

HDGT: Heterogeneous Driving Graph Transformer for Multi-Agent Trajectory Prediction via Scene Encoding

Xiaosong Jia, *Student Member, IEEE*, Penghao Wu, Li Chen, Yu Liu, Hongyang Li, *Senior Member, IEEE*, Junchi Yan, *Senior Member, IEEE*

Abstract—Encoding a driving scene into vector representations has been an essential task for autonomous driving that can benefit downstream tasks e.g. trajectory prediction. The driving scene often involves heterogeneous elements such as the different types of objects (agents, lanes, traffic signs) and the semantic relations between objects are rich and diverse. Meanwhile, there also exist relativity across elements, which means that the spatial relation is a relative concept and need be encoded in a ego-centric manner instead of in a global coordinate system. Based on these observations, we propose Heterogeneous Driving Graph Transformer (HDGT), a backbone modelling the driving scene as a heterogeneous graph with different types of nodes and edges. For heterogeneous graph construction, we connect different types of nodes according to diverse semantic relations. For spatial relation encoding, the coordinates of the node as well as its in-edges are in the local node-centric coordinate system. For the aggregation module in the graph neural network (GNN), we adopt the transformer structure in a hierarchical way to fit the heterogeneous nature of inputs. Experimental results show that HDGT achieves state-of-the-art performance for the task of trajectory prediction, on INTERACTION Prediction Challenge and Waymo Open Motion Challenge.

Index Terms—Autonomous Driving, Trajectory Prediction, Heterogeneous Graph Neural Network, Scene Understanding

arXiv:2205.09753v2 [cs.AI] 20 Jul 2023

1 INTRODUCTION

Autonomous driving is one of the most influential domains for AI applications in recent years. An autonomous driving system is made of a few key modules: perception/localization, prediction, planning, and control. The perception/localization module is responsible to extract visual and semantic information from raw sensor data and HD-Map. The outcomes of preceding functionalities are in compact yet highly unstructured forms, e.g. positions and states of surrounding vehicles, pedestrians and traffic lights, and relations between lanes. For downstream tasks such as trajectory prediction and planning, it is essential to obtain comprehensive semantic representations from those unstructured inputs. The conventional pipeline for the task is to first acquire a representation vector for each target agent, which should exhibit both ego moving patterns and states of surrounding environment. Then given the vector as input, the subsequent decoder could conduct prediction/planning in the desired form (raw trajectory/control signal/etc), combined with some prior assumptions or physical constraints. However, even deeply inspired by the success of deep learning literature in recent years, the encoding process needs to be investigated nonetheless. There is no de facto paradigm to encode a representative vector due to the following observations.

- X. Jia, P. Wu, H. Li and J. Yan are with Department of Computer Science and Engineering, and MoE Key Lab of Artificial Intelligence, Shanghai Jiao Tong University, and Shanghai AI Lab, Shanghai, China. L. Chen and Y. Liu are with Shanghai AI Lab, Shanghai, China. This work was partly done when X. Jia and P. Wu were interns at the Shanghai AI Lab.

E-mail: {jiaxiaosong, wupenghao, yanjunchi}@sjtu.edu.cn
{lichen, liuyu, lihongyang}@pjlab.org.cn,
Correspondence author: Junchi Yan

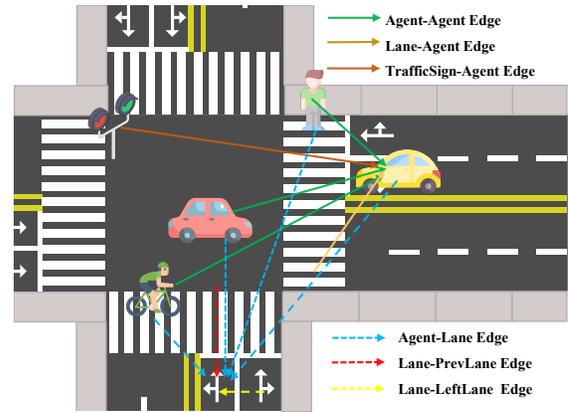


Fig. 1: Illustration of a heterogeneous driving graph. In the proposed framework, different types of agents and map elements are modeled as different types of nodes. The diverse semantic relations among them are modeled as different types of edges.

One inherent property that lies in the data is the *heterogeneity*, where we imply objects and their relations to have different types. The input data may include agents’ motion states, traffic lights with position, lane polylines with feasible movement connections, road line polygons with various types, etc. Fig. 1 illustrates an heterogeneous driving graph. It is non-trivial to encode those information comprehensively since they are highly unstructured, compared to the image input in the computer vision (CV) or the text input in the natural language processing (NLP). Thus, instead of directly the adopting widely used CNN or Transformers, a *good network design should consider a wide span of the input types and*

their semantic relations separately and explicitly.

There has been recent progress towards exploiting the heterogeneity of the driving scene. Trajectron++ [1] and HEAT-R [2] encode different types of agents with different parameters but they encode the map elements by rasterized representation, which leads to inferior performance compared to vector-based methods [3]. VectorNet based approaches [3]–[6] and SceneTransformer [7] put all agents and lanes in one single fully-connected graph and use shared parameters for information aggregation, which ignores different node types and semantic relations. LaneGCN based solutions [8]–[12] introduce a series of four graphs (LaneToAgent, LaneToLane, LaneToActor, and finally ActorToActor) according to distance heuristics and connective information of the lane. It demonstrates better performance compared to VectorNet and its extensions [3], [4], which suggests the importance of modeling heterogeneity. However, the sequential order of those four different relation aggregation functions in LaneGCN is designed manually by intuition.

The aforementioned issues make LaneGCN-like models not very flexible in face of datasets with different size, which makes it hard to simply apply the modern stacking-scaling [13]–[15] principle to increase the capacity of the model with more layers in the industry. To this end, we propose the **Heterogeneous Driving Graph Transformer**, namely **HDGT**, where the driving scene is modeled as a heterogeneous graph and with different types of nodes (agents or map elements) and the edges encode the the semantic relations between two nodes.

Specifically, we devise the aggregation and update function of the nodes and edges in a spirit from the recent success of Transformer [16]. To capture the heterogeneity, different sets of Transformer parameters are adopted for different node and edge types. Different from LaneGCN-like models, all nodes and edges are updated simultaneously. The same structure (Transformer) and simultaneous updating of HDGT make it easy to stack and scale when dealing with different data types and sizes. We conduct thorough ablation studies to demonstrate the importance of exploiting heterogeneity. We conduct experiments to verify the scalability of HDGT under different datasets with different size of models.

Besides heterogeneity, another import property of the driving scene data is the *relativity*, where we indicate each agent in the driving scene should process surrounding elements in its local coordinate system. If we adopt global reference, models would simply overfit the value of the absolute coordinates on the training scenes, and perform poorly on the unseen scenes. For example, to predict a vehicle’s future trajectory, the model should consider whether it would go straight or turn left or turn right, which is in relative reference instead of a global one. Hence, instead of processing the absolute coordinates, it is more reasonable to unanimously encode the relative relations among elements in a local view. Such a view shifts from other elements to the agent need be incorporated into the network explicitly. Most existing works choose a reference point [4], [7], [8] (for example, the ego agent) and then all the other elements including vehicles and lanes are in the reference point’s coordinate system. Though this way achieves relativity for the ego agent, we note that the representations of non-ego elements are not in their local coordinate systems. Consequently, as suggested in [17], to obtain the best results for multi-agent prediction, those methods have to forward multiple times with each agent as the reference agent, which brings notable additional cost. In this paper, *we aim to design a model where all agents’ are processing information in their local coordinate*

system and they are treated symmetrically so that predictions for the entire scene could be done in a single forward.

To this end, in our proposed HDGT approach, the spatial features of nodes as well as their in-edges¹ are in each node’s local, ego-centric system. As a result, the aggregation and updating of node features are done in their own reference, which captures the relative spatial state of other objects when their orientation is upfront and position is (0, 0). Then, for an edge ($u \rightarrow v$), since its feature represents u ’s states in v ’s coordinate system, we update it with the node feature of u (which is in u ’s coordinate system) combining with the coordinate transformation information from u to v . In the aforementioned two steps, we update node features in their local coordinate system and edge features back to its in-node’s coordinate system, which exploits the spatial relativity between elements and keeps all nodes and edges symmetric. Ablation studies in the experiment part will demonstrate the effectiveness of the proposed design.

Towards a practical multi-agent trajectory prediction engine for autonomous driving, the main highlights of this paper are:

- We devise an unified heterogeneous GNN for trajectory forecasting with HD-map and propose a hierarchical way of adopting Transformer to fit the heterogeneous nature of the inputs. The explicit modeling of all semantics and relations in the scene could improve the prediction accuracy by a large margin. The unified Transformer structure and simultaneous updating of all elements make it easy to stack and scale.
- Instead of choosing one reference agent, in a similar spirit with concurrent works [18], [19], we propose a symmetric way of encoding the spatial relationship between elements in the scene, which notably improves the generalization ability and enables the model to predict multiple agents’ future in one forward without performance drop.
- We examine the proposed HDGT on two challenging large-scale datasets and it achieves state-of-the-art performance on the corresponding online leaderboards. Thorough ablation studies are performed to evaluate the effectiveness of each designed component.

In particular, our extensive experimental results show that the proposed HDGT achieves new records on two challenging benchmarks, namely INTERACTION Prediction Challenge [20] and Waymo Open Motion Challenge [17]. Throughout 11/3/2021 (our final submission to INTERACTION) to 4/21/2022 (we published this report in arxiv), we still ranked the 1st and 2nd respectively in terms of minADE/minFDE metrics. Source code and trained models are publicly available at: <https://github.com/OpenPerceptionX/HDGT>.

2 RELATED WORK

2.1 Heterogeneous Graph Neural Networks

2.1.1 General Heterogeneous GNNs

Graph, as a common structure, is composed of nodes and edges, where a node usually represents an entity or element and a edge represents some kind of relation between the two connected nodes. To encode the structure and information of a graph, Graph neural network (GNN) [21] is proposed, which usually contains an aggregation function and an update function. The aggregation

1. In this paper, the heterogeneous graph we build is directional, which means the edge feature from A to B is different from B to A . We denote in-edges of A as all the edges going into A .

function is applied on each node/edge to collect information from their neighbors while the update function is to generate new representations for each node/edge according to their collected information and old representations. Two most widely used models are GCN [21] and GAT [22], which adopts the Laplacian matrix of the graph and attention for aggregation and updating respectively. For graphs with different types of nodes and edges, heterogeneous graph neural network is further introduced. R-GCN [23] uses GCN [21] like aggregating functions but with different parameters for different types, while HGN [24] and HetSANN [25] adopt GAT [22] like aggregation function with different parameters for different types. HetGNN [26] adopts the random walk to generate neighbors for nodes and uses distinct RNNs for different types of nodes. HGT [27] apply different Transformers for different relations on the academic graph.

2.1.2 Heterogeneous GNNs for Trajectory Prediction

In the heterogeneous graph of a driving scene, a node could either be an agent (vehicle, pedestrian, cyclist, etc) or a map element (lane, traffic light, stop sign, etc). A node usually contains some information (called node feature) such as agents' state (position, velocity), lane's type and the corresponding composing polyline, etc. One edge usually contains information (called edge feature) e.g. relative position between two entities or relations (Lane-LeftLane, Agent-Lane, etc). Trajectron++ [1] and HEAT-R [2] encode different types of agents with different parameters of GNNs or Transformer. However, as pioneering works, they encode the map elements by rasterized images and use a CNN to extract representation vectors of maps. According to the recent progress in the trajectory prediction community, it is inefficient and leads to much inferior performance compared to vector-based methods [3]. In fact, vector-based methods have dominated leaderboards of multiple public datasets Waymo [17], nuScenes [28], Argoverse [29] and INTERACTION [20]. Different from existing works, in HDGT, we constructs one heterogeneous graph including both agents and map elements and an unified Transformer structure with type-specific parameters is adopted for the aggregation and update function of heterogeneous graph.

