

Post-interactive Multimodal Trajectory Prediction for Autonomous Driving

Ziyi Huang^a, Yang Li^{b,*}, Dushuai Li^b, Yao Mu^c, Hongmao Qin^b and Nan Zheng^{d,*}

^aCollege of Computer Science and Electronic Engineering, Hunan University, Changsha, 410082, China

^bCollege of Mechanical and Vehicle Engineering, Hunan University, Changsha, 410082, China

^cDepartment of Computer Science, The University of Hong Kong, Hong Kong, 999077, China

^dDepartment of Civil Engineering, Monash University, Wellington Rd, Clayton VIC, 3800, Australia

ARTICLE INFO

Keywords:

Autonomous driving
Multimodal trajectory prediction
Post-interaction
Hypergraph neural network
Transformer

ABSTRACT

Modeling the interactions among agents for trajectory prediction of autonomous driving has been challenging due to the inherent uncertainty in agents' behavior. The interactions involved in the predicted trajectories of agents, also called post-interactions, have rarely been considered in trajectory prediction models. To this end, we propose a coarse-to-fine Transformer for multimodal trajectory prediction, *i.e.*, Pioformer, which explicitly extracts the post-interaction features to enhance the prediction accuracy. Specifically, we first build a Coarse Trajectory Network to generate coarse trajectories based on the observed trajectories and lane segments, in which the low-order interaction features are extracted with the graph neural networks. Next, we build a hypergraph neural network-based Trajectory Proposal Network to generate trajectory proposals, where the high-order interaction features are learned by the hypergraphs. Finally, the trajectory proposals are sent to the Proposal Refinement Network for further refinement. The observed trajectories and trajectory proposals are concatenated together as the inputs of the Proposal Refinement Network, in which the post-interaction features are learned by combining the previous interaction features and trajectory consistency features. Moreover, we propose a three-stage training scheme to facilitate the learning process. Extensive experiments on the Argoverse 1 dataset demonstrate the superiority of our method. Compared with the baseline HiVT-64, our model has reduced the prediction errors by 4.4%, 8.4%, 14.4%, 5.7% regarding metrics minADE₆, minFDE₆, MR₆, and brier-minFDE₆, respectively.

1. Introduction

Trajectory prediction for autonomous vehicles (AVs) aims to forecast the future paths or motions of the traffic participants based on past trajectories and environmental clues such as lanes and obstacles (Huang et al., 2022; Li et al., 2020), which is essential to the safe and efficient decision and planning of AVs, especially in dynamic urban conditions (Li et al., 2023; Huang et al., 2023b). With accurate trajectory predictions, the vehicle can make proactive decisions that can maximize traffic efficiency while satisfying safety constraints. Trajectory prediction methods can be divided into deterministic prediction and multimodal prediction (Huang et al., 2023a), in which the former provides only one prediction for each agent, and the latter can generate multiple predictions each time. To handle the uncertainty and the multimodality involved in driving behavior, this study focuses on multimodal trajectory prediction that can generate multiple traffic rule-compliant predictions based on observed data and context information.

Interaction modeling has been crucial for trajectory predictions of multiple agents in strong-interactive traffic scenes (Ma et al., 2019; Huang et al., 2023a). Social Force is a conventional method that uses attractive or repulsive forces to model the interactions among agents, while it cannot handle complex interactions (Helbing and Molnar, 1995). Recently, learning-based methods have been widely used for trajectory predictions, which use Recurrent Neural Networks (RNN) (Alahi et al., 2016; Chandra et al., 2019), Graph Convolutional Networks (GCNs) (Jiang et al., 2019; Xu et al., 2022a; Jia et al., 2023) or Graph Attention Networks (GATs) (Velickovic et al., 2017), to learn the spatial and

*This work was supported by the National Natural Science Foundation of China with 52302493 and 52372411, Hunan Provincial Natural Science Foundation of China with 2023JJ10008, and State Key Laboratory of Advanced Design and Manufacturing Technology for Vehicle with 72275004.

*Corresponding author

✉ hzy03@hnu.edu.cn (Z. Huang); lyxc56@gmail.com (Y. Li); lidushuai@163.com (D. Li); muyao@connect.hku.hk (Y. Mu); qinhongmao@hnu.edu.cn (H. Qin); Nan.Zheng@monash.edu (N. Zheng)

ORCID(s):

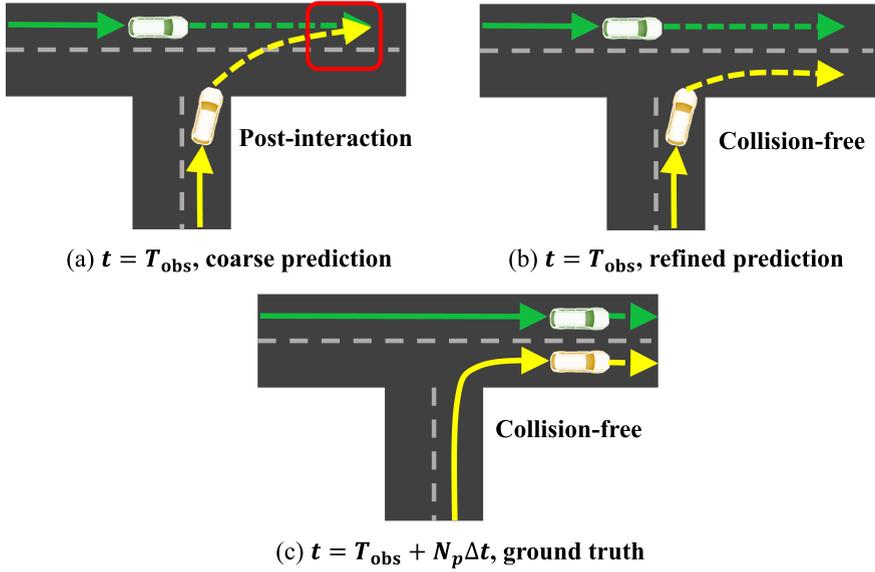


Figure 1: Illustration of trajectory predictions that consider post-interaction behavior. This study proposes a method, i.e., Pioformer, for trajectory predictions of multiple agents. (a) Current time step $t = T_{\text{obs}}$. We observe the post-interaction in the coarse predicted trajectory, which refers to the overlap between the predicted trajectories of two agents at future time step $t = T_{\text{obs}} + N_p \Delta t$. (b) We use Pioformer to make predictions, which utilizes the post-interaction features embedded in coarse trajectories to refine the trajectory predictions. There are no collisions in the refined predicted trajectory. (c) Ground truth trajectories show that there are no collisions between the two agents at time step $t = T_{\text{obs}} + N_p \Delta t$. This indicates that our model can generate refined predictions that align well with the ground truth trajectory.

temporal interactions among multiple traffic agents (Mohamed et al., 2020; Gao et al., 2020; Zeng et al., 2021; Ye et al., 2022; Xu et al., 2022a). For instance, HiVT (Zhou et al., 2022) extracts local context and models global interaction based on GAT and proposes translation-invariant scene representation and rotation-invariant spatial learning modules. Over the years, various variants of GNN have emerged (Xu et al., 2022a; Jia et al., 2023), leading to continuous improvements in prediction performance.

