [28] propose a pedestrian-dependent spatio-temporal graphical representation that can effectively represent the importance of pedestrians in congested environments. Social-BiGAT [29] proposes the Graph Attention Network (GAT) [30] to learn feature representations and performs reversible transformations between the scene and the response's underlying noise vector. Similarly, SR-LSTM [31] employs Graph Neural Networks (GNN) to model the interconnections among agents and predicts their deterministic future trajectories. Agent-Former [32] utilizes the Transformer network [4] to learn spatial-temporal information of agents and applies CVAE for multimodal trajectory predictions. In contrast, [3] proposes a model called DCENet that utilizes self-attention and LSTM to model interactions between agents and a CVAE framework to predict a set of possible trajectories conditioned on its observed trajectory and the learned dynamic context for each agent. Considering its superior performance, we utilize DCENet as our prediction model baseline. In comparison to the original DCENet, we add the agent's shape and pose information to the dynamic maps, which can model interactions more accurately. We also adopt DCENet to moving camera scenes, i.e., real autonomous driving scenarios. In this way, in contrast to training the prediction model using ground truth trajectories, our model utilizes camera data only, and the prediction is conducted on detection and tracking results.

### C. Joint Tracking and Forecasting

Recently, a few studies have investigated the possibility of joint MOT and trajectory prediction [1, 33, 34]. Weng et al. [34] propose a novel data association method that utilizes

GNNs to model interactions between new detections and trajectories. After message passing by GNNs, the affinity matrix between new detections and trajectories is learned for association. Instead of processing the MOT task first and then the trajectory prediction task, Weng et al. [34] process the two tasks in parallel, so the trajectory prediction task does not explicitly depend on the MOT results. Liang et al. [33] change the order of tasks for MOT and trajectory prediction. Unlike the traditional approach, it first carries out joint detection and prediction tasks before updating the object trajectory. Zhang et al. [1] propose to extract the motion information based on the affinity cues among detection results and predict trajectories directly based on detection results. Inspired by the query-based end-to-end object detection with transformers [35], end-to-end perception and motion prediction is achieved by extending the object query with recurrent temporal information, such as ViP3D [36] and UniAD [37]. In contrast to those works, our work utilizes camera data only to build an end-to-end detection, tracking, and prediction framework. The trajectories extracted using monocular images have noise and introduce a greater challenge to the task of trajectory prediction. As most trajectory prediction methods are trained and inferenced on ground truth trajectories, the effect of noise in the detection and tracking tasks on trajectory prediction is not considered. To the opposite, we compare the performance differences between generative and deterministic trajectory prediction models when dealing with such noisy trajectories.

## III. METHODOLOGY

### A. Problem Formulation

The goal of the end-to-end framework ODPT is to take a monocular image sequence as input for the 3D object detection module and output a set of 3D bounding boxes (bbx) at frame $t$, denoted as $S = \{s_1^t, ..., s_J^t\}$. After performing data association and motion refinement in the MOT module, which takes the bbx as input, we obtain a series of smooth trajectories denoted as $\mathbb{T} = \{\tau_1, ..., \tau_N\}$, where $\tau_i \in \mathbb{R}^{T \times 2}$, and the refined 3D bounding boxes denoted as $\mathbb{S} = \{s_1^t, ..., s_N^t\}$. Here, $i \in \{1, ..., N\}$, $N \leq J$ represents the total number of detected and tracked agents in the given scenario, and $T$ is the observed time horizon. $T \geq 2$, to make sure that we have enough observed steps to derive the speed and pose information. Subsequently, the trajectories $\mathbb{T}$ and detection bbx $\mathbb{S}$ serve as the input to the trajectory prediction module. In this module, we predict a set of possible future trajectories denoted as $\{\hat{Y}_{i,1}^{T+1:T'}, ..., \hat{Y}_{i,K}^{T+1:T'}\}$ conditioned on the detected trajectories and bbx for each agent $i$. Here, $K$ represents the number of predicted trajectories and $T' - T$ represents the predicted time horizon. In the following sections, we explain each module of ODPT in detail.