2.2 Transformer

2.2.1 General Transformers

Transformer [16], as a set function to aggregate information of elements, has show promising success in multiple fields. The two core components of Transformer is the dot-product attention and the feed-forward network. The dot-product attention mechanism computes the pair-wise relativity of two elements by dot-product and updates the representations of each element by the normalized weighted sum of other elements. Emerging from the NLP field, BERT [13] and GPT series [30], [31] demonstrate the strong scalability of Transformers when handling huge amount of data, which could benefit from large scale pretraining. Later, Transformer is also shown to be able to scale well on image data, which starts the era of Vision Transformer (ViT) [32]. Besides, in fields including speech [33], multimodal learning [34], reinforcement learning [35], and etc, Transformer also achieves SOTA performance, which exhibits its wide applicability.

2.2.2 Transformers for Trajectory Prediction

Transformers have also recently been adopted in autonomous driving, especially for trajectory prediction. However, their map

elements are separately encoded through rasterized representation by a CNN, which loses parts of heterogeneity. VectorNet encoder [3] adopts attention mechanism on a fully connected graph with agents and map elements and SceneTransformer [7] additionally apply attention on the temporal axis. IDE-Net [36] explicitly extracts the interaction between agents by Transformer. However, the above works basically update representations of all agents and map elements with the same set of sharing parameters, ignoring the heterogeneity of nodes and relations. In our HDGT, instead of building a fully connected graph and sharing parameters for all nodes, we connect nodes by their semantic relation and conduct aggregation and updating of node and edge embeddings with type-specific parameters. Additionally, all nodes are processing information in their local coordinate systems. Experiments show that HDGT outperforms VectorNet/DenseTNT [5] and SceneTransformer [7] by a notable margin.

2.3 Trajectory Prediction for Autonomous Driving

2.3.1 Scene Encoder For Trajectory Prediction

Pioneering methods for trajectory prediction focus on the interaction among agents (called Social Force [37]) and do not take the HD-Map into consideration. Social LSTM [38] combines LSTM with social pooling. Social GAN [39] computes relative positions between the ego and all other people, rather than only considering people inside the grid as did in Social pooling. The authors in [40] consider the edge relations to better capture the interaction. The work [40] designs a GNN for agents detected from raw sensor data to capture their interactions. Trajectron [41] proposes a graph-structured encoder for the ease of predicting distributions of multi-agent trajectory. ST-GCN [42] adopts GCN on both spatial and temporal dimension.

However, in the driving scene, the movement of agents is not solely influenced by other agents. The map elements including lanes, traffic lights, and stop signs could all have significant influence on agents' future trajectories. Thus, it is necessary to take the map into consideration in a prediction model. Since the HD-map data is highly unstructured (polylines and their types), for deep learning based methods, one popular choice is to use the rasterized images to represent the scene [43]–[48]. Specifically, they project map elements into a top-down view image according to the 2D coordinates; different types of elements may be painted in different channels. In this way, the convolutional neural network (CNN) could directly applied on the rasterized image to extract the map information. As for the encoder for agents, one choice is to apply RNN and GNN similar to pioneering works and then combine it with the rasterized representation of map. Following this branch, Trajectron++ [1] improves Trajectron with the rasterized HD-Map and heterogeneous agents type. [2] builds an agent-only heterogeneous graph neural network whose edges including the relative information. [49] designs a memory network to store the scene knowledge. [50] uses multiple sensors' data for trajectory prediction. Different from the aforementioned works, [51] adopts the rasterization for agents as well and further formulates the output on the rasterized images. With the blossom of CNN structure in vision, this raster-based methods could outperform those without map information.

Nevertheless, as demonstrated in [3], the raster-based methods lose too much information during painting and are inefficient to generate good representation, since the road elements are sparse in image space. Besides, it is difficult for CNN to capture

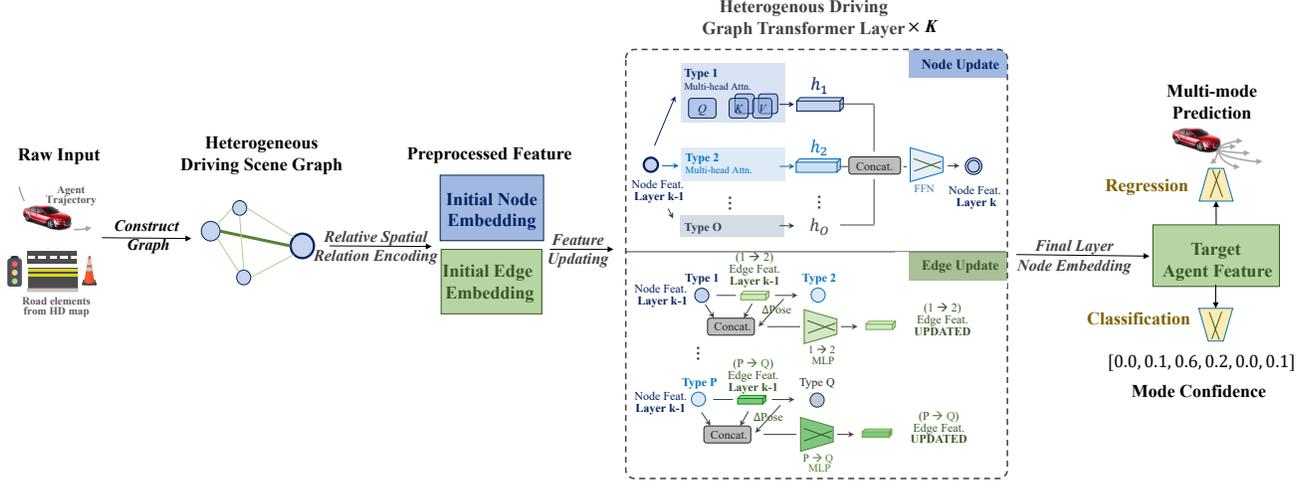


Fig. 2: **The pipeline of HDGT.** We first construct heterogeneous driving scene graph according to agents’ and map elements’ types and relations. Then, we initialize the node and edge feature of the heterogeneous graph by relative spatial encoding. Next, we update the node and edge features by the Heterogeneous Driving Graph Transformer Layer. The updating is conducted K times by K layers. Finally, we take the target agent nodes’ feature at the last layer and feed them the regression and classification MLP to obtain the multi-modal trajectory prediction with confidence scores.

global information due to its limited effective perception field issue [52]. Most recently, encoding the map in a vector-based spirit [3], [8] has proven great performance in multiple competitions including Waymo [17], nuScenes [28], Argoverse [29] and INTERACTION [20]. Works in this branch typically employ 1D CNN or LSTM [53] to process the temporal data, adopt the PointNet [54] to process polylines, and use graph neural network (GNN) to handle the relation among all agents and elements. VectorNet [3] uses sub-graph for lane and agent encoding and a fully connected global graph to capture their relations. Scene-Transformer [7] proposes a factorized spatial-temporal network which applies Transformer on the fully-connected spatial/temporal graph alternatively. LaneGCN [8] introduces a series of four different information aggregations modules for the driving scene. TPCN [55] treats all coordinate data as point clouds and achieves comparable performance with SOTA works.

2.3.2 Output Decoder for Trajectory Prediction

Besides those works focusing on extracting high-quality representations of the scene (i.e. the encoder part), there are also works explore the output formulation and training paradigm for the trajectory prediction, i.e., the decoder part. Social LSTM [38] assumes an independent bivariate Gaussian distribution for the future position of each time-step of each agent. Social GAN [39] further introduces GAN to improve the diversity of prediction. Trajectron [41] combines recurrent sequence modeling and variational deep generative modeling. Trajectron++ [1] additionally discusses the paradigm of future-conditional predictions. Prank [56] produces the conditional distribution of agent’s trajectories plausible in the given scene. [57] aims to model the joint distribution over future trajectories via an implicit latent variable model with GNN for agents only. TNT [4] introduces goal-based predictions for diversity and DenseTNT [5] further designs dense goal sets and a corresponding optimization pipeline to further improve its performance. LaneRCNN [9] proposes local LaneROI and based on it goal-based prediction is used. [10] proposes to predict multi-modal future conditioned on lane connecting topology. [12]

proposes the graph-based heatmap output paradigm to optimize the missing rate while [11] designs a hierarchical heatmap paradigm and a learnable trajectory recombination method for joint multi-agent trajectory prediction. DCMS [58] introduces dual consistency and multi-pseudo-target supervision for training. [59] conducts refinements with temporal priors.

In this work, we focus on the encoder part, which aims to extract the massive diverse information from the scene comprehensively. Comparing to the existing works, we are the first to comprehensively model all node types and all semantic and geometric relations into one single heterogeneous graph and update the node and edge features in a unanimous and parallel way but with distinct parameters for different node or edge types. Additionally, we propose to process all nodes in their own local coordinate system, which enhances the generalization ability of the model and enables the single-forward multi-agent prediction feasible. Notably, concurrent works Multipath++ [18]² has similar designs where all the features in each agent’s local coordinate system. HiVT [19] proposes to process the scene in a hierarchical way with a local part and a global interaction module by Transformer. With similar ideas of the local coordinate system, they both achieve new records on the public leaderboards, which demonstrates importance of relativity as mentioned in Sec. 1.

3 METHODOLOGY

3.1 Problem Formulation and Approach Overview

We introduce our method. The input of the trajectory prediction includes two parts: agents’ history states and map elements. Agents’ history states is composed of past L time-steps in a certain frequency (for example, 10 Hz) where at each time-step agents’ coordinate, velocity, heading, etc are provided. Besides, the size of bounding box (width, length, height) and the type (pedestrian, cyclist, vehicle) of agents are given as well. Map elements are

2. Specifically, Multipath++ [18]’s preprint was initially uploaded to Arxiv on 29/11/2021. While we submitted our results to INTERACTION Leaderboard on 03/11/2021 as recorded on the INTERACTION Leaderboard.

surrounding objects and signs which reflect traffic rules and thus have nonnegligible influence on agents' future movements. The commonly given map elements include: lane, road line, crosswalk, stop sign, traffic light, speed bump, curb, etc. Each map element contains a polyline or a polygon to indicate its physical position as well as other information: for traffic lights, their state (green, yellow, red) at each time-step are given; for traffic lights and stop signs, the lanes which they have influence on are marked; for lanes, their predecessors, successors, and neighbors are given.

The goal of the trajectory prediction task is to predict the future trajectories of agents in the next T time-steps. To better capture the uncertainty of predicting future, multiple different predictions together with confidence scores are allowed for each agent.

The pipeline of HDGT is shown in Fig. 2, which involves:

- Sec. 3.2 describes how we construct the heterogeneous graph of the driving scene based on the raw input.
- Sec. 3.3 gives how we transform the the node and edge feature into local reference system.
- Sec. 3.4 describes how Transformer structure is adopted to aggregate and update the node and edge features.
- Sec. 3.5 describes the heads and loss used to do the prediction based on the extracted agent node features.

3.2 Construction of the Heterogeneous Graph

Preliminary In this work, we represent each driving scene with multiple agents and map elements as one directed heterogeneous graph. A directed graph is represented as $\mathcal{G} = \{\mathcal{V}, \mathcal{E}\}$, where \mathcal{V} denotes the set of nodes and \mathcal{E} denotes the set of edges. A heterogeneous graph has a node type mapping function of $\tau(v) : \mathcal{V} \rightarrow \text{Type}_v$ and an edge type mapping function $\phi(e) : \mathcal{E} \rightarrow \text{Type}_e$, where Type_v and Type_e denote the set of node and edge type respectively. In HDGT, the type of an edge $e : u \rightarrow v$ is determined by the type of source node u and destination node v and their relations. When building the heterogeneous graph, all possible combinations of semantics relations are considered.