Modeling the interactions among agents for trajectory prediction of autonomous driving has been challenging due to the inherent uncertainty in agents' behavior. The interactions involved in the predicted trajectories of agents, also called post-interactions, have rarely been considered in trajectory prediction models. The predicted trajectories sometimes include unexpected interactions such as traffic conflicts, lane deviations, or even collisions, which need more careful treatments. It is unreasonable to directly remove those collided predictions since the crash could probably happen in some situations. Instead, we come to the idea of enhancing the trajectory prediction by utilizing those interactions involved in the predicted trajectories. Why can future interactions benefit trajectory prediction? The idea comes from the common sense that the participating agents are more likely to make proactive decisions to resolve the risks once unexpected interactions are predicted. That is, those unexpected future interactions are less likely to turn into real risks and most trajectories would be collision-free and traffic rules-compliant at last. Moreover, we can infer that the future predictions can affect current decisions, and thus influence future motions of participating agents and then benefit the trajectory predictions. The interactions within the predicted trajectories, also called post-interactions, are expected to be helpful in trajectory predictions. Current methods for interaction modeling rarely consider the post-interaction and its potential impacts on trajectory predictions. In addition, standard GNN architectures for interaction modeling usually focus on low-order interactions and may fail to capture high-order interactions between the agents which are subtle and uncertain and directly affect the future motions of the agents. To address this, this study proposes a post-interactive multimodal trajectory prediction model, in which the high-order post-interaction features are learned to improve the prediction accuracy in strongly-interactive scenarios. Furthermore, we demonstrate how the multimodal trajectory can be integrated with the downstream planning under uncertainty and improve robustness for autonomous driving.

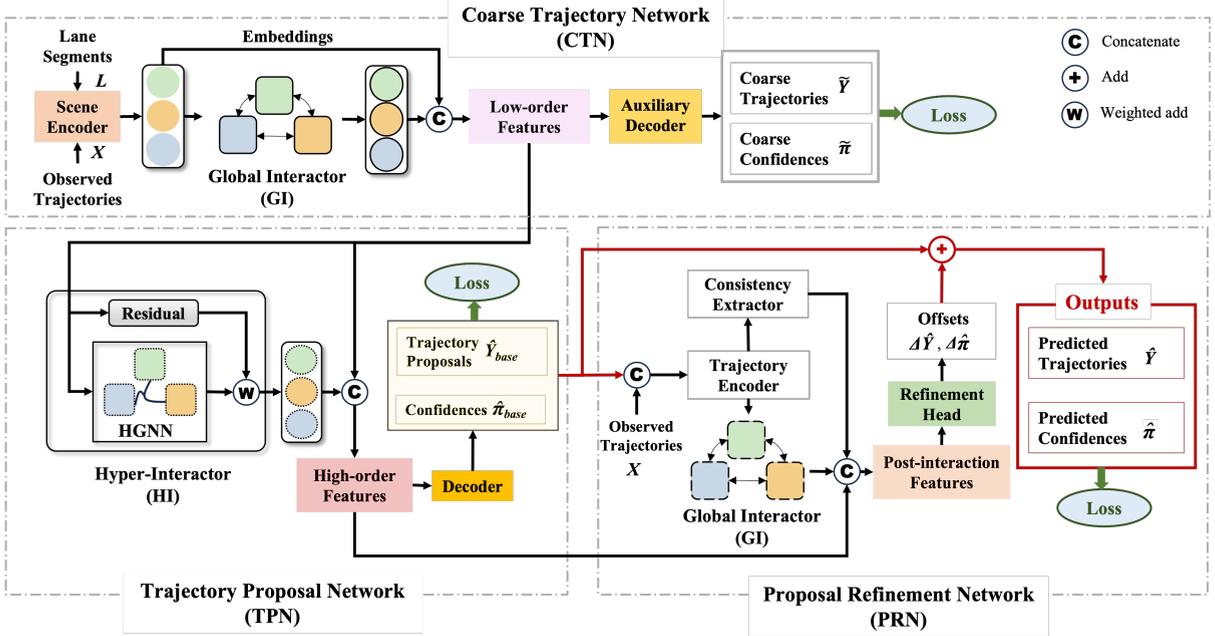


Figure 2: The architecture of Pioformer. The entire model can be divided into three networks: Coarse Trajectory Network (CTN), Trajectory Proposal Network (TPN) and Proposal Refinement Network (PRN). The top one is the CTN, which includes Scene Encoder, Global Interactor (GI), and Auxiliary Decoder, providing coarse predictions and low-order features (contextual information and simple pairwise interaction features). The bottom-left one is the TPN, which includes Hyper-Interactor (HI) and Decoder, generating trajectory proposals and corresponding confidences. HI leverages the low-order features to extract high-order post-interaction features (interactions beyond pairwise relationships) and refine the trajectories at the latent space. The bottom-right one is the PRN, which takes the trajectory proposals combined with observed trajectories as input. It further extracts post-interaction features of different trajectories and explores spatial-temporal consistency features for individual trajectories to generate offsets and refine initial proposals at the trajectory level. Both TPN and PRN are proposed post-interactive networks.

1.1. Overview of this study

This study proposes a post-interactive multimodal trajectory prediction model for autonomous driving, *i.e.*, Pioformer, which aims to model the post-interaction information among agents to enhance the trajectory predictions (see Fig. 1). As shown in Fig. 2, we employ a basic trajectory prediction model as a Coarse Trajectory Network, *i.e.* backbone, to extract the local and global context and low-order interaction features. Next, a Hyper-Interactor module is built based on Hypergraph Neural Network (HGNN) to further learn the high-order post-interactions based on the previous low-order features. We construct the hypergraph topology based on the affinity of agents and then employ a mechanism similar to message-passing in ordinary GNNs to extract post-interaction features. The high-order features are then used to generate multimodal trajectory proposals and their confidences for each agent. The trajectory proposals are concatenated with the observed trajectories and used in Proposal Refinement Network to further extract the post-interaction features and correct unreasonable trajectories. Furthermore, a three-stage training scheme is designed to facilitate the training of Pioformer, as shown in Fig. 3. Extensive experimental results have shown the effectiveness of post-interaction features in enhancing the accuracy of multimodal trajectory prediction for multiple agents. Different from (Zhou et al., 2022; Liu et al., 2024), we propose a post-interaction aware trajectory prediction framework that employs hypergraphs to model high-order post-interaction patterns among vehicles and introduce a three-stage training scheme to facilitate the training process. Our key contributions are summarized as follows:

- We propose a multimodal trajectory prediction framework called Pioformer to leverage the post-interaction features to enable better predictions. A three-stage training scheme is introduced to facilitate the training.
- An HGNN-based Hyper-Interactor module is built to capture the high-order post-interaction features and adaptively refine the coarse trajectories in the latent space. We also design a Proposal Refinement Network,