### B. QD-3DT

The goal of the QD-3DT module [2] is to provide the 3D information of all tracked objects by inputting consecutive frames of monocular images and GPS/IMU information from the ego vehicle. The GPS/IMU data is used to obtain



Fig. 2: Comparison between the original and the refined dynamic maps with the agents' shape and pose information. (a) No agents' shape and pose information (b). Only considering the agents' shape information, and (c) Considering both the agents' shape and pose information.

localization information about the ego vehicle's motion. To achieve this, we transfer the 3D information, including the shape, pose, and position of all neighboring agents, from the cameras to the local frame of the ego agent. During the perception process, the monocular images are first processed through a backbone network, such as VGG16 [38], and a Region Proposal Network (RPN) [8], to generate 2D regions of interest (ROIs). These ROIs are then fed into two multi-head networks, which output similarity feature embeddings and 3D layouts. To track 3D object instances over time, multimodal similarity metrics between the tracked trajectories and detected objects are computed utilizing 3D information, motion information, and feature embeddings. Additionally, motion-aware data association and depth-ordering matching techniques are used to mitigate occlusion problems. Finally, the tracking module refines the 3D information of the objects. It is worth noting that in our approach, we directly use the image pixel coordinates without normalization to calculate the IOU between the detected and ground truth bbx. This differs from the original setting of QD-3DT. We made this choice based on empirical findings that using normalized coordinates changes the image scale and leads to larger detection errors [2].

### C. DCENet++

We utilize the DCENet model [3] for the trajectory prediction module and adapt it from a bird's-eye view to the ego perspective of the mobile cameras. With the 3D tracking module from QD-3DT, we not only estimate the 3D object center $\{x, y, z\}$ but also the object dimensions $D = \{l, h, w\}$ and object pose $\theta$ for each agent. Therefore, compared to the original dynamic maps that use an approximation of the agent's shape, we use the detected shape and pose information to refine the dynamic maps, allowing for a more accurate modeling of the neighboring agents' position, velocity, and pose information. In this work, we assume that all agents are moving on the ground surface and our prediction task is focused on the 2D positions of the $x$- and $y$-coordinates. Consequently, the dynamic maps for all the objects are modelled based on the projection on the ground plane, leaving object height $h$ and altitude $z$ out of consideration. Fig. 2 (a) shows the original dynamic map [3], which assumes that each agent has the same size and orientation. This approximation is based on the observation that pedestrian size varies little within the pedestrian trajectory dataset (e.g., [39]). Since pedestrians are relatively

small and occupy only one grid cell in the dynamic maps, orientation information can be disregarded in the pedestrian dataset. However, when adapting DCENet to autonomous driving scenarios, we need to consider the significant shape differences among heterogeneous types of agents, such as vehicles and cyclists. To ensure proper alignment, especially for large agents occupying multiple grid cells in the dynamic maps, the occupied grid cells should align with the agent's pose. Fig. 2 (b) and (c) illustrate the dynamic maps that include shape information alone and both shape and pose information, respectively. These refined dynamic maps enable us to handle objects with varying shapes and poses. We refer to DCENet with the refined dynamic maps as DCENet++.

### D. Joint 3D Tracking and Forecasting

With the QD-3DT perception module, which includes data association and tracking, we obtain a set of tracked trajectories $\mathbb{T} = \{\tau_1^{1:T}, ..., \tau_N^{1:T}\}$ from previous frames, as well as a set of detections $\mathbb{S} = \{s_1^T, ..., s_N^T\}$ at frame $T$. We configure the batch size of DCENet++ to match the number of objects detected in the current frame, the same as $N$. This choice allows us to focus solely on the interactions among agents that are present concurrently and successfully detected in the given frames. To model the interactions between the ego agent and its neighboring agents, we employ the extended dynamic maps. At each time step $t \in \{1, ..., T\}$, we project the neighboring agents onto the grid cells of the dynamic map, centered at the current position $\{x, y\}$ of the ego agent. The projection is based on the detected 3D shape $\{l, h, w\}$ obtained from $\mathbb{S}$. Next, we map the position, velocity, and pose information derived from the tracked trajectories $\mathbb{T}$ onto dedicated channels in the dynamic maps, following the approach described in [3]. Simultaneously, the offset sequence $\Delta X_i^{1:T-1} = \{\Delta x_i^1, ..., \Delta x_i^{T-1}\} \in \mathbb{R}^{(T-1) \times 2}$ for each agent's trajectory is combined with the sequence of dynamic maps at the corresponding time steps. These combined inputs serve as the joint condition for the prediction module of DCENet++. Finally, DCENet++ predicts multimodal trajectories $\{\hat{Y}_{i,1}^{T+1:T'}, ..., \hat{Y}_{i,K}^{T+1:T'}\}$ for all agents.