For an agent node, instead of connecting it to all the other counterparts like VectorNet [3], LaneGCN [8] constrains each node only connecting to those whose distance is smaller than a threshold. However, in some scenarios, vehicles traversed a long way in the prediction horizon. For example, in Waymo Open Motion, some vehicles move more than 300 meters in the 8 seconds prediction horizon. It is computationally expensive and unnecessary to set a fixed large threshold for *all* agents. For example, for pedestrians, vehicles very far away would not influence their future behaviors.

For all agents' node, we calculate their distance with others by their position at the final observed time-step; for all nodes with the polyline shape, we calculate by the middle point of the polyline; for all nodes with a polygon shape, we calculate by their geometric centre. Thus, similar to [9], for each agent, we set the distance threshold flexible to be its speed times the prediction horizon with a type-specific buffer value.

For a lane or traffic signs node, we inversely connect them to the agent nodes according to the threshold mentioned in the paragraph above. Besides, for lane nodes, we connect them with their surrounding lane nodes and there are four types - left, right, entry, and exit. Additionally, since the length of each lane could vary a lot, we divide those lanes whose length is larger than a threshold into short ones with similar length so that all lane nodes have similar inputs. We give the pseudo code for the construction of the driving scene heterogeneous graph in Alg. 1

Algorithm 1: Construct Heterogeneous Graph for Driving Scene Encoding

Input: N_a agents' position, N_l lanes' positions and their connective topology, N_{tr} traffic signs positions

Output: Node set \mathcal{V} , adjacency matrix $\mathbf{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$

```

1 begin
2   Divide lanes with length larger than threshold  $\epsilon_{\text{lane}}$ 
   into  $\lceil \frac{\text{length}}{\epsilon_{\text{lane}}} \rceil$  pieces.
3   Node set  $\mathcal{V} \leftarrow \text{Agents} \cup \text{Lanes} \cup \text{Traffic Signs}$ 
4   Initialize Adjacency Matrix  $\mathbf{A} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{V}|}$  with zeros.
5   for  $i = 1, 2, \dots, N_a$  do
6     Retrieve agent node  $v_{a,i}$ 's position  $\text{pos}_{a,i}$ ,
       distance threshold  $\epsilon_{a,i}$ , and index  $\text{idx}_{a,i}$  in  $\mathcal{V}$ 
7     for  $j = 1, 2, \dots, |\mathcal{V}|$  do
8       Retrieve node  $v_j$ 's position  $\text{pos}_j$ 
9       if  $\text{distance}(\text{pos}_{a,i}, \text{pos}_j) < \epsilon_{a,i}$  then
10         $\mathbf{A}[\text{idx}_{a,i}, j] \leftarrow \text{Relation}(\text{Agent-Type}(j))$ ;
11         $\mathbf{A}[j, \text{idx}_{a,i}] \leftarrow \text{Relation}(\text{Type}(j)\text{-Agent})$ ;
12  for  $i = 1, 2, \dots, N_l$  do
13    Retrieve lane node  $v_{l,i}$ 's connective topology
       $\text{Connected}$  and index  $\text{idx}_{l,i}$  in  $\mathcal{V}$ ;
14    for  $\text{Lane node } v_{l,j} \in \text{Connected}(v_{l,i})$  do
15      Retrieve  $v_{l,j}$ 's connective relation with  $v_{l,i}$  -
       $\text{CR}$  and index  $\text{idx}_{l,i}$  in  $\mathcal{V}$ ;
16       $\mathbf{A}[\text{idx}_{l,i}, \text{idx}_{l,j}] \leftarrow \text{Relation}(\text{Lane-CR})$ ;
17  return  $\mathcal{V}, \mathbf{A}$ ;
```

Accordingly, the node types include agents, lanes, and traffic signs (stop signs, road lines, cross walks). The edge types include Agent-Lane, Agent-Trafficsign, Lane-Agent, Lane-NextLane, Lane-PreviousLane, Lane-LeftLane, Lane-RightLane, and Trafficsign-Agent.

3.3 Relative Spatial Relation Encoding

After constructing the heterogeneous graph of the driving scene, we need to initialize the feature embedding of all nodes and edges.

As mentioned in Sec. 1, it is undesirable to directly feed raw data to the neural network since the raw spatial data are usually in a global coordinate system defined by datasets. Considering same set of spatial relations could be represented differently by the coordinate system with different origins and orientations, the spatial data should be normalized into a canonical form. Otherwise, non-normalized inputs would inappropriately force the network to overfit the absolute value of training data and thus unfavorably cause poor generalization. To normalize the data, one popular choice is to select one agent as a global reference [3], [8]. Though this way achieves relativity for the ego agent, we note that the representations of non-ego elements are not in their local coordinate systems. For example, assume there are three vehicles in the scene. If we choose vehicle 1's current coordinate as the origin and heading as the positive y-axis, then in vehicle 1's view, all vehicles moves relative to its current state. However, in vehicle 2 and vehicle 3's view, they observe other vehicle's moving in a biased coordinate system which is un-normalized and could lead to worse generalization ability. Consequently, as suggested in [17], to obtain the best results for multi-agent prediction, those methods

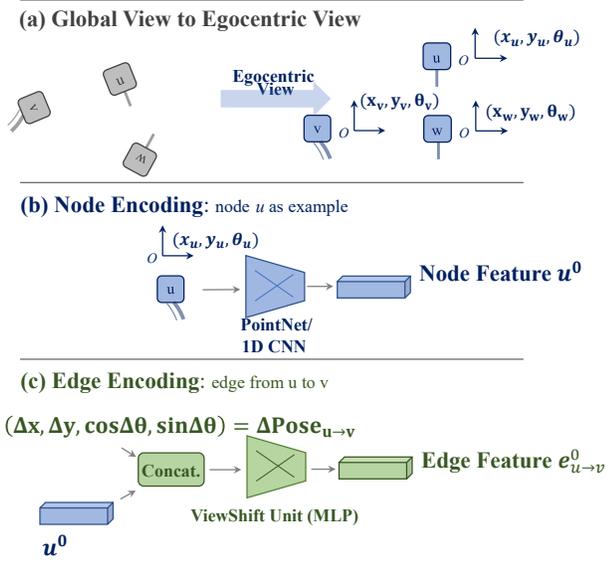


Fig. 3: Process of the Relative Spatial Relation Encoding which encodes both node and edge features in their own local reference system. (a) Spatial information encoding of the node, where we transform all nodes in their ego-centric coordinate system. (b) To transform node feature into a single compact vector instead of one trajectory/polyline/polygon, we adopt 1D CNN or PointNet. (c) Transform source node’s representation vector into target node’s in-edge feature by concatenating with the delta between two coordinate systems and feeding them into an MLP.

have to forward multiple times with each agent as the reference agent, which brings notable additional cost. Thus, we should design a model where all agents are processing information in their local coordinate system so that they could all be treated symmetrically into a canonical form. In this way, predictions for the entire scene could be done in a single forward.

For nodes which represent either agents or map elements, we initialize their feature embeddings in their own ego-centric coordinate system. Specifically, for each node v , we define its unique reference as a tuple $(x_{\text{ref}}^v, y_{\text{ref}}^v, \theta_{\text{ref}}^v)$ where x_{ref}^v and y_{ref}^v represents the origin of the coordinate system and θ_{ref}^v represents the orientation. For all coordinates feature of node v , we subtract $\begin{bmatrix} x \\ y \end{bmatrix}$ by $\begin{bmatrix} x_{\text{ref}}^v \\ y_{\text{ref}}^v \end{bmatrix}$ so that all these coordinates are relative to itself. Then, for all coordinate and velocity features, we rotate them by $-\theta_{\text{ref}}^v$, i.e., we let the rotation matrix $R = \begin{bmatrix} \cos -\theta_{\text{ref}}^v & -\sin -\theta_{\text{ref}}^v \\ \sin -\theta_{\text{ref}}^v & \cos -\theta_{\text{ref}}^v \end{bmatrix} = \begin{bmatrix} \cos \theta_{\text{ref}}^v & \sin \theta_{\text{ref}}^v \\ -\sin \theta_{\text{ref}}^v & \cos \theta_{\text{ref}}^v \end{bmatrix}$ dot product with either $\begin{bmatrix} x \\ y \end{bmatrix}$ or $\begin{bmatrix} \text{velocity}_x \\ \text{velocity}_y \end{bmatrix}$. Besides, for agents’ node with heading feature, we subtract the them with θ_{ref}^v . In this way, all the node features are composed of information relative to itself. Note that for agent nodes, we set $(x_{\text{ref}}^v, y_{\text{ref}}^v)$ as its position at the final observed time-step and θ_{ref}^v as its heading at the final observed time-step, which means each agents’ history trajectory and velocity are relative to its current state. During training, for labels - future trajectories of agent, we conduct the transformation as well so that the regression target are relative as well. For map elements composed of a polyline, we select the middle point’s position as $(x_{\text{ref}}^v, y_{\text{ref}}^v)$ and the direction of the vector from its start point to its end point as the θ_{ref}^v . For

map elements composed of a polygon, we set $(x_{\text{ref}}^v, y_{\text{ref}}^v)$ as the geometric centre of the polygon and θ_{ref}^v as the direction of the longest edge of the polygon. In this way, all the map elements are represented in a canonical form. Fig. 3(a) illustrates how we transform all node features in their local reference system.

Until now, all node features are in their ego-centric coordinate system. However, they are still unstructured, i.e., they could be a sequence of history states for agents, a polyline or a polygon. For the ease of efficient parallel operations in the graph neural network, we transform them into single vectors with the same length first. Following the common practice in [3], [8]: for all agent nodes, we adopt a shared 1D-CNN on each of their features to encode their temporal dynamics. The 1D CNN is composed of a series of 1D convolution layers, downsampling and finally pooling to obtain a single vector representation. For all map element nodes, we adopt a simplified version of PointNet. Specifically, for a set of points from either a polyline or a polygon, a PointNet layer is composed of a shared MLP for each point, an max-pooling operation over all points to obtain the overall representation, and finally concatenating the pooled vector with the vector after MLP. By stacking several layers of PointNet (here we set 3 layers) and finally adopting a max-pooling operation, the geometry information of a polyline or a polygon is encoded in a single vector. Fig. 3(b) shows the process that we encode node features in their ego-centric view by 1D CNN or PointNet. Additionally, sub-type information such as type of lanes (single, double, etc) and road lines (broke-broke, solid-broke, etc), we adopt learnable embeddings for them and then concatenate them with the geometry feature obtained from 1D CNN or PointNet. We denote the initial node feature set as $\{v^0\}$ for all nodes v in \mathcal{V} .