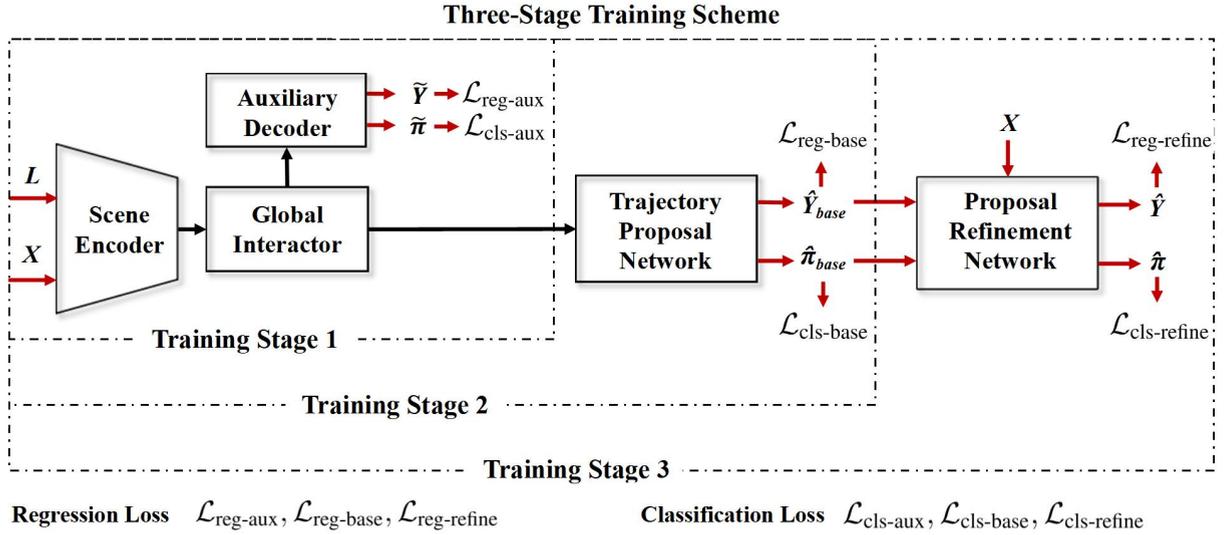


Figure 3: The overview of the three-stage training scheme. The three training stages are executed sequentially from left to right with each subsequent training stage reloading weights from the previous stage. In the first stage, we train CTN to achieve coarse trajectory prediction. In the second stage, we train CTN and TPN together, with the decoder in CTN acting as an auxiliary decoder while the decoder in TPN generates trajectory proposals. In the third stage, we train CTN, TPN and PRN together, with PRN refining the trajectory proposals to produce the final prediction. Each of the three networks has its corresponding trajectory regression loss and confidence classification loss.

which can further utilize the post-interaction features and exploit spatial-temporal consistency to refine the initial trajectory proposals explicitly.

- Experiments demonstrate that our method can handle the trade-off between accuracy and model size (see Fig. 4) and outperform the baseline HiVT on the benchmark *Argoverse 1* dataset, especially in strong-interactive scenarios.

1.2. Organization of the Paper

The subsequent sections of this paper are organized as follows. Section 2 provides an overview of related works. Section 3 presents the problem statement of trajectory prediction. Building on this foundation, Section 4 elaborates on the modules of Pioformer and introduces a three-stage training scheme. Section 5 presents extensive experimental results and analysis. The final section, Section 6 draws comprehensive conclusions of this paper.

2. Related Works

This section introduces the related works on multimodal trajectory prediction, interaction modeling, coarse-to-fine trajectory prediction, and hypergraph neural networks.

2.1. Multimodal Trajectory Prediction

Due to the inherent uncertainties in driving behaviors and the surrounding environment, it is hard to generate one trajectory prediction that matches the ground truth well. Multimodal trajectory prediction (MTP) aims to generate multiple plausible and socially acceptable future predictions for each agent (Huang et al., 2023a), which has been widely studied in recent years (Gupta et al., 2018a; Mohamed et al., 2020; Gu et al., 2022). Current MTP methods can be categorized into two groups, i.e., the noise-based method and the anchor-based method. The noise-based methods typically use generative models to generate multimodal predictions, such as Social GAN (Gupta et al., 2018b) that is based on Generative Adversarial Network (GAN) (Goodfellow et al., 2014), SocialVAE (Xu et al., 2022c) that is based on Conditional Variational Autoencoder (CVAE) (Sohn et al., 2015), STGlow (Liang et al., 2023) that is

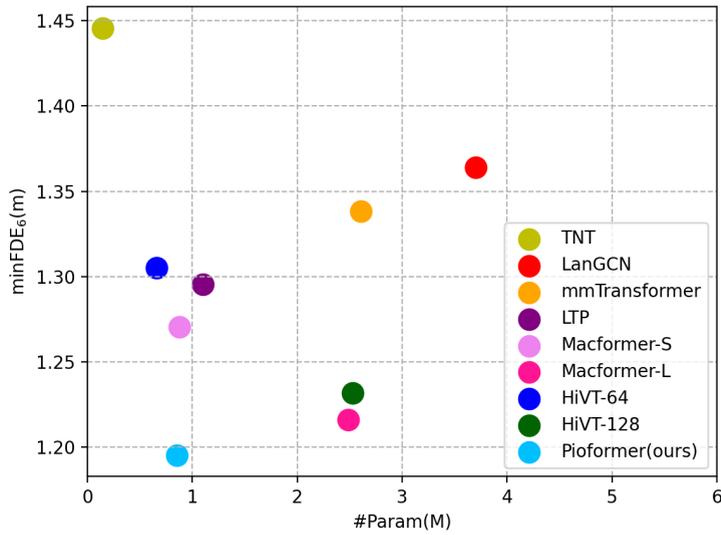


Figure 4: Overview of error-and-size trade-off for the task of trajectory prediction on *Argoverse 1* leaderboard. With compact size, our models outperform most of the state-of-the-art models in prediction accuracy. Especially, the size of Pioformer is approximately one-third that of other models (e.g. HiVT-128 and Macformer-L) that achieve similar accuracy.

based on Normalising Flow (NF) (Papamakarios et al., 2019), and MID (Gu et al., 2022) that is based on Denoise Diffusion Probabilistic Model (DDPM) (Ho et al., 2020). However, those generative models often suffer from issues like training instability and mode collapse or may lack fine-grained details (Kossale et al., 2022; Jiao et al., 2023). The anchor-based frameworks use anchors, such as the endpoints and the prototype trajectories, to guide the model to predict trajectories with controlled behaviors (Huang et al., 2023a; Zhao et al., 2021). Compared with the noise-based methods, the anchor-based methods can generate more controllable predictions conditioned on prior knowledge and thus mitigate the mode collapse issues (Chai et al., 2020; Wang et al., 2023). However, if the anchors are not well defined, the unreachable endpoints or unreasonable prototype trajectories can degrade the performance of the anchor-based approaches (Li et al., 2023; Gu et al., 2021). Different from previous studies, our proposed model can be considered as a dynamic anchor-based method. Rather than relying on predefined anchors, we adopt the Trajectory Proposal Network to generate proposals that serve as anchors. Subsequently, these anchors are further refined by the Proposal Refinement Network, resulting in the final predicted trajectories.

2.2. Interaction Modeling

Modeling the interactions among agents has been essential to accurate trajectory prediction in complicated traffic scenes. The Social Force Model is a popular approach used in crowd simulation that uses handcrafted features for interaction modeling, but it is hard to capture complex interaction behaviors. LSTM-based frameworks have been widely used for interaction modeling, such as Social LSTM (Alahi et al., 2016), TraPHic (Chandra et al., 2019), TrafficPredict (Ma et al., 2019). However, the social pooling strategies in LSTM-based models cannot capture the important interactions between the agents that are located far from each other. Recently, GNNs have become increasingly popular in modeling the spatial-temporal dependencies among agents. For instance, VectorNet (Gao et al., 2020) incorporates the vectorized scene context and agent dynamics for trajectory prediction based on graph neural networks. Social-STGCNN (Mohamed et al., 2020) models the interactions as a graph and proposes a weighted adjacency matrix in which the kernel function quantitatively measures pedestrian interactions. LaneGCN (Liang et al., 2020) extends graph convolutions with multiple adjacency matrices and along-lane dilation to effectively capture the lane graph’s complex topology and long-range dependencies. LanceRCNN (Zeng et al., 2021) learns a local lane graph representation per actor (LaneRoI) to encode its past trajectories and the local map topology. Recent works have expanded the conventional GCN and GAT to learn the interaction features. For instance, HDGT (Jia et al., 2023) utilizes a heterogeneous graph to model the traffic scenes and heterogeneous interactions. Most GNN-based methods

extract low-order pairwise interactions between agents and often overlook the future interactions involved in agents' future trajectories, called post-interactions. Different from previous studies, we build a trainable HGNN to learn the high-order, group-wise interactions that efficiently address complex relational dynamics, especially post-interactions.