## IV. EXPERIMENT

### A. Dataset

We evaluate our ODTP on nuScenes [40], which is one of the most commonly used large-scale real-word datasets for autonomous driving. The ego vehicle is equipped with multiple sensors, such as LiDAR, monocular camera, and radar. In this paper, we focus on the camera images. In total, the dataset contains 1000 driving scenes in Boston and Singapore, including 700 scenes for training, 150 scenes for validation and 150 scenes for test.

### B. Evaluation Metrics

**MOT metrics.** We adhere to the tracking metrics established by nuScenes, specifically AMOTA and AMOTP [41]. AMOTA stands for averaged multi-object tracking accuracy (MOTA) [42] at various recall thresholds. MOTA provides a comprehensive assessment by taking into account false positives, missed targets, and identity switches. AMOTP refers to the average multi-object tracking precision (MOTP) [42]. MOTP quantifies the misalignment between the annotated and predicted bounding boxes, offering insights into the accuracy of object localization.

**Trajectory prediction metrics.** We employ two widely used metrics to evaluate the trajectory prediction task: Average Displacement Error (ADE) and Final Displacement Error (FDE) [43]. ADE calculates the average Euclidean distance between the predicted trajectory and the corresponding ground truth trajectory, while FDE measures the Euclidean distance between their final positions. Consistent with previous works [27, 44, 45], we select the minimum ADE ($ADE_K$) and FDE ($FDE_K$) from the best prediction among $K$ trajectory samples for each agent.

### C. Experimental Setting

Following the VeloLSTM module in QD-3DT [2], we set the observation time horizon to 2.5 seconds with a frame rate of 2 Hz. The prediction time horizon is set to 4 seconds with the same frame rate.

For the 3D detection and tracking module, we adopt the approach presented in QD-3DT [2]. We utilize a pretrained Faster R-CNN [8] model on ImageNet [46] from TorchVision [47] for 2D detection and 3D center estimation. However, to address accumulated tracking errors, we modify the default setting in QD-3DT. Instead of continuously predicting the object state until it goes beyond the tracking range or its lifespan ends (e.g., ten time steps), we use the predicted bbx at the next step once to compute the affinity between the trajectory and the detected object state. Moreover, different from the default setting in QD-3DT, we use the image pixel coordinates without normalization to calculate the IOU between the detected and ground truth bbx.

During training, the VeloLSTM module is trained for 100 epochs using ten sample frames per object trajectory, with a batch size of 128. For trajectory prediction model DCENet++, we employ the Adam optimizer [48] with early stopping, setting the patience parameter to ten to prevent overfitting. The initial learning rate is set to $10^{-4}$, and we decay the learning rate by a factor of 0.5 every 20 epochs.

## V. RESULTS

In this section, we begin by evaluating the performance of the perception module. Subsequently, we assess the performance of the trajectory prediction module. Finally, we present the qualitative performance of the ODPT framework for end-to-end object detection, tracking, and prediction.

### A. Perception Performance

Firstly, we compare our setting with the default setting in terms of computing the IOU between the detected bounding box and the ground truth annotation, as well as the affinity between the trajectory and the detected object state. Table I illustrates the results of this comparison. The improved AMOTA and the decreased AMOTP demonstrate that our

pixel-based IOU computation and the one-time update of the detected object state yield better performance compared to the default setting in QD-3DT.

TABLE I: Perception performance for multi-object tracking.

| IOU | state update | AMOTA↑ | AMOTP↓ |
|---|---|---|---|
| default | default | 0.233 | 1.528 |
| pixel | default | 0.235 | 1.517 |
| pixel | one time | **0.243** | **1.512** |

Next, we explore different IOU thresholds to strike a balance between precision and recall for object detection. As depicted in Table II, we observe that an IOU threshold of 0.5 yields slightly better results compared to the other thresholds. Consequently, we adopt an IOU threshold of 0.5 as the default setting for subsequent experiments.