For edges in the directed graph, we initialize their feature embeddings in their target nodes’ coordinate system. Formally, for an edge $u \rightarrow v$, we refer u as its origin node and v as its target node. The edge feature represents node u ’s information in the view of node v . For example, one edge feature $e_{u \rightarrow v}$ might represent a lane u in the left in agent v ’s view. Recall that we already have the representation vector of u , as shown in Fig. 3(b). However, the vector of node $u - u^0$, is in u ’s local coordinate system instead of v ’s. Thus, we need to transform u^0 into the local coordinate system of v . Denote the reference of node u and v as $(x_{\text{ref}}^u, y_{\text{ref}}^u, \theta_{\text{ref}}^u)$ and $(x_{\text{ref}}^v, y_{\text{ref}}^v, \theta_{\text{ref}}^v)$ respectively. To transform coordinate from u to v ’s coordinate system, the required information is $(x_{\text{ref}}^v - x_{\text{ref}}^u, y_{\text{ref}}^v - y_{\text{ref}}^u, \cos(\theta_{\text{ref}}^v - \theta_{\text{ref}}^u), \sin(\theta_{\text{ref}}^v - \theta_{\text{ref}}^u)) = (\Delta x_{u \rightarrow v}, \Delta y_{u \rightarrow v}, \cos \Delta \theta_{u \rightarrow v}, \sin \Delta \theta_{u \rightarrow v})$, denoted as $\Delta \text{Pose}_{u \rightarrow v}$. For simpleness, we adopt MLP layers to obtain the initial edge feature $e_{u \rightarrow v}^0$ by:

$$e_{u \rightarrow v}^0 = \text{MLP}_{r(u)}(\text{concat}[u^0, \Delta \text{Pose}_{u \rightarrow v}]) \quad (1)$$

where $r(u)$ means that the MLP used for the transformation depends on the type of source node u considering the fact that different node types (agent, lane, traffic sign) have different geometry characteristics in their node features. In this way, we provide in-edge features of v with information related to v ’s coordinate system and this operation is symmetric for all nodes. We denote this MLP as *ViewShift Unit* as demonstrated in Fig. 3(c)

In a nutshell, by the relative spatial relation encoding procedure mentioned in this section, we decouple the spatial encoding process into two parts: First, the motion/geometric characteristics of a single object – node feature encoding Second, the spatial relation between two objects – edge feature encoding. As a result, the inputs for the neural network are all in the ego-centric

coordinate system, which satisfies the i.i.d input assumption and is less prone to overfitting absolute coordinate value. To be more specific, assume two vehicles in the same scene have exactly the same patterns. One vehicle moves from (0, 0) to (5, 5) and the other moves from (100, 100) to (105, 105). Choosing any vehicle as the reference would cause the encoder having to deal with inputs at different scales while the encoder still needs to output similar features since they share the same patterns. In the proposed way, the node features are exactly the same while the in-edge features contain the relative relation. We could observe the symmetry within both node features and edge features, which is beneficial for the model’s generalization ability since the same type of inputs share parameters.

3.4 Heterogeneous Driving Graph Transformer Layer

After building the heterogeneous driving graph and initialize node and edge features with relative spatial encoding, the next step is to conduct information exchange among nodes and edges and update their corresponding features. In this way, the complex relations among agents and map elements could be encoded so that we could feed the predicted agents’ node feature to the decoder head to generate their future multi-mode trajectories.

3.4.1 Preliminaries

Transformer: Transformer [16] is a permutation-invariant set function which can fuse the information among a set in pair-wise way by dot-product attention. It mainly consists of two parts: the multi-head attention module to fuse information from others and the feed-forward module to update information of themselves. Formally, suppose we have a set of m features with dimension h : $\mathbf{Z} = \{z_i | i = 1, 2, \dots, m\} \in \mathbb{R}^{m \times h}$. In the attention mechanism, three linear transformations ($\mathbf{W}_q \in \mathbb{R}^{h \times h_q}$, $\mathbf{W}_k \in \mathbb{R}^{h \times h_k}$, $\mathbf{W}_v \in \mathbb{R}^{h \times h_v}$, $h_q = h_k$) are conducted to generate the three vectors for each element in the set, called query (Q), Key (K), Value (V):

$$\mathbf{Q}, \mathbf{K}, \mathbf{V} = \mathbf{Z}\mathbf{W}_q, \mathbf{Z}\mathbf{W}_k, \mathbf{Z}\mathbf{W}_v \quad (2)$$

Then, a dot-product operation followed by a Softmax operation is conducted between \mathbf{Q} and \mathbf{K} to generate pair-wise importance weights $\mathbf{A} \in \mathbb{R}^{m \times m}$ which is called attention matrix:

$$\mathbf{A} = \text{Softmax}(\mathbf{Q}\mathbf{K}^\top / \sqrt{h_q}) \quad (3)$$

where $\sqrt{h_q}$ is for the numerical stability of the Softmax. Finally, the weighted sum is performed with weights \mathbf{A} and value \mathbf{V} :

$$\mathbf{Z}' = \mathbf{A}\mathbf{V} \quad (4)$$

Thus, we define the attention function as $\mathbf{Z}' = \text{Attn}(\mathbf{z})$. The multi-head attention (MHA) extends the attention by applying multiple attention on the same set of inputs under different parameters (called different heads) in parallel and then applying a linear transformation, namely $\mathbf{W}_o \in \mathbb{R}^{n_h h_v \times h_v}$, where n_h is the number of heads, to integrate the heads:

$$\mathbf{Z}'' = \text{MHA}(\mathbf{Z}) = [\text{Attn}_1(\mathbf{Z}), \dots, \text{Attn}_{n_h}(\mathbf{Z})]\mathbf{W}_o \quad (5)$$

The feed-forward module is a 2-layer MLP with activation $\delta(\cdot)$:

$$\mathbf{Z}''' = \text{FFN}(\mathbf{Z}'') = \delta(\mathbf{Z}\mathbf{W}_1 + \mathbf{b}_1)\mathbf{W}_2 + \mathbf{b}_2 \quad (6)$$

The \mathbf{Z}''' is the output of one Transformer layer and it usually shares exactly the same hidden dimension with \mathbf{Z} for the ease

of adding residual connection. Note that within MHA and FFA, residual connections [60] and layer norm (LN) [61] are applied for the stability of training:

$$\begin{aligned} \mathbf{Z}'' &= \text{MHA}(\text{LN}(\mathbf{Z})) + \mathbf{Z} \\ \mathbf{Z}''' &= \text{FFN}(\text{LN}(\mathbf{Z}'')) + \mathbf{Z}'' \end{aligned} \quad (7)$$

Graph Neural Network (GNN): in HDGT, the information exchange and feature updating are fulfilled by GNN. We give the definition of the aggregation function and update function [62].

Definition 3.1. Aggregation Function: An aggregation function ρ is a map which takes a set as input, and reduces it to a single element which represents the aggregated information. The aggregation functions must be invariant to permutations of their inputs, and should take variable numbers of arguments (e.g., element-wise summation, mean, maximum, etc).

Definition 3.2. Update Function: An update function ϕ is a map conducted on all elements in the same set (for example, all agent nodes is a set while all lane nodes is another set) and serves as an integration of information from ego and neighbors.

In HDGT, for the heterogeneous graph, there are two types of aggregation function and two types of update function: ρ^V is applied on nodes and aggregates features from each node’s in-edges; ρ^E is applied on edges and aggregates features from each edge’s source node. ϕ^V and ϕ^E update node and edge features respectively based on aggregated information and their current features.

In the following sections, we introduce how we adopt Transformer as the aggregation and update function of GNN in the spirit of heterogeneity and relativity. In HDGT, there are K layers of GNN, which means the node and edge updating mentioned below are conducted K times to extract the high-order relation among agents and map elements.

3.4.2 Node Aggregation and Update

For the aggregation function of nodes ρ^V , we adopt the multi-head attention (MHA) of Transformer to collect information from each node’s in-edges which is their own ego-centric coordinate system. However, Transformer is initially designed for language data, which assumes all elements are homogeneous - words. This assumption does not hold true for the driving scene data since in-edges could be agents (dynamic states) or lanes (connective topology) or traffic signs (traffic rules). Therefore, it is unreasonable to let all in-edges share the same set of parameters for the aggregation function. Additionally, since an agent node’s neighbor is decided by its speed as well as a buffer value, we find that some agent nodes may have thousands of in-edges. Consequently, due to the sparse nature of the attention mechanism, it is difficult to aggregate information from rare but significant nodes (for example, crosswalk) and thus leads to degenerated performance.

To alleviate the aforementioned issues, we propose a hierarchical manner to adopt Transformer by first processing different types of semantic relations separately by attention and then integrating information from different types with an MLP. For the aggregation of in-edges’ features, each type of in-edges shares a set of parameters of MHA and does cross-attention by letting the node features as \mathbf{Q} and the in-edge features as \mathbf{K} and \mathbf{V} . As a result, there would be a set of attended features - each for a specific edge type. Finally, we feed the concatenation of the set into an MLP to obtain the aggregated features of each node.

Formally, for a node v with a set of m in-edges $\{e_i\}_{i=0}^{i=m}$ in the k^{th} layer, the input is node feature \mathbf{v}^{k-1} and its in-edge features $\{\mathbf{e}_i^{k-1}\}_{i=0}^{i=m}$ from the last layer. Note that for layer 1, we use the features initialized by relative spatial relation encoding in Sec. 3.3. Denote the features of all in-edges of v with a specific edge type j from $(k-1)^{th}$ layer as $\mathbf{E}_j^{k-1} \in \mathbb{R}^{m_j \times h}$ where m_j is the number of type- j in-edges and h is the hidden dimension, we calculate the attention vector of this specific type for this node as:

$$\mathbf{a}_j = \text{Softmax}(\mathbf{q}_j \mathbf{K}_j)^\top / \sqrt{h}, \in \mathbb{R}^{m_j} \quad (8)$$

where $\mathbf{q}_j = \mathbf{v}^{k-1} \mathbf{W}_{j,\text{query}} \in \mathbb{R}^h$, $\mathbf{K}_j = \mathbf{E}_j \mathbf{W}_{j,\text{key}} \in \mathbb{R}^{m_j \times h}$, $\mathbf{W}_{j,\text{query}}$ and $\mathbf{W}_{j,\text{key}}$ are the query and key linear transformation matrix for edge type j . In this way, the attention score is calculated in a type-specific way, which utilizes the prior of heterogeneity. For example, for an agent to drive safely, its attention for other agents and for lanes should be different and treated separately.

For all types of relations, we conduct MHA with different parameters in the aforementioned way, which outputs a set of aggregated features $\{\mathbf{v}_j^{k-1,\prime} \in \mathbb{R}^h | j = 1, \dots, n_r\}$ where n_r is the number of possible in-edge types for v 's node type. To integrate them, an MLP layer is adopted to obtain the final aggregated node feature $\mathbf{v}^{k-1,\prime\prime}$ as:

$$\mathbf{v}^{k-1,\prime\prime} = \text{MLP} \left(\text{concat} \left(\{\mathbf{v}_j^{k-1,\prime} | j = 1, \dots, n_r\} \right) \right) \quad (9)$$

Here, the operation of the aggregation function ρ^V is finished.

As for the update function of nodes ϕ^V , we adopt FFN to update all node features. Considering the heterogeneous nature of nodes, we adopt different parameters for different types of nodes:

$$\mathbf{v}^{k+1} = \text{FFN}_{\text{Type}(v)}(\mathbf{v}^{k-1,\prime\prime}) \quad (10)$$

where \mathbf{v}^{k+1} is the updated features of node v . Note that for the clarity of the expression, we do not write the residual connection and layer norm layers in the equations above while we do adopt them following [16].

3.4.3 Edge Aggregation and Update

Recall that edge feature $\mathbf{e}_{u \rightarrow v}$ represents element u 's representation in element v 's local view. Thus, to update $\mathbf{e}_{u \rightarrow v}$, we need to combine node u 's feature and transform it to node v 's view.

Formally, assume at layer k we have \mathbf{u}^{k-1} , \mathbf{v}^{k-1} , and $\mathbf{e}_{u \rightarrow v}^{k-1}$ which represents the node feature of u and v and their edge feature on last layer respectively. Note that for layer 1, we use the features initialized by relative spatial relation encoding in Sec. 3.3.

For the aggregation function of edge ρ^E , we retrieve node u 's feature \mathbf{u}^{k-1} and the transform between node u and v 's coordinate system $\Delta \text{Pose}_{u \rightarrow v}$ as defined in Sec. 3.3.

For the update function of edge ϕ^E , we concatenate the aggregated information \mathbf{u}^{k-1} and $\Delta \text{Pose}_{u \rightarrow v}$ with its previous feature $\mathbf{e}_{u \rightarrow v}^{k-1}$ and feed them into an MLP to obtain the updated edge feature $\mathbf{e}_{u \rightarrow v}^k$. Again, for different edge types, we adopt different MLPs hence the heterogeneity is ensured:

$$\mathbf{e}_{u \rightarrow v}^k = \text{MLP}_{\text{Type}(u \rightarrow v)} \left(\text{concat} \left[\mathbf{u}^{k-1}, \Delta \text{Pose}_{u \rightarrow v}, \mathbf{e}_{u \rightarrow v}^{k-1} \right] \right) \quad (11)$$

Similarly, we adopt the residual connection and layer norm in the update function as well and we do not write it for brevity.