2.3. Coarse-to-fine Trajectory Prediction

Coarse-to-fine trajectory prediction refers to a hierarchical trajectory prediction approach where an initial coarse prediction is refined by incorporating more detailed features to improve accuracy. The refinement module has been applied extensively in computer vision (CV). It is supported by its implementation in several tasks such as Cascade RCNN (Cai and Vasconcelos, 2018) for object detection and RefineNet (Lin et al., 2020) for semantic segmentation. These studies showcase the potential of the refinement module in improving prediction accuracy at a fine-grained level. In this task, trajectory prediction refinement aims to correct the predicted trajectories. TPNNet (Fang et al., 2020) uses a two-stage framework for multimodal motion prediction of vehicles and pedestrians, in which the first stage generates a candidate set of future trajectories and the second stage performs classification and refinement on the proposals and selects the best one as the final prediction. MTR (Shi et al., 2022a) and MTR-A (Shi et al., 2022b) refine the trajectories by iteratively gathering fine-grained trajectory features. DESIRE (Lee et al., 2017) refines the trajectories by incorporating the semantic context of scenes and the social interactions among agents. The Proposal Refinement Network of our proposed model is inspired by the LAformer (Liu et al., 2024), and it aims to learn the trajectory offsets for correcting the trajectories. Different from this work, our approach extends upon it by incorporating post-interaction features, refining the confidence of each mode, and reusing an existing module to enhance model robustness. In addition to explicit refinement, we also refine the predicted trajectories in the high-dimensional latent space, aiming to move the initial trajectories to more plausible positions as much as possible.

2.4. Hypergraph Neural Networks

HGNN applications have been extensively adopted in diverse fields such as CV (Han et al., 2023), natural language processing (NLP) (Ding et al., 2020), multimodal learning (Kim et al., 2020), and data mining (Luo et al., 2023). HGNNs use hypergraphs to capture intricate, high-order dependencies between nodes that go beyond pairwise relationships (Feng et al., 2018; Gao et al., 2023). Traffic scenarios involve numerous intricate high-order interactions among agents, such as interactions among multiple vehicles engaged in cooperative driving, as well as the sudden changing behaviors of vehicles in response to dynamic events. These can be effectively modeled using HGNN, where nodes represent agents and hyperedges are high-order interactions among agents (Xu et al., 2022a,b). Most research primarily focuses on modeling pairwise low-order interactions while neglecting high-order interactions, failing to comprehensively capture multi-faceted interactions among agents. We introduce HGNN to extract high-order post-interaction features and we make adaptive adjustments to accommodate the specific characteristics of autonomous driving and propose a corresponding three-stage training scheme to facilitate learning. Different from works such as GroupNet (Xu et al., 2022a), since our task focuses on vehicle trajectory prediction rather than pedestrian trajectory prediction. The interaction intensity between vehicles is significantly weaker than that between pedestrians. Therefore, we employ a single-scale hypergraph instead of a multi-scale hypergraph, reducing the computational complexity of the model.

3. Problem Formulation

In this section, we formulate the multimodal trajectory prediction problem as a sequence-to-sequence model, and our goal is to generate multiple predictions for each agent based on the observed trajectories.

3.1. Multimodal Trajectory Prediction

We focus on a multimodal trajectory prediction model, which aims to generate a distribution of future trajectories for each agent. The trajectory prediction model leverages the observed trajectories of multiple agents and lane segment information to predict the trajectory distributions, as shown in Fig. 2. Then, the best-predicted trajectory is selected from the trajectory distribution by minimizing the prediction errors, as shown in Fig. 5. That is, the trajectory mode that is located closest to the ground truth trajectory is chosen as the best-predicted one.

We denote the observed state of the agent i at time step t as \mathbf{x}_i^t , the observed trajectory is written as $\mathbf{X} = \{\mathbf{x}_i^t\}$, where $i \in \{1, 2, \dots, M\}$ and $t \in \{1, 2, \dots, T_{\text{obs}}\}$, M is the number of agents, T_{obs} is the observed time steps. Similarly, we

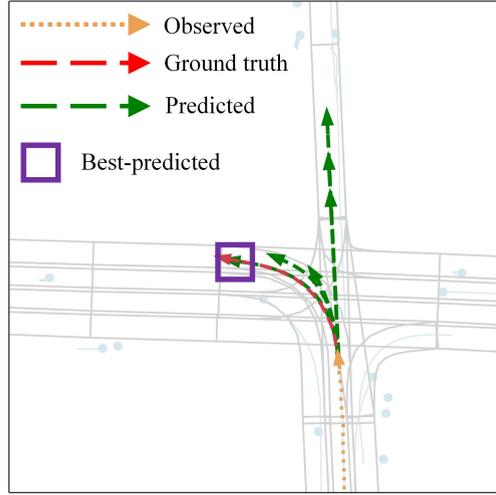


Figure 5: An example of multimodal trajectory prediction in a traffic scenario. We predict multiple possible trajectories based on the observed trajectories of agents and lane information.

denote the predicted state of the agent i at time step t of the mode k as $\hat{\mathbf{y}}_{i,k}^t$, and the predicted trajectories $\hat{\mathbf{Y}} = \{\hat{\mathbf{y}}_{i,k}^t\}$, where $t \in \{T_{\text{obs}} + 1, \dots, T_{\text{pred}}\}$, $k \in \{1, 2, \dots, K\}$, and K is the number of prediction modes. Then, the trajectory prediction problem can be represented as,

$$\mathbf{f} : (\mathbf{X}, L) \rightarrow \hat{\mathbf{Y}} \quad (1)$$

where \mathbf{f} is the multimodal trajectory prediction model, \mathbf{X} denotes the observed trajectories, L is the lane segment encodings, $\hat{\mathbf{Y}}$ denotes the predicted trajectories.

To find the best-predicted trajectory from K modes, the accumulated prediction error of each prediction mode is calculated based on the ground truth trajectory. The trajectory mode that has the minimum accumulated prediction error is selected as the best-predicted trajectory,

$$k_i^* = \underset{k}{\operatorname{argmin}} \sum_{t=T_{\text{obs}}+1}^{T_{\text{pred}}} \|\mathbf{y}_i^t - \hat{\mathbf{y}}_{i,k}^t\| \quad k \in \{1, 2, \dots, K\}, \quad (2)$$

where k_i^* is the best-predicted mode of agent i , \mathbf{y}_i^t is the ground truth state of agent i at time step t , $\hat{\mathbf{y}}_{i,k}^t$ is the predicted state of agent i of mode k at time step t .

3.2. Scene Representation

A translation-invariant scene representation is employed in this study, similar to (Zhou et al., 2022). We represent the agent trajectories and lane segments using relative positions to ensure translation invariance. Therefore, the entire scene can be represented as a set of vectors.