TABLE II: Comparison of IOU.

| IOU | AMOTA↑ | AMOTP↓ |
|---|---|---|
| 0.4 | 0.236 | 1.516 |
| 0.5 | **0.243** | **1.512** |
| 0.7 | **0.243** | 1.515 |

### B. Trajectory Prediction Performance

We conduct experiments by varying the dimensions of the latent variable in the CVAE-based DCENet++ model. As presented in Table III, increasing the dimension from 2 to 32 leads to a reduction in trajectory prediction errors, as measured by $ADE_{10}$ and $FDE_{10}$. However, when we further increase the dimension to 64, the prediction performance deteriorates. Therefore, for subsequent experiments, we maintain a fixed dimension of 32 for the latent variable.

TABLE III: Different dimensions of latent variable $z$.

| Methods | $ADE_{10}\downarrow$ | $FDE_{10}\downarrow$ |
|---|---|---|
| $z_{dim}=2$ | 0.82 | 1.54 |
| $z_{dim}=32$ | **0.79** | **1.50** |
| $z_{dim}=64$ | 0.82 | 1.53 |

Additionally, we perform an ablation study on the dimension and pose information in the dynamic maps of DCENet++. Comparing it to the baseline model, DCENet, which lacks dimension and pose information for aligning agents in the dynamic maps, the performance of DCENet++ is better in terms of both ADE and FDE, as shown in Table IV. Interestingly, we discover that utilizing either dimension or pose information alone does not result in a clear improvement in performance. This is because either the lack of dimension or pose information could lead to sub-optimal alignments in the dynamic maps.

After adjusting the hyperparameters of DCENet++, we conducted a performance comparison with two well-known trajectory prediction models using the nuScenes motion prediction dataset, as shown in Table V. It should be noted that, for a fair comparison, all these models do not use any

TABLE IV: The ablation study for the dynamic maps.

| Methods | dimension | pose | $ADE_{10}\downarrow$ | $FDE_{10}\downarrow$ |
|---|---|---|---|---|
| DCENet | - | - | 0.80 | 1.51 |
| DCENet+ | √ | - | 0.80 | 1.51 |
| DCENet+ | - | √ | 0.80 | **1.50** |
| DCENet++ | √ | √ | **0.79** | **1.50** |

map with scene context information. One of the models we compared DCENet++ to is AgentFormer [32], which is a transformer and CVAE-based model capable of generating multimodal predictions for each agent. In the multimodal prediction task, DCENet++ outperforms AgentFormer in predicting both five and ten modalities. Furthermore, DCENet++ exhibits superior performance compared to AgentFormer in single-modal prediction as well. Given these promising results, we proceeded to compare DCENet++ with the deterministic model SR-LSTM [31], which utilizes GNN to model interactions among agents. In terms of both ADE and FDE, DCENet++ surpasses SR-LSTM.

TABLE V: Evaluation of trajectory prediction on nuScenes.

| Methods | $ADE_1\downarrow$ | $FDE_1\downarrow$ | $ADE_5\downarrow$ | $FDE_5\downarrow$ | $ADE_{10}\downarrow$ | $FDE_{10}\downarrow$ |
|---|---|---|---|---|---|---|
| AgentFormer [32] | 7.91 | 4.55 | 1.67 | 2.62 | 1.06 | 1.56 |
| SR-LSTM [31] | 1.29 | 2.57 | - | - | - | - |
| DCENet++ | **0.97** | **1.86** | **0.86** | **1.62** | **0.79** | **1.50** |

In the following analysis, we examine the impact of noisy MOT training data by comparing the performance of DCENet++ when trained on ground truth (GT) trajectories versus MOT trajectories. The results presented in Table VI clearly demonstrate that, as anticipated, training the prediction model DCENet++ using GT trajectories and subsequently testing it on MOT trajectories leads to a significant drop in performance. Namely, the prediction errors increases by over 200% when compared to the realistic setting of both training and testing DCENet++ using the MOT trajectory data. This indicates that a model trained solely on GT trajectories may struggle to generalize effectively in real-world driving scenarios, where the perception module unavoidably produces noisy trajectories during the observation period. However, this issue can be largely mitigated when DCENet++ is trained and tested on MOT trajectories. Remarkably, the MOT-trained DCENet++ exhibits impressive generalization capabilities when deployed in testing on GT trajectories, performing only slightly worse than the ideal scenario where both training and testing are conducted on GT trajectories.