In summary, after conduct the aforementioned operations for all nodes and edges for once, higher order of relation information with semantics could be extracted. Note that all aggregation functions and update functions are similar to the ones used

in Transformer [16], which enables the model to exploit the scalability by simply stacking layers while the modification for heterogeneity and relativity injects the task-specific prior into the structure. In Fig. 2, Heterogeneous Driving Graph Transformer Layer part gives an visualization about how we conduct node and edge feature updating at each layer.

3.5 Output Head and Training Objective

For the output head, we adopt the commonly used regression-MLP plus classification-MLP combination. It falls in the formulation of the MTP loss [47]. Different from the original one used in [47] which selects the winner trajectory by a distance function that considers an angle between the last points of the two trajectories as seen from the actor position, we follow the recent practice in [3], [8] and adopt the average Euclid distance as the selection criterion for the winner. Specifically, for each target agent, we use the node feature at HDGT's last layer as its hidden representation. An MLP module is used for regression; it outputs a tensor of size $N \times K \times T \times 2$, where N is the number of target agents, K is the number of required mode, T is the prediction length, and 2 corresponds to the dimension size of (x, y) . Another MLP is used to generate the confidence of each mode for each agent - a tensor of size $N \times K$. For different types of agents (vehicles/pedestrians/cyclists), different parameters for the two heads are used to capture the diverse moving behaviors. Denote an agent's outputs as:

$$O_{\text{reg}} = \{(\mathbf{p}_1^k, \mathbf{p}_2^k, \dots, \mathbf{p}_T^k)\}_{k \in [1, K]}, \quad (12)$$

$$O_{\text{cls}} = \{c^k\}_{k \in [1, K]}, \quad (13)$$

where $\mathbf{p}_i^k = (x_i^k, y_i^k)$ is the predicted coordinate of this agent at time step i in the k^{th} mode, and c_k is the confidence.

As for the training objective, for each agent, we only backpropagate loss for the regression head with the lowest regression loss to avoid mode collapse and encourage diversity. For the classification head, we use the cross-entropy loss. The mode with the lowest regression loss is set as the positive sample, and the rest modes are set as negative samples consequently.

Mathematically, for the n^{th} agent, we find the index k_n of predicted mode with lowest regression loss:

$$k_n = \underset{k \in [1, K]}{\text{argmin}} \frac{1}{2T} \sum_{t=1}^T [d(x_{n,t}^k, x_{n,t}^*) + d(y_{n,t}^k, y_{n,t}^*)], \quad (14)$$

where $x_{n,t}^*, y_{n,t}^*$ is the ground truth coordinate at time step t of the n^{th} agent, and we adopt the widely used smooth L1 loss as [63]:

$$d(x_1, x_2) = \begin{cases} 0.5(x_1 - x_2)^2 & \text{if } \|x_1 - x_2\|_1 < 1, \\ \|x_1 - x_2\|_1 - 0.5 & \text{otherwise.} \end{cases} \quad (15)$$

The final training objective can be calculated as:

$$\begin{aligned} \mathcal{L} &= \lambda \mathcal{L}_{\text{cls}} + \mathcal{L}_{\text{reg}} \\ &= \lambda \frac{1}{N} \sum_{n=1}^N \text{CrossEntropy} \left(\{c^k\}_n, I(k_n) \right) \\ &\quad + \frac{1}{2NT} \sum_{n=1}^N \sum_{t=1}^T [d(x_{n,t}^{k_n}, x_{n,t}^*) + d(y_{n,t}^{k_n}, y_{n,t}^*)], \end{aligned} \quad (16)$$

where $I(k_n)$ is an one-hot vector of length K whose k_n^{th} entry is 1 and λ is the hyperparameter set to be 0.1 throughout the paper.

4 EXPERIMENTS

We depict the experimental results in details. Sec. 4.1 describes the datasets, metrics and the implementation details; Compared with the peer methods in Sec. 4.2, the numerical analysis is provided in Sec. 4.3 where state-of-the-art works are compared. Ablation studies are given in Sec. 4.4 with more detailed study and comparison of different Relative Encoding Strategies on Sec. 4.5 and Heterogeneous Graph Design in Sec. 4.6. Sec. 4.7 investigates the model capacity of HDGT. The qualitative visualization is presented subsequently in Sec. 4.8 to give the audience a glimpse of what the model learns in various complicated scenarios.

4.1 Protocols

4.1.1 Metrics

Similar to [3], [8], we use minADE, minFDE, and MR as the validation metric. Note that different benchmarks have slightly different definitions of the aforementioned metrics. Some other metrics such as minJointADE (minSADE), minJointFDE (minSFDE), SMR [20], [57], and mAP [17] are also used. Therefore, we use the following metrics in the validation set and report the leaderboard results in spirit of the official metrics respectively.

Formally, denote the multi-mode prediction of the model as $O = \{(\mathbf{p}_1^k, \mathbf{p}_2^k, \dots, \mathbf{p}_T^k)\}_{k \in [1, K]}$ where $\mathbf{p}_t^k = (x_t^k, y_t^k)$ is the predicted coordinate of the agent at time step t in the k^{th} mode. Denote the ground-truth trajectory of the agent as $O^* = (\mathbf{p}_1^*, \mathbf{p}_2^*, \dots, \mathbf{p}_T^*)$.

- 1) **minADE** (Minimum Average Displacement Error): the minimum value of the Euclidean distance between the prediction and ground truth averaged by the prediction length T , for K required predictions. It measures how good the predictions are in an average spirit under the Euclidean space. The detailed calculation is:

$$\text{minADE} = \min_{k=1, \dots, K} \frac{1}{T} \sum_{t=1}^T \|\mathbf{p}_t^k - \mathbf{p}_t^*\|_2 \quad (17)$$

- 2) **minFDE** (Minimum Final Displacement Error): similar to minADE, despite that it only calculates the error at the final time-step T . It focuses more on the accuracy of the predicted goal point, which emphasizes on the importance of predicting agents' intentions. The detailed calculation is:

$$\text{minFDE} = \min_{k=1, \dots, K} \|\mathbf{p}_T^k - \mathbf{p}_T^*\|_2 \quad (18)$$

- 3) **MR** (Missing Rate): the ratio of whether the Euclidean distance between the prediction and ground truth at the final time-step T for all K predictions is larger than 2m. Different from Euclid-based metrics like minADE and minFDE, it only requires the model to predict the coarse-grained intention. If any one of the predicted mode falls in the neighborhood of the ground-truth final point, it is counted as a hit (a 0 for the miss rate). The detailed calculation is:

$$\text{MR} = \begin{cases} 0, & \exists k \in \{1, \dots, K\}, \|\mathbf{p}_T^k - \mathbf{p}_T^*\|_2 \leq 2 \\ 1, & \text{Otherwise} \end{cases} \quad (19)$$

Note that all three aforementioned metrics is averaged over all target agents.

4.1.2 Datasets

We use two popular and large-scale datasets for evaluation.

- 1) **INTERACTION Dataset** It consists of various highly interactive driving situations, including highway ramps, roundabouts, and intersections, recorded using drones or fixed cameras worldwide [20]. We use its most recent 1.2 version which provides the official train/val/test split and the labels of test set are held out for the online competition. In this dataset, target agents are only vehicles, which has 40K tracks in total. The input data is 1 second history in 10 Hz and the required output is 3 seconds future in 10Hz. At most 6 modes are allowed. Note that there are four tracks in the INTERACTION challenges with different data and settings. We adopt their settings by: for the two conditional tracks which additionally give the future trajectory of the ego agent, we ignore the additional information for the consistency with other datasets and experiments. For the two joint tracks which require joint multi-agent prediction, we let all predictions with the same index as one modality.
- 2) **Waymo Open Motion Dataset** It provides a large scale dataset with annotations for objects with interacting behaviors over a wide range of road geometries [17] over a diverse set of locations, with 1750km unique roadways. It contains 7.64M agent tracks in total. We use its most recent 1.2 version with road connection information. It also provides the official train/val/test split and has an online leaderboard for test set. The target agents are divided into 3 types: vehicles, pedestrians, cyclists. The input is 1.1 seconds history in 10 Hz and required output is 8 seconds future in 2Hz. At most 6 modes are allowed. They provide up to 8 objects to predict in each scene, where their selection of target is biased to require objects that do not follow a constant velocity model or straight paths. We follow the official train/val/test split of the datasets above and report the test set results on the online leaderboards.

4.1.3 Implementation Details

The model is implemented by Pytorch and DGL. We use a hidden dimension of size 128. We use three Heterogeneous Driving Graph Transformer Layers on INTERACTION, similar to the number of GNN layers in [6], [64]. On Waymo Open Motion, we have six IHeterogeneous Driving Graph Transformer Layers (12.1 Million parameters) containing the same magnitude of parameters with [7] (15.3 Million parameters). Following the popular recipe for Transformer-based model, we use the AdamW optimizer with initial learning rate 5e-4, wight decay 1e-4, and batch size of 64. For both datasets, the number of training epochs is 30 with 1 epoch warmup and linearly decay to 0. The type-specific agent distance buffer hyper-parameter is: vehicle - 30 meters, pedestrian - 10 meters, cyclist - 20 meters, which is set empirically. For those agents with missed time-step (marked with empty), we interpolate the trajectory if it is in the observed horizon and mask its regression loss if it is in the future horizon. Different from the original Transformer [16], we do not set dropout since it slows down the converge speed with no explicit performance gain. As for the data augmentation, agent drop is adopted [7], which randomly drop 10% agents of a scene during training. For ResNet-like 1D CNN, similar to [8], we have three stages with each stage have two blocks. Different from [8], we do not adopt an additional FPN as we observe no performance gain. We use two layers of PointNet

TABLE 1
Results on the INTERACTION Leaderboard, *test set*.

Method	minADE↓	minFDE↓	MR↓
DenseTNT [5]	0.4342	0.7952	0.0596
MultiModalTransformer [6]	0.2130	0.5511	0.0511
GOHOME [12]	0.2005	0.5988	0.0491
HDGT (Ours)	0.1676	0.4776	0.0556

(a) Single Agent Track

Method	minSADE↓	minSFDE↓	SMR↓
ReCoG2 [64]	0.4668	1.1597	0.2377
DenseTNT [5]	0.4195	1.1288	0.2240
THOMAS [12]	0.4164	0.9679	0.1791
HDGT (Ours)	0.3030	0.9580	0.1938

(c) Multi-Agent Track

Method	minADE↓	minFDE↓	MR↓
DenseTNT [5]	0.5452	0.6375	0.0348
GOHOME [12]	0.2020	0.5902	0.0284
MultiModalTransformer [6]	0.1989	0.5799	0.0429
HDGT (Ours)	0.1553	0.4577	0.0276

(b) Conditional Single Agent Track

Method	minSADE↓	minSFDE↓	SMR↓
ReCoG2 [64]	0.3295	0.8693	0.1498
THOMAS [12]	0.3148	0.7162	0.1067
DenseTNT [5]	0.2786	0.8916	0.1502
HDGT (Ours)	0.2255	0.7875	0.1322

(d) Conditional Multi-Agent Track

TABLE 2
Results on Waymo Motion Leaderboard, *test set*.

	minADE↓	minFDE↓	MR↓	mAP↑
SimpleCNNOnRaster [65]	0.7400	1.4936	0.2091	0.2136
ReCoAt [66]	0.7703	1.6668	0.2437	0.2711
DenseTNT [5]	1.0387	1.5514	0.1573	0.3281
Kraken-NMS (Yandex SDG)	0.6732	1.3947	0.1850	0.3646
SceneTransformer [7]	0.6117	1.2116	0.1564	0.2788
*MultiPath++ [18]	<i>0.5557</i>	<i>1.1577</i>	<i>0.1340</i>	<i>0.4092</i>
HDGT (ours)	0.5933	1.2055	0.1511	0.2854

for polylines and polygons, similar to [3]. All experiments are conducted on NVIDIA Tesla V100.