- **Agent:** the position of agent i at time t can be represented as $\mathbf{p}_i^t = (p_{i,x}^t, p_{i,y}^t)$. The trajectory of agent i can then be represented as $\{\mathbf{p}_i^t - \mathbf{p}_i^{t-1}\}_{t=1}^{T_{\text{obs}}}$.
- **Lane:** a lane segment with the starting point \mathbf{p}_l^s and the ending point \mathbf{p}_l^e can be represented as $\mathbf{p}_l^e - \mathbf{p}_l^s$.
- **Scene elements:** the spatial relationship between scene elements can also be represented as relative position vectors, and $\mathbf{p}_j^t - \mathbf{p}_i^t$ denotes the position vector of agent j relative to agent i at time step t , and $\mathbf{p}_l^s - \mathbf{p}_i^t$ denotes the position vector of lane segment l related to agent i at time step t .
- **Local region:** like (Zhou et al., 2022), we divide the scene into local regions, selecting one agent as the central agent, which is located in the center of the local region. We use the latest trajectory segment $\mathbf{p}_i^{T_{\text{obs}}} - \mathbf{p}_i^{T_{\text{obs}}-1}$

of central agent i as the reference vector and all vectors are rotated according to the orientation angle θ_i of the reference vector. The rotation matrix is represented as \mathbf{R}_i^\top matrix parameterized by θ_i .

4. Methodology

This section introduces our model, Pioformer, and three main parts, including Coarse Trajectory Network, Trajectory Proposal Network, and Proposal Refinement Network.

4.1. Overall Framework

Fig. 2 illustrates our model, which can be divided into three networks, i.e., Coarse Trajectory Network (CTN), Trajectory Proposal Network (TPN) and Proposal Refinement Network (PRN). The top one is the CTN, which includes Scene Encoder, Global Interactor (GI), and Auxiliary Decoder, providing coarse predictions and low-order interaction features (contextual information and simple pairwise interaction features). The bottom-left one is the TPN, which includes Hyper-Interactor (HI) and Decoder, generating trajectory proposals and corresponding trajectory confidences based on low-order interaction features. In particular, the Hyper-Interactor module in TPN is built to extract high-order post-interaction features (interactions beyond pairwise relationships) leveraging the low-order interaction features. The bottom-right one is the PRN, which takes the trajectory proposals from TPN and observed trajectories as input and output refined trajectories. We note that PRN is designed to further extract post-interaction features of multiple agents and explore spatial-temporal consistency features to refine trajectory proposals.

4.2. Coarse Trajectory Network (CTN)

CTN aims to encode the observed trajectories and context information and then decode coarse-grained trajectory predictions while providing low-order interaction features for subsequent networks. In this paper, we adopt HiVT as our CTN, which consists of three components: Scene Encoder, Global Interactor and Auxiliary Decoder. This section will provide a brief introduction to these modules, for detailed information please refer to (Zhou et al., 2022).

4.2.1. Agent-Agent Interaction

In each local region, we extract the features of the central agent i , denoted as \mathbf{z}_i^t , as well as the features of any neighboring agents j w.r.t. the agent i , denoted as \mathbf{z}_{ij}^t , at time step t :

$$\mathbf{z}_i^t = \phi_{\text{center}} \left(\left[\mathbf{R}_i^\top \left(\mathbf{p}_i^t - \mathbf{p}_i^{t-1} \right), \mathbf{a}_i \right] \right), \quad (3)$$

$$\mathbf{z}_{ij}^t = \phi_{\text{nbr}} \left(\left[\mathbf{R}_i^\top \left(\mathbf{p}_j^t - \mathbf{p}_j^{t-1} \right), \mathbf{R}_i^\top \left(\mathbf{p}_j^t - \mathbf{p}_i^t \right), \mathbf{a}_j \right] \right), \quad (4)$$

where $\phi_{\text{center}}(\cdot)$ and $\phi_{\text{nbr}}(\cdot)$ are both MLP blocks, \mathbf{a}_i denotes the semantic attributes of agent i . We then utilize a rotation-invariant cross-attention mechanism to obtain spatial interaction features among the agents. Subsequently, we need to capture temporal dependencies by aggregating each agent's embeddings in chronological order using a Transformer module to obtain spatio-temporal interaction embeddings.

4.2.2. Agent-Lane Interaction

We need to leverage features from the surrounding lane segments to guide the future motion of the agents:

$$\mathbf{z}_{il} = \phi_{\text{lane}} \left(\left[\mathbf{R}_i^\top \left(\mathbf{p}_l^e - \mathbf{p}_l^s \right), \mathbf{R}_i^\top \left(\mathbf{p}_l^s - \mathbf{p}_i^{T_{\text{obs}}^t} \right), \mathbf{a}_l \right] \right), \quad (5)$$

where \mathbf{z}_{il} represents the embedding between the agent i and the lane segment l , $\phi_{\text{lane}}(\cdot)$ is an MLP block, \mathbf{a}_l is the semantic attributes of the lane segment. Similarly, the agent-lane attention is calculated using the cross-attention module. The final output ψ_i^{local} , referred to as the local embedding of the agent i is the fusion of spatial-temporal interaction features in the local region.

4.2.3. Global Interactor (GI)

This module aims to further extract long-range dependencies and establish geometric relationships between local regions to bridge the differences between their respective coordinate systems. For agents i and j , we continue to employ an MLP block $\phi_{\text{rel}}(\cdot)$ to extract embeddings representing the pairwise interactions between the two agents, denoted as \mathbf{s}_{ij} :

$$\mathbf{s}_{ij} = \phi_{\text{rel}} \left(\left[\mathbf{R}_i^T (\mathbf{p}_j^{\text{obs}} - \mathbf{p}_i^{\text{obs}}), \cos(\Delta\theta_{ij}), \sin(\Delta\theta_{ij}) \right] \right), \quad (6)$$

where $\mathbf{p}_j^{\text{obs}} - \mathbf{p}_i^{\text{obs}}$ represents the differences between agent i 's and agent j 's coordinate frames, and $\Delta\theta_{ij}$ denotes $\theta_j - \theta_i$. Subsequently, the cross-attention mechanism is employed again to further obtain the global embedding $\boldsymbol{\psi}_i^{\text{global}}$. Up to now, we have successfully captured both local and long-range dependencies from the agent's observed trajectory and obtained low-order interaction features.

4.2.4. Auxiliary Decoder

We employ the Laplace Mixture Model (LMM) to parameterize the distribution of future trajectories, and the predicted multimodal trajectory distribution of agent i is denoted as $\tilde{\mathbf{Y}}_i$:

$$\tilde{\mathbf{Y}}_i = \sum_{k=1}^K \tilde{\pi}_{i,k} \mathfrak{L}(\tilde{\boldsymbol{\mu}}_{i,k}, \tilde{\boldsymbol{\mathbf{b}}}_{i,k}), \quad (7)$$

where $\mathfrak{L}(\cdot, \cdot)$ denotes Laplace distribution, the predicted probability that agent i adopts mode k is $\tilde{\pi}_{i,k}$, the predicted location and uncertainty parameters of each Laplace component are denoted as $\tilde{\boldsymbol{\mu}}_{i,k}$ and $\tilde{\boldsymbol{\mathbf{b}}}_{i,k}$. The decoder takes the previous embeddings as input and employs three separate MLP blocks to generate the three parameters of the Laplace distribution. We note that this auxiliary decoder is only employed during the training process and is not used in the inference stage.