TABLE VI: Comparison between using ground truth and MOT trajectory data. Using the MOT data for both training and testing is referred as the baseline of prediction errors.

| Model | Training | | Testing | | Errors | |
|---|---|---|---|---|---|---|
| | MOT | GT | MOT | GT | $ADE_{10}\downarrow$ | $FDE_{10}\downarrow$ |
| DCENet++ | | √ | √ | | 2.54 (+222%) | 4.59 (206%) |
| DCENet++ | √ | | √ | | 0.79 | 1.50 |
| DCENet++ | √ | | | √ | 0.47 (-40%) | 0.94 (-37%) |
| DCENet++ | | √ | | √ | 0.44 (-44%) | 0.90 (-40%) |

Fig. 3: Qualitative results of DCENet++ (first row) and SR-LSTM [31] (second row). On the left is the ego-driving perspective and on the right is the bird's-eye view. The squares represent predicted trajectories and circles denote history trajectories.



Fig. 4: Qualitative results of DCENet++ (first row) and AgentFormer [32] (second row). On the left is the ego-driving perspective and on the right is the bird's-eye view. The squares represent predicted trajectories and circles denote history trajectories.



Fig. 5: Qualitative results of DCENet++ from multi-camera perspectives. On the left is the ego-driving perspective and on the right is the bird's-eye view. The squares represent predicted trajectories and circles denote history trajectories.

## C. Qualitative Results

We present the performance analysis of the ODPT in various driving scenarios. To visualize the results, we employ the visualization method proposed by Hu et al. [2]. Specifically, we display the detection results obtained from the camera view and plot the tracked and predicted trajectories using a bird's-eye view perspective above the ego vehicle. Fig. 3 and 4 showcase the outputs of DCENet++ and SR-LSTM, and DCENet++ and AgentFormer models, respectively, revealing that they generate realistic predictions. However, compared to the other models, the predictions from DCENet++ exhibit smoother trajectories even though the input trajectory data from the MOT module is noisy.

Furthermore, Fig. 5 demonstrates the adaptability of ODPT to monocular images captured from various perspectives. This demonstrates the versatility and effectiveness of the ODPT approach across different camera viewpoints.

**Limitations.** In spite of the promising performance demonstrated above, in this work, we need to separately train the perception module and the trajectory prediction module, which is time-consuming. Also, because they are separately trained, we could not share the intermediate feature maps to unify the encodings for both the perception and trajectory prediction tasks. Moreover, whether the multimodal trajectory prediction in this work can consolidate the data association in MOT through providing a more time and space consistent object motion could be further explored. We leave this as our future work. Last but not least, when the perception model mis-detects agents due to, e.g., occlusions and lighting conditions, the trajectory prediction module fails to anticipate the movements of these mis-detected agents. To mitigate the issue of occlusions and detection limitations, one potential solution could be implementing cooperative perception by sharing detection information among agents using the vehicle-to-vehicle communication network [49, 50].

## VI. CONCLUSION

In this paper, we propose a framework called ODTP that combines the perception module of the monocular Quasi-Dense 3D Object Tracking with the trajectory module of DCENet. This framework enables end-to-end detection, tracking, and prediction for autonomous driving. We enhance the DCENet model by extending the dynamic maps to include agents' shape and pose information, which is termed DCENet++. This enhancement allows for more accurate mapping of interactions among agents. Furthermore, we demonstrate that training the trajectory prediction module using multi-object tracking data helps the prediction module better adapt to cope with the noisy data perceived in real-world driving scenarios.

## REFERENCES

[1] P. Zhang, L. Bai, Y. Wang, J. Fang, J. Xue, N. Zheng, and W. Ouyang, "Towards trajectory forecasting from detection," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, pp. 1–12, 2023.

[2] H.-N. Hu, Y.-H. Yang, T. Fischer, T. Darrell, F. Yu, and M. Sun, "Monocular quasi-dense 3d object tracking," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2022.