4.2 Compared Baselines

We briefly introduce the methods we compare with on the public leaderboard of INTERACTION and Waymo Motion.

- **DenseTNT** [5] adopts VectorNet [3] as the encoder which models the scene as a homogeneous fully connected graph and proposes an optimization-based technique to select the best goal point from a predicted dense candidate set. Compared to direct decode trajectory, their approach could generate more diverse modes and thus have better MR.
- **MultiModalTransformer** [6] builds two fully connected graph (AgentAgent and AgentMap) and adopts the Transformer to encode the two graphs. The decoder is classification plus regression MLP (same with ours).
- **GOHOME** [12] encodes the scene with a LaneGCN [8]-like, which builds a series of four graphs (LaneToAgent, LaneToLane, LaneToActor, and finally ActorToActor) and extracts information with GCN [21]. Based on the graphs, they generate heatmap outputs representing the future position probability distribution of agents. Similar to DenseTNT, they could generate more diverse multi-modal predictions by an additional optimization step.
- **ReCoG2** [64] builds a heterogeneous graph contains two types of nodes: vehicles and lanes, which is first encoded by RNN and CNN respectively. Then, they adopt GCN [21] for the graph and uses a RNN as decoder. Compared to HDGT, their map is encoded in the rasterized way and the GCN they used did not consider the heterogeneity of nodes.
- **THOMAS** [12]: extends GOHOME in a hierarchical way of heatmap generation for fine-grained prediction and an attention-based decoder for consistent multi-agent prediction.
- **SimpleCNNOnRaster** [65] and **ReCoAt** [66] adopt CNN on rasterized images to encode map while [65] adopts rasteriza-

tion for agents as well and [66] adopts RNN for agents. Their decoders are both classification plus regression MLP.

- **SceneTransformer** [7] builds a fully connected graph for all agents and map elements and they adopt Transformer on both spatial and temporal axis in an factorized way.
- **Multipath++** [18] is a concurrent work which extract agents' surrounding elements in their relative reference in encoder and proposes an ensemble method for decoding.

4.3 Overall Performance Evaluation

Throughout 11/3/2021 (as of date for our final submission to INTERACTION public evaluation) to 4/21/2022 (as of the date for our initial submission of this paper), in terms of minADE/minFDE metric, we rank the **best** result in all four tracks on the INTERACTION Leaderboard³ and **second** on the Waymo Open Motion Leaderboard⁴. Table 1 and Table 2 report the detailed results and comparison to previous state-of-the-arts.

4.3.1 Results on INTERACTION

Table 1 shows that HDGT outperforms other candidates by a large margin in terms of minADE/minFDE. Specifically, compared to MultiModalTransformer [6] and ReCoG2 [64] which has two separate graphs and fully connected graph, HDGT outperforms them in all metrics, which demonstrates the importance of modelling symmetric relativity and heterogeneity as in HDGT. As for the MR metric, we notice that DenseTNT and GOHOME has better ranks compared to their minADE/minFDE rank, which shows the effectiveness of their optimization-based decoding THOMAS has the best MR in the multi-agent settings, which comes from their consistent multi-agent trajectory decoder. Note that HDGT's output head is just the simple regression plus classification MLP scheme, which still has competitive MR results and thus demonstrates the strong representation vector from the encoder.

4.3.2 Results on Waymo Open Motion

In Table 2, * denotes that the concurrent MultiPath++ involves ensemble with results shown in *italic*; our HDGT ranks the first in the single model in terms of minADE/minFDE/MR as shown in **bold**.

We can observe that models with rasterized images for map encoding including SimpleCNNOnRaster [65] and ReCoAt [66] have relative large performance gap compared to other methods, which verify the advantages of vector-based methods over

3. <http://challenge.interaction-dataset.com/prediction-challenge/intro>

4. <https://waymo.com/open/challenges/2021/motion-prediction/>

TABLE 3
Results of Model Component Ablation Study.

Model variants	Relative	HD-Map	Semantic	Heterogeneous	minADE ↓	minFDE ↓	MR ↓
1					0.2287	0.5338	0.0153
2	✓				0.1889	0.4572	0.0093
3	✓	✓			0.1526	0.3726	0.0044
4	✓	✓	✓		0.1328	0.3379	0.0038
5	✓	✓	✓	✓	0.1071	0.2945	0.0023

raster-based methods. Again, DenseTNT [5] demonstrates strong performance on predicting agents’ intentions, i.e., high MR and mAP while it has low scores on distance-based metrics - minADE/minFDE, which is a widely known property of goal-based methods. As for the SceneTransformer which adopts Transformer for scene encoding as well, HDGT outperforms it with lower cost. From the model complexity perspective, SceneTransformer has 15.3 million parameters while the submitted version of HDGT has 12.1 million parameters. From the computation efficiency perspective, SceneTransformer conducts attention mechanism on both spatial and temporal axes with two factorized fully connected graphs while HDGT only operates on one single heterogeneous graph sparsely connected according to semantic relations. As for the concurrent work Multipath++, we observe its strong performance, which show the effectiveness of their ensemble technique.

In summary, HDGT, as a backbone for driving scene encoding, could outperform state-of-the-art backbones on L2-norm based metrics with simple output heads. For other types of metrics like MR/mAP, more advanced decoder as well as techniques like ensemble and test-time-augmentation could be explored.

4.4 Ablation Study

We conduct the ablation study to validate the effectiveness of each design in HDGT and the results are shown in Table 3. The experiments are conducted on the validation set of INTERACTION dataset and we keep all the other settings the same. Note that all models have enough capacities to overfit the training set and we report their best performance on the validation set. In Table 3, *Relative* means using local coordinate system for each node and its in-edges, otherwise we set the autonomous vehicle as the global reference. *HD-Map* means using HD-Map otherwise only agents’ information. *Semantic* means building graph according to distance threshold and lane connection information otherwise building a fully-connected graph. *Heterogeneous* means using different parameters for different node and edge types otherwise sharing the parameters. From the results, we can conclude that:

- *Relative representation in the allocentric view is better.* In model 1, we set the autonomous vehicle as the global reference while in model 2 we adopt the proposed node-centric local coordinate system. The improvement of the model 2 shows the necessity of learning relative spatial relations instead of absolute coordinates, which aligns well with the conclusions in [2], [18], [67].
- *Road semantics from HD-map are important.* In model 3, we add HD-Map and construct a homogeneous fully-connected graph similar to VectorNet [3]. It has better performance compared to model 2, which proves the importance of HD-map information for the trajectory prediction. For example, drivers tend to drive following the lane centerline. When it comes to parking, they would stop nearby the curb.

- *Semantic relations are helpful.* In model 4, instead of a fully-connected graph, we build the semantic graph as proposed in HDGT but the GNN parameters for different types of nodes are still shared. As a result, it brings improvement compared to model 3 while worse than model 5.
- *Heterogeneous graph network parameters are vital.* Model 5 is the full version of HDGT and we can find that specific parameters for different types of nodes and edges could boost the performance significantly compared to model 4.

4.5 Comparison of Relative Encoding Strategy

In Sec. 3.3, we discuss the importance to encode the agent in the relative reference. Here, we compare three strategies to model the relativity among agents:

- 1) **Fixed Reference:** it select one agent (usually the ego agent) as the global reference, i.e, adopting its coordinate and heading at the last observed time-step as the origin and positive x-axis of the global coordinate system. It is widely used in VectorNet [3] series of work, LaneGCN [8] series of work, and SceneTransformer [7]. This strategy could achieve decent performance for the ego agent. However, since other agents are not in their local reference system, people have to forward model *number of agents* times with each agent as the global reference to achieve the best performance for multi-agent prediction [17].
- 2) **Relative Edge:** One modification of Fixed reference is to have each agent in their local coordinate system with in-edges also transformed into its target nodes’ coordinate system. It is adopted in methods including concurrent work Multipath++ [18] and HEAT-I-R [2]. One potential downside of this strategy is that for those in-edges with the same source node, their features are encoded separately which ignores the fact that they are indeed one object.
- 3) **Pose Change:** to this end, we encode the relative feature by combing the source nodes’ feature with Δ Pose. The concurrent work HiVT also designs a Global Interaction Module which encodes the pair-wise coordinate transformation.

We conduct experiments with the above three strategies on the validation set of INTERACTION dataset, while keeping other designs of HDGT unchanged. The results are in Table 4 and we can find that fixing reference has explicit performance compared to encoding each agent in their own coordinate system. When it comes to the Relative Edge and Pose Change, we can observe that Pose Change demonstrates better performance, verifying the effectiveness of the proposed Relative Relation Encoding technique.

4.6 Analysis of Heterogeneous Graph Design

In this section, we examine the following design choice:

TABLE 4
Results of Different Relative Encoding Strategies. Averaged over 3 runs.

Relative Encoding	minADE ↓	minFDE ↓	MR ↓
Fixed Reference	0.1549±0.0032	0.4079±0.0095	0.0051±0.0000
Relative Edge	0.1158±0.0021	0.3296±0.0060	0.0025±0.0000
Pose Change (Ours)	0.1081±0.0030	0.2945±0.0061	0.0023±0.0000

TABLE 5
Performance by Different Numbers of Heterogeneous Driving Graph Transformer Layers.

Dataset Layer #	INTERACTION				Waymo			
	minADE ↓	minFDE ↓	MR ↓	Inference Time (second)	minADE ↓	minFDE ↓	MR ↓	Inference Time (second)
1	0.1430	0.3997	0.0043	95	0.6462	1.3099	0.1865	901
2	0.1149	0.3243	0.0025	100	0.6013	1.2004	0.1648	977
3	0.1071	0.2945	0.0023	107	0.5835	1.1736	0.1601	1038
6	0.1082	0.3035	0.0022	122	0.5673	1.1507	0.1505	1146

TABLE 6
Results of Heterogeneous Graph Design.

Model	minADE ↓	minFDE ↓	MR ↓
HDGT (Ours)	0.1071	0.2945	0.0023
Merge Lane Connectivity	0.1622	0.4003	0.0078
Homogeneous Map Node	0.1272	0.3171	0.0031
GCN Function	0.1882	0.4721	0.0114

- **Merge Lane Connectivity:** we merge the four types of lane connectivity edges into one to examine the necessity of the modeling the lane topology with separate parameters.
- **Homogeneous Map Node:** we merge the node type of all lanes and traffic signs while keep all the edges unchanged to examine the necessity to process lanes and traffic signs with different parameters.
- **GCN Function:** we adopt the GCN [21] like aggregation and update function as in LaneGCN [8] to compare with Transformer in HDGT. The in-edges’ features of each node are fed into an MLP and then an summation operation is conducted to aggregate information. For update function of nodes, we concatenate the aggregated information with the old node features and feed them into an MLP. For different node and edge types, we still adopt different parameters.

We compare the three aforementioned design choice with the proposed one on the validation set of INTERACTION dataset. The results are in Table 6 and we can conclude that: (i) Change the Transformer to GCN-like function leads to the largest performance drop, which demonstrates the advantages of adopting Transformer structure. (ii) Ignoring the lane topology leads to the second largest performance drop. The result is in accordance with expectation since vehicles tends to move along the lane and lane topology could provide useful information to reduce the uncertainty of the prediction, which aligns with the conclusions in existing works [8], [10]. (iii) Merge the lane and traffic signs leads to the least performance drop. It might come from the fact during the initialization of the feature, the sub-type information has been encoded by learnable embeddings, which could be used to distinguish lanes and traffic signs. However, compared to the original design, it still leads to performance drop which might be because of the different characteristics of the two types of objects.