4.3. Trajectory Proposal Network (TPN)

TPN includes Hyper-Interactor (HI) and Main Decoder, which aims to model high-order interactions using HGNNs based on the low-order interaction features provided by the CTN and output multimodal trajectory proposals and related trajectory confidences, (see Fig. 6). HI takes previous low-order embeddings as input, which are subsequently decoded into real trajectories in CTN. This implies that these embeddings represent the future trajectories in a high-dimensional space. We then use an MLP block to fuse the previous embeddings, thereby generating features relevant to agents' future trajectories. With these embeddings, we can proceed with post-interaction modeling and learning.

4.3.1. Hyper-Interactor (HI)

We use hypergraph neural networks instead of the graph neural networks to learn the high-order post-interaction features, facilitating an in-depth exploration of both low-order and high-order interaction features.

Hypergraph. Mathematically, a hypergraph is a mathematical structure that extends the concept of a simple graph. In a simple graph, the edges represent pairwise connections between nodes. However, in a hypergraph, the hyperedges are capable of connecting any number of nodes, thus allowing for a more flexible and expressive representation of relationships among nodes. Formally, a hypergraph can be represented as: $\mathcal{G} = (\mathcal{V}, \mathcal{E})$, where $\mathcal{V} = \{v_1, v_2, \dots, v_{N_V}\}$ denotes the vertex set which contains N_V vertices, $\mathcal{E} = \{e_1, e_2, \dots, e_{N_E}\}$ denotes the hyperedge set which contains N_E hyperedges. The topological structure of a hypergraph can be represented using an incidence matrix $\mathbf{H} \in \mathbb{R}^{|\mathcal{V}| \times |\mathcal{E}|}$ which is a binary matrix with entries defined as

$$h(v, e) = \begin{cases} 1, & \text{if } v \in e \\ 0, & \text{if } v \notin e, \end{cases} \quad (8)$$

where $v \in \mathcal{V}$ and $e \in \mathcal{E}$ represent vertex and hyperedge respectively (Feng et al., 2018). To make it clear, we provide an example in Fig. 7. It consists of 5 vertices represented by rounded rectangles and 3 hyperedges represented by curves. The rows and columns of the incidence matrix \mathbf{H} represent vertices and hyperedges, with the elements indicating whether the vertex belongs to the hyperedge. For instance, hyperedge e_1 connects v_1, v_2 and v_5 , then the first column of the matrix would be represented as $(1, 1, 0, 0, 1)$.

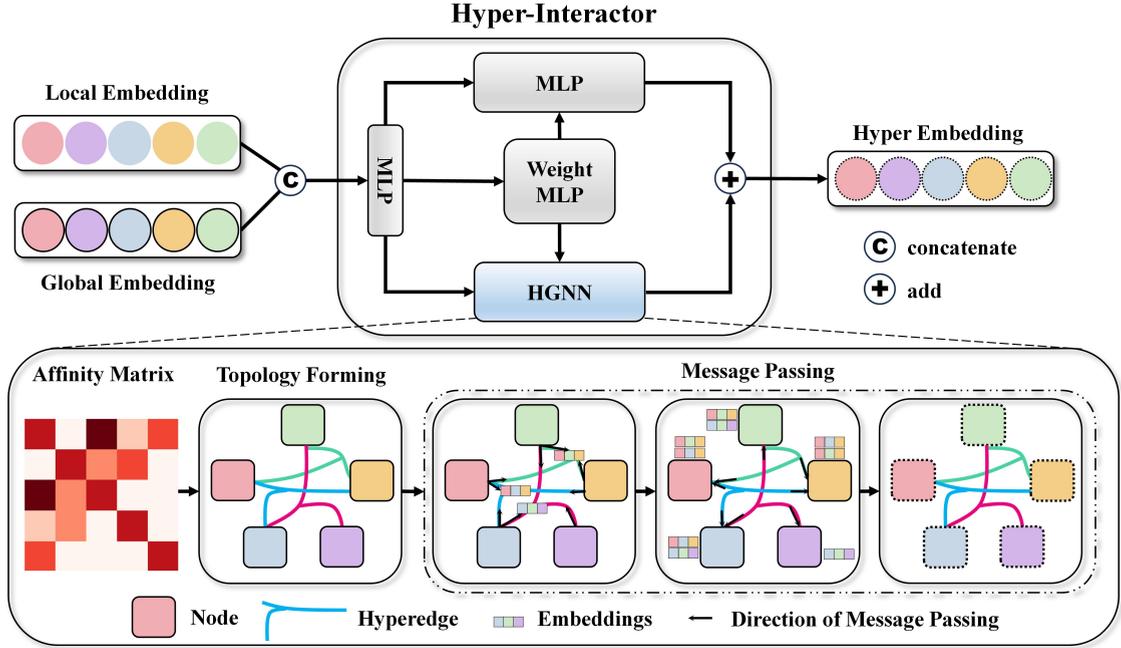


Figure 6: The architecture of the Hyper-Interactor module. The core of Hyper-Interactor is an HGNN, which is aimed at capturing high-order interaction features. The hypergraph topology is adaptively constructed by calculating the affinity between agent nodes, which is then used for message propagation. Additionally, an MLP block is utilized as a residual branch to provide low-order interaction features. The high-order and low-order features are then integrated in a weighted manner to obtain the final output.

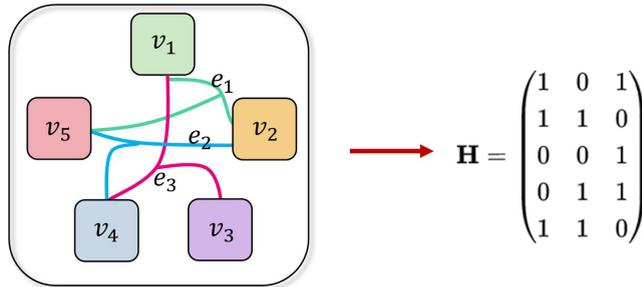


Figure 7: An example of a hypergraph. It consists of 5 vertices represented by rounded rectangles and 3 hyperedges represented by curves. The rows and columns of the incidence matrix \mathbf{H} represent vertices and hyperedges, with the elements indicating whether the vertex belongs to the hyperedge. For instance, hyperedge e_1 connects v_1 , v_2 and v_5 , then the first column of the matrix would be represented as $(1, 1, 0, 0, 1)$.

Hypergraph Neural Networks. HGNNs are a class of neural networks that leverage the capacity of hypergraphs to capture complex, high-order dependencies among nodes, going beyond pairwise relationships. HGNNs employ an iterative message-passing strategy for information propagation, similar to GNNs, facilitating the integration of the local and global hypergraph structures in the update of node features. During each message-passing iteration, nodes engage in information exchange with their connected hyperedges. In turn, the hyperedges transmit this rich information to other connected nodes.