[3] H. Cheng, W. Liao, X. Tang, M. Y. Yang, M. Sester, and B. Rosenhahn, "Exploring dynamic context for multi-path trajectory prediction," in *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021, pp. 12 795–12 801.

[4] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.

[5] D. P. Kingma, S. Mohamed, D. J. Rezende, and M. Welling, "Semi-supervised learning with deep generative models," in *Advances in neural information processing systems*, 2014, pp. 3581–3589.

[6] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich feature hierarchies for accurate object detection and semantic segmentation," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2014, pp. 580–587.

[7] R. Girshick, "Fast r-cnn," in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.

[8] S. Ren, K. He, R. Girshick, and J. Sun, "Faster r-cnn: Towards real-time object detection with region proposal networks," *Advances in neural information processing systems*, vol. 28, 2015.

[9] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You only look once: Unified, real-time object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.

[10] J. Redmon and A. Farhadi, "Yolo9000: better, faster, stronger," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.

[11] A. Bewley, Z. Ge, L. Ott, F. Ramos, and B. Upcroft, "Simple online and realtime tracking," in *2016 IEEE international conference on image processing (ICIP)*. IEEE, 2016, pp. 3464–3468.

[12] L. D. Stone, R. L. Streit, T. L. Corwin, and K. L. Bell, *Bayesian multiple target tracking*. Artech House, 2013.

[13] H. W. Kuhn, "The hungarian method for the assignment problem," *Naval research logistics quarterly*, vol. 2, no. 1-2, pp. 83–97, 1955.

[14] N. Wojke, A. Bewley, and D. Paulus, "Simple online and realtime tracking with a deep association metric," in *2017 IEEE international conference on image processing (ICIP)*. IEEE, 2017, pp. 3645–3649.

[15] S. Scheidegger, J. Benjaminsson, E. Rosenberg, A. Krishnan, and K. Granström, "Mono-camera 3d multi-object tracking using deep learning detections and pmbm filtering," in *2018 IEEE Intelligent Vehicles Symposium (IV)*. IEEE, 2018, pp. 433–440.

[16] S. Sharma, J. A. Ansari, J. K. Murthy, and K. M. Krishna, "Beyond pixels: Leveraging geometry and shape cues for online multi-object tracking," in *2018 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2018, pp. 3508–3515.

[17] J. Luiten, T. Fischer, and B. Leibe, "Track to reconstruct and reconstruct to track," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1803–1810, 2020.

[18] P. Li, T. Qin *et al.*, "Stereo vision-based semantic 3d object and ego-motion tracking for autonomous driving," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 646–661.

[19] A. Osep, W. Mehner, M. Mathias, and B. Leibe, "Combined image- and world-space tracking in traffic scenes," in *2017 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2017, pp. 1988–1995.

[20] Z. Lu, V. Rathod, R. Votel, and J. Huang, "Retinatrack: Online single stage joint detection and tracking," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 14 668–14 678.

[21] M. Chaabane, P. Zhang, J. R. Beveridge, and S. O'Hara, "Deft: Detection embeddings for tracking," *arXiv preprint arXiv:2102.02267*, 2021.

[22] P. Li and J. Jin, "Time3d: End-to-end joint monocular 3d object detection and tracking for autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 3885–3894.

[23] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," *Advances in neural information processing systems*, vol. 27, 2014.

[24] K. Sohn, X. Yan, and H. Lee, "Learning structured output representation using deep conditional generative models," in *Proceedings of the 28th International Conference on Neural Information Processing Systems-Volume 2*, 2015, pp. 3483–3491.

[25] H. Cheng, M. Liu, L. Chen, H. Broszio, M. Sester, and M. Y. Yang, "Gatraj: A graph-and attention-based multi-agent trajectory prediction model," *arXiv preprint arXiv:2209.07857*, 2022.

[26] M. Liu, H. Cheng, L. Chen, H. Broszio, J. Li, R. Zhao, M. Sester, and M. Y. Yang, "Laformer: Trajectory prediction for autonomous driving with lane-aware scene constraints," *arXiv preprint arXiv:2302.13933*, 2023.

[27] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi, "Social gan: Socially acceptable trajectories with generative adversarial networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 2255–2264.