4.7 Exploration on Model Capacity

We examine the scalability of HDGT. We conducted experiments on the validation set of INTERACTION and Waymo Open Motion by models with different number of Heterogeneous Driving Graph Transformer Layers. Additionally, we report the inference time on the each datasets’ validation set for 1 epoch to estimate

the influence of different number of layers. For Waymo, the minADE/minFDE are calculated in their official way. From Table 5, we can observe that for INTERACTION (40K tracks), the performance saturates with around 3 Heterogeneous Driving Graph Transformer Layers. When it comes to the larger Waymo Open Motion (7.64M tracks), the performance could be still improved even with 6 layers. This shows that HDGT enjoys the strong scalability of Transformer structure, which brings huge success in the vision [68] and language [31] fields. The simple and yet unified structure of HDGT makes it easy to extend its capacity when there is more data, which fits perfectly for the autonomous driving industry.

4.8 Visualization

We visualize the prediction results of HDGT in some scenarios on the Waymo Open Motion dataset in Fig. 4. One can find that to predict agents’ future trajectories accurately, information about diverse elements and their relations are necessary such as agents, lanes, stop-signs, traffic lights and *etc.* Thus, the explicit modeling of the heterogeneous nature of the driving scene in HDGT is significant and beneficial for the trajectory prediction.

5 CONCLUSION

We have proposed a principled approach to extract better representation for trajectory prediction. It models the driving scene as a heterogeneous graph and adapt transformer as the graph aggregation and update function.

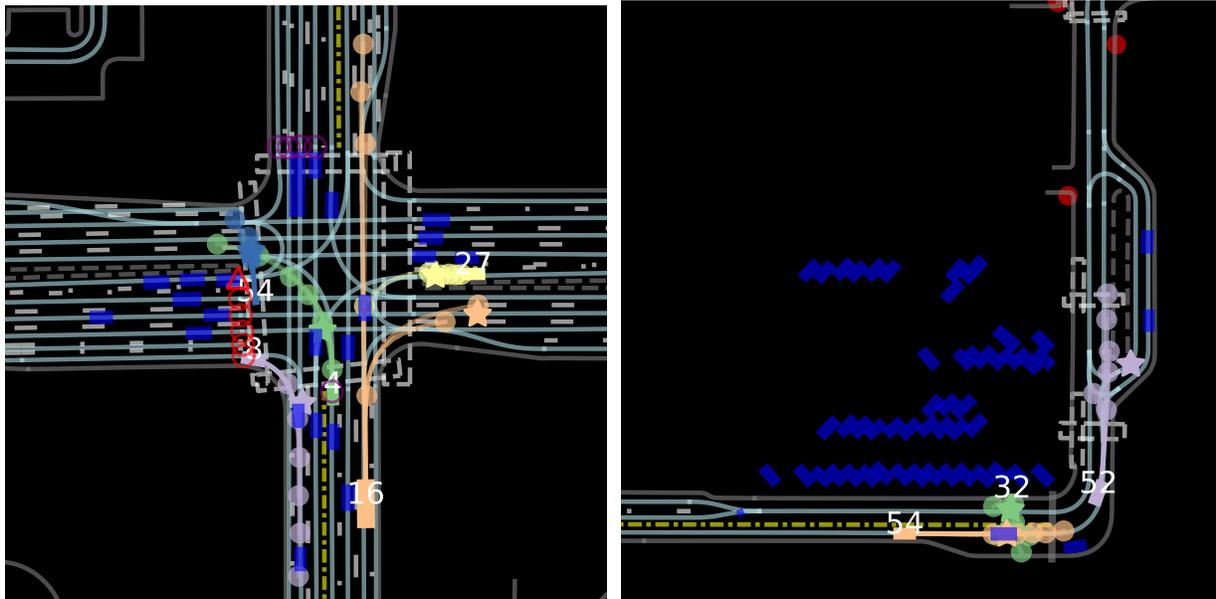
The spatial features are normalized into the local coordinate system when being aggregated. The proposed approach has achieved state-of-the-art results on two recent large-scale and competitive benchmarks. Thorough ablation studies validate the effectiveness of each module in our method.

ACKNOWLEDGEMENT

This work was partly supported by NSFC (62222607,62206172), Shanghai Municipal Science and Technology Major Project (2021SHZDZX0102). The authors are thankful to the anonymous reviewers for their valuable comments.

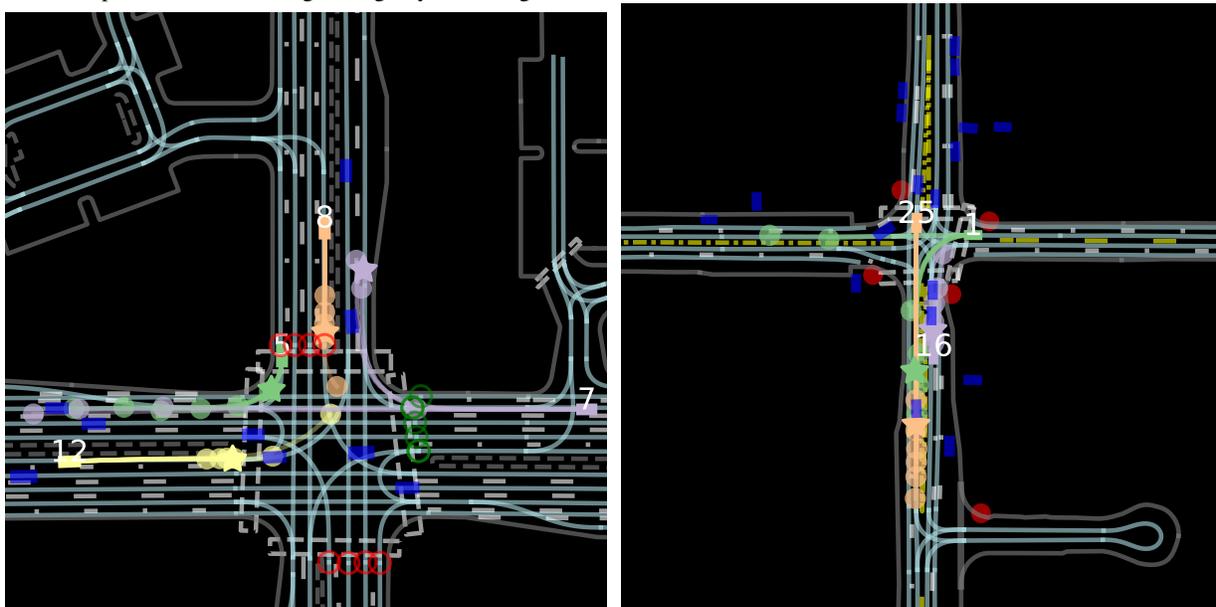
REFERENCES

- [1] T. Salzmann, B. Ivanovic, P. Chakravarty, and M. Pavone, “Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data,” p. 683–700, 2020.
- [2] X. Mo, Z. Huang, Y. Xing, and C. Lv, “Multi-agent trajectory prediction with heterogeneous edge-enhanced graph attention network,” vol. 23, no. 7, 2022, pp. 9554–9567.



(a) **P54** was passing the road through the crosswalk while **V4** desired to turn left. The model outputs their different potential futures. For **V16**, as its target is unclear by observation, the model outputs diverse futures: go straight, yield, or right turn.

(b) Since **P32** could either walk along curb or pass the road, the model outputs the proper reactions of **V54** under both conditions. For **V52**, its possibility for going at different lanes are all considered.



(c) **P7** could either go straight through the green light or turn right. Thus, the prediction of **P5** is influenced. For **V8** and **V12**, they were far away from the traffic light and thus their predictions include both passing and yielding.

(d) **V1**'s future had high uncertainty (going straight or left turn) and it was also influenced by either **V25** or **V16** under different situations. Also, the existing of stop sign (red double ring) results in conservative predictions for **V16**.

Fig. 4: **Visualization of Prediction Results on Waymo Opem Motion Dataset.** \star represents the ground-truth final position and \circ denotes predicted final positions of modes (lower confidence, more transparent). Abbreviation: V-vehicle, P-pedestrian, and C-cyclist.

- [3] J. Gao, C. Sun, H. Zhao, Y. Shen, D. Anguelov, C. Li, and C. Schmid, "Vectornet: Encoding hd maps and agent dynamics from vectorized representation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 11 525–11 533.
- [4] H. Zhao, J. Gao, T. Lan, C. Sun, B. Sapp, B. Varadarajan, Y. Shen, Y. Shen, Y. Chai, C. Schmid *et al.*, "Tnt: Target-driven trajectory prediction," in *Conference on Robot Learning*, 2021, pp. 895–904.
- [5] J. Gu, C. Sun, and H. Zhao, "Densentn: End-to-end trajectory prediction from dense goal sets," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 15 303–15 312.
- [6] Z. Huang, X. Mo, and C. Lv, "Multi-modal motion prediction with transformer-based neural network for autonomous driving," *arXiv preprint arXiv:2109.06446*.
- [7] J. Ngiam, B. Caine, V. Vasudevan, Z. Zhang, H.-T. L. Chiang, J. Ling, R. Roelofs, A. Bewley, C. Liu, A. Venugopal *et al.*, "Scene transformer: A unified multi-task model for behavior prediction and planning," in *International Conference on Learning Representations*, 2022, pp. 1–15.
- [8] M. Liang, B. Yang, R. Hu, Y. Chen, R. Liao, S. Feng, and R. Urtasun, "Learning lane graph representations for motion forecasting," in *European Conference on Computer Vision*. Springer, 2020, pp. 541–556.
- [9] W. Zeng, M. Liang, R. Liao, and R. Urtasun, "Lanercnn: Distributed representations for graph-centric motion forecasting," in *2021 IEEE/RSJ*