The initial step for constructing an HGNN is to define the topology of the hypergraph. As interactions between agents in traffic scenarios are unknown, it is hard to define the topology. Therefore, following GroupNet (Xu et al., 2022a), we can use an affinity matrix to calculate the correlation between different agents to identify a suitable topology for building a hypergraph. With cosine similarity adopted as a metric, the entries of the affinity matrix $\mathbf{A} \in \mathbb{R}^{M \times M}$

are defined as

$$\mathbf{A}_{i,j} = \frac{\boldsymbol{\tau}_i^\top \boldsymbol{\tau}_j}{\|\boldsymbol{\tau}_i\| \|\boldsymbol{\tau}_j\|}, \quad (9)$$

where $\boldsymbol{\tau}_i$ is the future trajectory feature of agent i . We can further use the affinity matrix to search for high-density submatrices. We assume that the agent nodes located on the same hyperedge should have a higher affinity with each other. The search for hyperedges with a predetermined number of connected nodes S can be viewed as solving an optimization problem. We aim to find an optimal hyperedge that can maximize the norm of the affinity matrix,

$$\begin{aligned} e_i &= \arg \max_{\Omega \subseteq \mathcal{V}} \|\mathbf{A}_{\Omega,\Omega}\|_{1,1}, \\ \text{s.t. } |\Omega| &= S, v_i \in \Omega, i = 1, \dots, M, \end{aligned} \quad (10)$$

where e_i is the optimal hyperedge of agent i , $\|\cdot\|_{1,1}$ is an entrywise matrix norm that measures the sum of the absolute values of all elements in the matrix, and $|\Omega|$ is the cardinal of the vertex set Ω . The first condition in the formula imposes a limitation on the number of agents on each hyperedge. The following condition guarantees that no agent is left out. We can use an exhaustive approach to solve the optimization problem directly when the number of agents in practical traffic scenarios is small, usually in the range of a few dozen.

With the topological structure of the hypergraph, we can utilize it for message-passing. The message-passing process in a hypergraph is similar to that in a simple graph. During the node-to-hyperedge phase, node embeddings are aggregated into hyperedges, which can be represented as:

$$\mathbf{e}_i = \phi_{\text{hyperedge}} \left(\sum_{v_j \in e_i} \mathbf{v}_j \right). \quad (11)$$

where the embedding of each node is denoted as \mathbf{v}_j , \mathbf{e}_i is the embedding of hyperedge e_i and $\phi_{\text{hyperedge}}(\cdot)$ is an MLP to obtain the interaction embeddings.

During the hyperedge-to-node phase, nodes update their embeddings by incorporating the interaction embeddings learned from the hyperedges. The updated embedding, \mathbf{v}_i , can be obtained through:

$$\mathbf{v}_i \leftarrow \phi_{\text{node}} \left(\left[\mathbf{v}_i, \sum_{e_j \in \Omega_i} \mathbf{e}_j \right] \right), \quad (12)$$

where Ω_i is a set of hyperedges containing node i , and $\phi_{\text{node}}(\cdot)$ is an MLP to obtain updated embeddings. This message-passing process facilitates the iterative refinement of node embeddings.

Consequently, the HGNN yields embeddings denoted as \mathbf{r}^{high} that represent the high-order interactions of various agents and their latent impact in the future.

In many traffic scenarios, especially in straight-line situations with few surrounding agents, the post-interactions among agents are relatively slight. In such scenarios, the initially predicted trajectories may already meet the requirements. Employing the HGNN for processing may generate overfitting issues and result in undesired lane changes. Thus, we use a simple MLP $\phi_s(\cdot)$ as a residual block to leverage the embedding $\boldsymbol{\tau}_i$ to compute low-order interaction features:

$$\mathbf{r}_i^{\text{low}} = \phi_s(\boldsymbol{\tau}_i). \quad (13)$$

Where $\mathbf{r}_i^{\text{low}}$ is the low-order interaction features.

To adaptively incorporate the high-order and low-order interaction features, a weighting mechanism is adopted to incorporate the results of the HGNN and the simple MLP block:

$$\boldsymbol{\omega}_i = \sigma(\phi_{\text{weight}}(\boldsymbol{\tau}_i)), \quad (14)$$

$$\boldsymbol{\psi}_i^{\text{hyper}} = \mathbf{r}_i^{\text{high}} \odot \boldsymbol{\omega}_i + \mathbf{r}_i^{\text{low}} \odot (1 - \boldsymbol{\omega}_i), \quad (15)$$

where $\phi_{\text{weight}}(\cdot)$ is an MLP block, $\sigma(\cdot)$ denotes the sigmoid function, \odot denotes the element-wise product. ψ_i^{hyper} is the weighted post-interaction embedding of agent i . This weighting mechanism enables the flexibility to selectively leverage the output of the MLP and the output of the HGNN to accommodate the interactions in different traffic scenarios.

4.3.2. Main Decoder

Hyper-Interactor module captures the dynamic nature of high-order post-interactions, enabling refinement of the original trajectory embeddings based on the interactions between the initially predicted trajectories in the latent space. We subsequently employ a decoder to generate trajectory proposals, which serves as the main decoder. Similar to the auxiliary decoder, the trajectory proposal $\hat{\mathbf{Y}}_{\text{base},i}$ generated by the main decoder can be expressed as:

$$\hat{\mathbf{Y}}_{\text{base},i} = \sum_{k=1}^K \hat{\pi}_{\text{base},i,k} \mathfrak{G}(\hat{\boldsymbol{\mu}}_{\text{base},i,k}, \hat{\mathbf{b}}_{\text{base},i,k}), \quad (16)$$

where $\hat{\pi}_{\text{base},i,k}$, $\hat{\boldsymbol{\mu}}_{\text{base},i,k}$ and $\hat{\mathbf{b}}_{\text{base},i,k}$ are the predicted parameters associated to the trajectory proposals. The main decoder leverages the high-order features along with the previously obtained low-order features to generate more accurate trajectory proposals, which serve as anchor trajectories for the Proposal Refinement Network.

4.4. Proposal Refinement Network (PRN)

PRN takes the trajectory proposals and observed trajectories as input, and aims to learn the post-interaction features among multiple agents and explores spatial-temporal consistency features to generate offsets and refine trajectory proposals.

4.4.1. Post-interaction Embedding

The observed trajectory \mathbf{X}_i is concatenated with the trajectory proposal $\hat{\mathbf{Y}}_{\text{base},i}$, and then sent to a trajectory encoder to represent the trajectory embedding η_i :

$$\eta_i = \mathbf{f}_{\text{encode}}([\mathbf{X}_i, \hat{\mathbf{Y}}_{\text{base},i}]), \quad (17)$$

where $\mathbf{f}_{\text{encode}}(\cdot)$ denotes the trajectory encoder that is built with a bi-directional GRU (Chung et al., 2014).

The trajectories from the previous decoder may have misalignment or inconsistency with observed trajectories, such as sudden acceleration, excessively large steering angles, etc. To address this, we employ a three-layer MLP block to build the consistency extractor $\phi_{\text{consist}}(\cdot)$ based on the trajectory embedding η_i ,

$$\psi_i^{\text{consist}} = \phi_{\text{consist}}(\eta_i), \quad (18)$$

where ψ_i^{consist} denotes the trajectory consistency embedding, which is expected to capture the spatial-temporal consistency features of the entire trajectory. This can be considered as interactions between future trajectories and observed trajectories, which enables a better understanding of the agent's interaction behavior. The trajectory consistency embedding is used in the subsequent trajectory refinement process in order to improve reasonableness and adherence to physical constraints.

To further explore post-interaction features, we continue to leverage trajectory embeddings η_i over the entire time horizon. Due to the use of relative coordinates to represent contextual scenes, additional information about the geometric relationships between agents must be provided following a straightforward encoding process. Thus, the trajectory embedding η_i is also fed into Global Interactor module to obtain the spatial-temporal post-interaction features:

$$\psi_i^{\text{post-int}} = \mathbf{f}_{\text{global}}(\eta_i), \quad (19)$$

where $\psi_i^{\text{post-int}}$ represents post-interaction embeddings, $\mathbf{f}_{\text{global}}(\cdot)$ is the Global Interactor. By reusing the previous modules, we can reduce the number of parameters and allow Global Interactor to be trained sufficiently.