[28] O. D. Altan, G. Wu, M. J. Barth, K. Boriboonsomsin, and J. A. Stark, "Glidepath: Eco-friendly automated approach and departure at signalized intersections," *IEEE Transactions on Intelligent Vehicles*, vol. 2, no. 4, pp. 266–277, 2017.

[29] V. Kosaraju, A. Sadeghian, R. Martín-Martín, I. Reid, H. Rezatofighi, and S. Savarese, "Social-bigat: Multimodal trajectory forecasting using bicycle-gan and graph attention networks," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[30] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio, "Graph attention networks," *arXiv preprint arXiv:1710.10903*, 2017.

[31] P. Zhang, W. Ouyang, P. Zhang, J. Xue, and N. Zheng, "Sr-lstm: State refinement for lstm towards pedestrian trajectory prediction," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 085–12 094.

[32] Y. Yuan, X. Weng, Y. Ou, and K. M. Kitani, "Agentformer: Agent-aware transformers for socio-temporal multi-agent forecasting," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 9813–9823.

[33] M. Liang, B. Yang, W. Zeng, Y. Chen, R. Hu, S. Casas, and R. Urtasun, "Pnpnet: End-to-end perception and prediction with tracking in the loop," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 553–11 562.

[34] X. Weng, Y. Yuan, and K. Kitani, "Ptp: Parallelized tracking and prediction with graph neural networks and diversity sampling," *IEEE Robotics and Automation Letters*, vol. 6, no. 3, pp. 4640–4647, 2021.

[35] N. Carion, F. Massa, G. Synnaeve, N. Usunier, A. Kirillov, and S. Zagoruyko, "End-to-end object detection with transformers," in *European conference on computer vision*. Springer, 2020, pp. 213–229.

[36] J. Gu, C. Hu, T. Zhang, X. Chen, Y. Wang, Y. Wang, and H. Zhao, "Vip3d: End-to-end visual trajectory prediction via 3d agent queries," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5496–5506.

[37] Y. Hu, J. Yang, L. Chen, K. Li, C. Sima, X. Zhu, S. Chai, S. Du, T. Lin, W. Wang *et al.*, "Planning-oriented autonomous driving," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 17 853–17 862.

[38] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[39] S. Pellegrini, A. Ess, K. Schindler, and L. Van Gool, "You'll never walk alone: Modeling social behavior for multi-target tracking," in *2009 IEEE 12th international conference on computer vision*. IEEE, 2009, pp. 261–268.

[40] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, "nuscenes: A multimodal dataset for autonomous driving," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2020, pp. 11 621–11 631.

[41] X. Weng and K. Kitani, "A baseline for 3d multi-object tracking," *arXiv preprint arXiv:1907.03961*, vol. 1, no. 2, p. 6, 2019.

[42] K. Bernardin, A. Elbs, and R. Stiefelhagen, "Multiple object tracking performance metrics and evaluation in a smart room environment," in *Sixth IEEE International Workshop on Visual Surveillance, in conjunction with ECCV*, vol. 90, no. 91. Citeseer, 2006.

[43] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese, "Social lstm: Human trajectory prediction in crowded spaces," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 961–971.

[44] Y. Chai, B. Sapp, M. Bansal, and D. Anguelov, "Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction," *arXiv preprint arXiv:1910.05449*, 2019.

[45] A. Sadeghian, V. Kosaraju, A. Sadeghian, N. Hirose, H. Rezatofighi, and S. Savarese, "Sophie: An attentive gan for predicting paths compliant to social and physical constraints," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 2019, pp. 1349–1358.

[46] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.

[47] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *Advances in neural information processing systems*, vol. 32, 2019.

[48] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[49] Y. Yuan, H. Cheng, and M. Sester, "Keypoints-based deep feature fusion for cooperative vehicle detection of autonomous driving," *IEEE Robotics and Automation Letters*, vol. 7, no. 2, pp. 3054–3061, 2022.

[50] R. Xu, H. Xiang, X. Xia, X. Han, J. Li, and J. Ma, "Opv2v: An open benchmark dataset and fusion pipeline for perception with vehicle-to-vehicle communication," in *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022, pp. 2583–2589.