- International Conference on Intelligent Robots and Systems*. IEEE, 2021, pp. 532–539.
- [10] N. Deo, E. Wolff, and O. Beijbom, “Multimodal trajectory prediction conditioned on lane-graph traversals,” in *Conference on Robot Learning*, 2022, pp. 203–212.
- [11] T. Gilles, S. Sabatini, D. Tsishkou, B. Stanculescu, and F. Moutarde, “Thomas: Trajectory heatmap output with learned multi-agent sampling,” in *International Conference on Learning Representations*, 2022, pp. 1–15.
- [12] T. Gilles, S. Sabatini, D. V. Tsishkou, B. Stanculescu, and F. Moutarde, “Gohome: Graph-oriented heatmap output for future motion estimation,” *2022 International Conference on Robotics and Automation*, pp. 9107–9114, 2021.
- [13] J. D. M.-W. C. Kenton and L. K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of NAACL-HLT*, 2019, pp. 4171–4186.
- [14] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” in *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M. Balcan, and H. Lin, Eds., vol. 33. Curran Associates, Inc., 2020, pp. 1877–1901.
- [15] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” *International Conference on Learning Representations*, pp. 3383–3395, 2021.
- [16] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” vol. 30, 2017.
- [17] S. Ettinger, S. Cheng, B. Caine, C. Liu, H. Zhao, S. Pradhan, Y. Chai, B. Sapp, C. R. Qi, Y. Zhou *et al.*, “Large scale interactive motion forecasting for autonomous driving: The waymo open motion dataset,” pp. 9710–9719, 2021.
- [18] B. Varadarajan, A. Hefny, A. Srivastava, K. S. Refaat, N. Nayakanti, A. Cornman, K. Chen, B. Douillard, C. P. Lam, D. Anguelov, and B. Sapp, “Multipath++: Efficient information fusion and trajectory aggregation for behavior prediction,” in *International Conference on Robotics and Automation*. IEEE Press, 2022, p. 7814–7821.
- [19] Z. Zhou, L. Ye, J. Wang, K. Wu, and K. Lu, “Hivt: Hierarchical vector transformer for multi-agent motion prediction,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 8823–8833.
- [20] W. Zhan, L. Sun, D. Wang, H. Shi, A. Clause, M. Naumann, J. Kummerle, H. Konigshof, C. Stiller, A. de La Fortelle *et al.*, “Interaction dataset: An international, adversarial and cooperative motion dataset in interactive driving scenarios with semantic maps,” *arXiv preprint arXiv:1910.03088*, 2019.
- [21] T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” in *International Conference on Learning Representations*, 2017.
- [22] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Liò, and Y. Bengio, “Graph attention networks,” in *International Conference on Learning Representations*, 2018, pp. 3383–3395.
- [23] M. Schlichtkrull, T. N. Kipf, P. Bloem, R. v. d. Berg, I. Titov, and M. Welling, “Modeling relational data with graph convolutional networks,” in *European semantic web conference*. Springer, 2018, pp. 593–607.
- [24] X. Wang, H. Ji, C. Shi, B. Wang, Y. Ye, P. Cui, and P. S. Yu, “Heterogeneous graph attention network,” in *The world wide web conference*, 2019, pp. 2022–2032.
- [25] H. Hong, H. Guo, Y. Lin, X. Yang, Z. Li, and J. Ye, “An attention-based graph neural network for heterogeneous structural learning,” *Conference on Artificial Intelligence*, pp. 4132–4139, 2020.
- [26] C. Zhang, D. Song, C. Huang, A. Swami, and N. V. Chawla, “Heterogeneous graph neural network,” in *Proceedings of ACM SIGKDD international conference on knowledge discovery & data mining*, 2019, pp. 793–803.
- [27] Z. Hu, Y. Dong, K. Wang, and Y. Sun, “Heterogeneous graph transformer,” in *Proceedings of The Web Conference 2020*, 2020, pp. 2704–2710.
- [28] H. Caesar, V. Bankiti, A. H. Lang, S. Vora, V. E. Liong, Q. Xu, A. Krishnan, Y. Pan, G. Baldan, and O. Beijbom, “nusenes: 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.
- [29] M.-F. Chang, J. Lambert, P. Sangkloy, J. Singh, S. Bak, A. Hartnett, D. Wang, P. Carr, S. Lucey, D. Ramanan *et al.*, “Argoverse: 3d tracking and forecasting with rich maps,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8748–8757.
- [30] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, I. Sutskever *et al.*, “Language models are unsupervised multitask learners,” *OpenAI blog*, vol. 1, no. 8, p. 9, 2019.
- [31] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell *et al.*, “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [32] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*, 2021, pp. 1–15.
- [33] A. Baevski, Y. Zhou, A. Mohamed, and M. Auli, “wav2vec 2.0: A framework for self-supervised learning of speech representations,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 12 449–12 460, 2020.
- [34] W. Kim, B. Son, and I. Kim, “Vilt: Vision-and-language transformer without convolution or region supervision,” in *International Conference on Machine Learning*, 2021, pp. 5583–5594.
- [35] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch, “Decision transformer: Reinforcement learning via sequence modeling,” *Advances in neural information processing systems*, vol. 34, pp. 15 084–15 097, 2021.
- [36] X. Jia, L. Sun, M. Tomizuka, and W. Zhan, “Ide-net: Interactive driving event and pattern extraction from human data,” *IEEE Robotics and Automation Letters*, vol. 6, 2021.
- [37] D. Helbing and P. Molnar, “Social force model for pedestrian dynamics,” *Physical review E*, vol. 51, no. 5, p. 4282, 1995.
- [38] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 961–971.
- [39] 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/CVF Conference on Computer Vision and Pattern Recognition*, 2018, pp. 2255–2264.
- [40] B. Ivanovic, E. Schmerling, K. Leung, and M. Pavone, “Generative modeling of multimodal multi-human behavior,” in *2018 IEEE/RSSJ International Conference on Intelligent Robots and Systems*. IEEE, 2018, pp. 3088–3095.
- [41] B. Ivanovic and M. Pavone, “The trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs,” in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2019, pp. 2375–2384.
- [42] A. Mohamed, K. Qian, M. Elhoseiny, and C. Claudel, “Social-stgcnn: A social spatio-temporal graph convolutional neural network for human trajectory prediction,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 14 424–14 432.
- [43] C. Tang and R. R. Salakhutdinov, “Multiple futures prediction,” in *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [44] N. Lee, W. Choi, P. Vernaza, C. B. Choy, P. H. Torr, and M. Chandraker, “Desire: Distant future prediction in dynamic scenes with interacting agents,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2017, pp. 336–345.
- [45] T. Zhao, Y. Xu, M. Monfort, W. Choi, C. Baker, Y. Zhao, Y. Wang, and Y. N. Wu, “Multi-agent tensor fusion for contextual trajectory prediction,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12 126–12 134.
- [46] Y. Chai, B. Sapp, M. Bansal, and D. Anguelov, “Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction,” in *Conference on Robot Learning*, 2020, pp. 86–99.
- [47] H. Cui, V. Radosavljevic, F.-C. Chou, T.-H. Lin, T. Nguyen, T.-K. Huang, J. Schneider, and N. Djuric, “Multimodal trajectory predictions for autonomous driving using deep convolutional networks,” in *2019 International Conference on Robotics and Automation*. IEEE, 2019, pp. 2090–2096.
- [48] J. Hong, B. Sapp, and J. Philbin, “Rules of the road: Predicting driving behavior with a convolutional model of semantic interactions,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8454–8462.
- [49] F. Marchetti, F. Becattini, L. Seidenari, and A. D. Bimbo, “Mantra: Memory augmented networks for multiple trajectory prediction,” in

Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 7143–7152.

- [50] C. Choi, J. H. Choi, J. Li, and S. Malla, “Shared cross-modal trajectory prediction for autonomous driving,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 244–253.
- [51] M. Bansal, A. Krizhevsky, and A. Ogale, “Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst,” in *Robotics: Science and Systems*, 2019.
- [52] W. Luo, Y. Li, R. Urtaşun, and R. Zemel, “Understanding the effective receptive field in deep convolutional neural networks,” *Advances in neural information processing systems*, vol. 29, 2016.
- [53] S. Hochreiter and J. Schmidhuber, “Long short-term memory,” *Neural Computation*, vol. 9, 1997.
- [54] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, “Pointnet: Deep learning on point sets for 3d classification and segmentation,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 652–660.
- [55] M. Ye, T. Cao, and Q. Chen, “Tpcn: Temporal point cloud networks for motion forecasting,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2021, pp. 11 318–11 327.
- [56] Y. Biktairov, M. Stebelev, I. Rudenko, O. Shliakhko, and B. Yangel, “Prank: motion prediction based on ranking,” vol. 33, 2020, pp. 2553–2563.
- [57] S. Casas, C. Gulino, S. Suo, K. Luo, R. Liao, and R. Urtaşun, “Implicit latent variable model for scene-consistent motion forecasting,” in *European Conference on Computer Vision*. Springer, 2020, pp. 624–641.
- [58] M. Ye, J. Xu, X. Xu, T. Cao, and Q. Chen, “Dcms: Motion forecasting with dual consistency and multi-pseudo-target supervision,” *arXiv preprint arXiv:2204.05859*, 2022.
- [59] X. Jia, L. Chen, P. Wu, J. Zeng, J. Yan, H. Li, and Y. Qiao, “Towards capturing the temporal dynamics for trajectory prediction: a coarse-to-fine approach,” in *Conference on Robot Learning*, ser. Proceedings of Machine Learning Research, K. Liu, D. Kulic, and J. Ichnowski, Eds., vol. 205. PMLR, 14–18 Dec 2023, pp. 910–920.
- [60] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016, pp. 770–778.
- [61] J. L. Ba, J. R. Kiros, and G. E. Hinton, “Layer normalization,” *NeurIPS Deep Learning Symposium*, 2016.
- [62] P. W. Battaglia, J. B. Hamrick, V. Bapst, A. Sanchez-Gonzalez, V. Zambaldi, M. Malinowski, A. Tacchetti, D. Raposo, A. Santoro, R. Faulkner et al., “Relational inductive biases, deep learning, and graph networks,” *arXiv preprint arXiv:1806.01261*, 2018.
- [63] R. Girshick, “Fast r-cnn,” in *Proceedings of the IEEE international conference on computer vision*, 2015, pp. 1440–1448.
- [64] X. Mo, Y. Xing, and C. Lv, “Recog: A deep learning framework with heterogeneous graph for interaction-aware trajectory prediction,” *arXiv preprint arXiv:2012.05032*, 2020.
- [65] S. Konev, K. Brodt, and A. Sanakoyeu, “Motioncnn: A strong baseline for motion prediction in autonomous driving,” *Workshop on Autonomous Driving, Computer Vision and Pattern Recognition*, 2021.
- [66] Z. Huang, X. Mo, and C. Lv, “Recoat: A deep learning framework with attention mechanism for multi-modal motion prediction,” *Workshop on Autonomous Driving, Computer Vision and Pattern Recognition*, 2021.
- [67] X. Jia, L. Sun, H. Zhao, M. Tomizuka, and W. Zhan, “Multi-agent trajectory prediction by combining egocentric and allocentric views,” in *Conference on Robot Learning*, 2022, pp. 1434–1443.
- [68] X. Zhai, A. Kolesnikov, N. Houlsby, and L. Beyer, “Scaling vision transformers,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 12 104–12 113.



Xiaosong Jia is currently a PhD student at Department of Computer Science and Engineering, Shanghai Jiao Tong University (SJTU), Shanghai. Before that, he earned B.E. in IEEE Honor class at SJTU. His research interests include autonomous driving and machine learning, with (co-) first-authored papers published in CoRL, NeurIPS, CVPR, RAL, TKDE, etc.



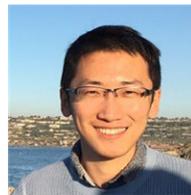
Penghao Wu was once an undergraduate student with the University of Michigan-Shanghai Jiao Tong University Joint Institute (UM-SJTU JI), Shanghai Jiao Tong University (SJTU), Shanghai where he obtained his B.E. in electrical and computer engineering. He is now a master student with department of computer science, UCSD, USA. His research interests include autonomous driving and computer vision. He has published first/co-authored paper in ECCV, NeurIPS, ICLR on autonomous driving.



Li Chen received the B.E. in mechanical engineering from Shanghai Jiao Tong University in 2019, and the M.S. in Robotics from University of Michigan, Ann Arbor, USA in 2020. He was once a research engineer with SenseTime after graduation, and is currently a researcher with Shanghai AI Laboratory, Shanghai, China. His research interests include autonomous driving, computer vision and machine learning. He has published first/co-authored paper in ECCV, CoRL, NeurIPS on autonomous driving.



Yu Liu is a Principal Investigator in Shanghai AI Lab and also the General Manager of ADG Business Unit and a Senior Director of research in SenseTime Group, Hong Kong. Before that, he obtained the PhD from Multimedia Lab (MMLab), the Chinese University of Hong Kong, and B.E. from Beihang University, Beijing, China in 2016. His research interests include Artificial General Intelligence, Deep RL and neural network understanding with big data.



Hongyang Li received PhD in computer science from Chinese University of Hong Kong, in 2020, and the B.E. degree in Computer Science from Dalian university of technology in 2014. He is currently a Research Scientist and leading the OpenDriveLab team at Shanghai AI Lab. His expertise focuses on perception and cognition, end-to-end autonomous driving and foundation model. He is also affiliated with Department of Computer Science and Engineering, Shanghai Jiao Tong University as a post-doc researcher.



Junchi Yan (S'10-M'11-SM'21) is currently an Associate Professor with Department of Computer Science and Engineering, Shanghai Jiao Tong University, Shanghai, China. Before that, he was a Senior Research Staff Member with IBM Research where he started his career since April 2011. His research interests include machine learning and computer vision. He regularly serves as Senior PC/Area Chair for NeurIPS, ICML, CVPR, AAAI, IJCAI, ACM-MM and Associate Editor for the Pattern Recognition Journal.