4.4.2. Trajectory Output

We have obtained multiple embeddings, including ψ_i^{local} , ψ_i^{global} , ψ_i^{hyper} , ψ_i^{consist} , $\psi_i^{\text{post-int}}$, which are fed into the Refinement Head in PRN to generate the trajectory offsets for correcting the trajectory proposals. Considering the uncertainties involved in multiple trajectory modes, the PRN uniformly corrects the trajectory offsets and confidences of each mode. Here, a three-layer MLP block $\phi_{\text{refine-traj}}(\cdot)$ is adopted to generate position offsets $\Delta\hat{\mathbf{Y}}_i$ for the trajectory of agent i , and a three-layer MLP $\phi_{\text{refine-conf}}(\cdot)$ is used to generate confidence offsets $\Delta\hat{\pi}_i$:

$$\Delta\hat{\mathbf{Y}}_i = \phi_{\text{refine-traj}}\left(\psi_i^{\text{local}}, \psi_i^{\text{global}}, \psi_i^{\text{hyper}}, \psi_i^{\text{consist}}, \psi_i^{\text{post-int}}\right). \quad (20)$$

$$\Delta\hat{\pi}_i = \phi_{\text{refine-conf}}\left(\psi_i^{\text{local}}, \psi_i^{\text{global}}, \psi_i^{\text{hyper}}, \psi_i^{\text{consist}}, \psi_i^{\text{post-int}}\right), \quad (21)$$

With the trajectory offset $\Delta\hat{\mathbf{Y}}_i$ and related confidence offset $\Delta\hat{\pi}_i$, we can compute the final trajectory $\hat{\mathbf{Y}}_i$ and confidence $\hat{\pi}_i$ of agent i :

$$\hat{\mathbf{Y}}_i = \hat{\mathbf{Y}}_{\text{base},i} + \Delta\hat{\mathbf{Y}}_i, \quad (22)$$

$$\hat{\pi}_i = \hat{\pi}_{\text{base},i} + \Delta\hat{\pi}_i. \quad (23)$$

where $\hat{\mathbf{Y}}_{\text{base},i}$ and $\hat{\pi}_{\text{base},i}$ are the trajectory proposal and base confidence from the TPN, respectively. As PRN can learn the refinement offsets based on those embeddings. Its refinement offsets will be relatively small if the prediction mode is reasonable; otherwise, the refinement offsets will become very large to correct the unreasonable modes. In this way, we can relocate the trajectories and further approximate the ground truth trajectories.

4.5. Three-stage Training Scheme

We propose a three-stage training scheme to facilitate the learning of interaction features at different levels while preventing TPN and PRN from interfering with CTN during the training process, as shown in Fig. 3. The three training stages are executed sequentially with each subsequent training stage reloading weights from the previous stage. In the first stage, we train CTN to achieve coarse trajectory prediction with an auxiliary decoder in an end-to-end manner. In the second stage, we train CTN and TPN together, and use a main decoder in TPN to generate trajectory proposals. In the third stage, we train CTN, TPN and PRN together, in which PRN is designed to refine the trajectory proposals to produce the final prediction.

Like the previous models (Liang et al., 2020; Zhou et al., 2022; Liu et al., 2024), we use two types of loss function, including regression loss and classification loss. During the calculation of the regression losses, only the error between the best-predicted trajectory and the ground truth is considered. Therefore, regardless of whether the multimodal trajectories are generated by the decoder or PRN, the rest trajectories with destinations farther from the ground truth destinations will not be used to calculate the regression loss. We note that each of the three networks has its corresponding trajectory regression loss and confidence classification loss.

In the first stage, the regression loss $\mathcal{L}_{\text{reg-aux}}$, uses a winner-takes-all strategy to minimize the error between the best-predicted and ground truth trajectories. We employ a negative log-likelihood loss for the auxiliary decoder:

$$\mathcal{L}_{\text{reg-aux}} = -\frac{1}{N} \sum_{t=T_{\text{obs}}+1}^{T_{\text{pred}}} \log P(\mathbf{Y}^t | \tilde{\boldsymbol{\mu}}_{k^*}^t, \tilde{\mathbf{b}}_{k^*}^t), \quad (24)$$

where N is the number of predicted time steps, k^* is the best-predicted mode, \mathbf{Y}^t denotes the ground truth trajectory at time step t , $\tilde{\boldsymbol{\mu}}_{k^*}^t$ and $\tilde{\mathbf{b}}_{k^*}^t$ are parameters of the Laplace distribution of the best-predicted trajectory at time step t . The classification loss $\mathcal{L}_{\text{cls-aux}}$ is defined with a cross-entropy loss, which optimizes the mixing coefficients to allow the model to discriminate the trajectory confidences of different modes:

$$\mathcal{L}_{\text{cls-aux}} = -\sum_{k=1}^K \pi_k \log(\tilde{\pi}_k), \quad (25)$$

where π_k denotes the ground truth confidence of mode k and $\tilde{\pi}_k$ denotes the predicted confidence of mode k . Overall, the loss of the first stage \mathcal{L}_{s1} is defined as the sum of the regression loss and the classification loss:

$$\mathcal{L}_{s1} = \mathcal{L}_{\text{reg-aux}} + \mathcal{L}_{\text{cls-aux}}. \quad (26)$$

We adhere to the practice of assigning equal weights to the two loss functions, like HiVT (Zhou et al., 2022).

In the second stage, we introduce the HI module along with a new decoder that serves as the main decoder, which predicts the trajectory proposal $\hat{\mathbf{Y}}_{\text{base}}$ and $\hat{\pi}_{\text{base}}$. We note that the decoder in the first stage serves as the auxiliary decoder to guide the training process, which predicts the coarse trajectory $\tilde{\mathbf{Y}}$ and $\tilde{\pi}$. We reload the weights obtained in the first stage and copy the weights of the auxiliary decoder to the new main decoder. In this stage, we train CTN and TPN together. The loss function in the second stage \mathcal{L}_{s2} is similar to that in the first stage, defined as:

$$\mathcal{L}_{s2} = \mathcal{L}_{s1} + \lambda_1(\mathcal{L}_{\text{reg-base}} + \mathcal{L}_{\text{cls-base}}), \quad (27)$$

where $\mathcal{L}_{\text{reg-base}}$ and $\mathcal{L}_{\text{cls-base}}$ are the regression loss and the classification loss for the main decoder in TPN, λ_1 is a hyperparameter used to balance the loss in stage 1 and the loss of the trajectory decoder in TPN. Similarly, $\mathcal{L}_{\text{reg-base}}$ and $\mathcal{L}_{\text{cls-base}}$ is defined with negative log-likelihood loss and cross-entropy loss:

$$\mathcal{L}_{\text{reg-base}} = -\frac{1}{N} \sum_{t=T_{\text{obs}}+1}^{T_{\text{pred}}} \log P\left(\mathbf{Y}^t \mid \hat{\boldsymbol{\mu}}_{\text{base},k^*}^t, \hat{\mathbf{b}}_{\text{base},k^*}^t\right), \quad (28)$$

$$\mathcal{L}_{\text{cls-base}} = -\sum_{k=1}^K \pi_k \log(\hat{\pi}_{\text{base},k}), \quad (29)$$

where $\hat{\boldsymbol{\mu}}_{\text{base},k^*}^t$ and $\hat{\mathbf{b}}_{\text{base},k^*}^t$ are distribution parameters of the best-predicted trajectory from the trajectory proposals. $\hat{\pi}_{\text{base},k}$ denotes the predicted confidence in the second stage.

In the third stage, we introduce PRN to refine the trajectories explicitly. In this stage, we aim to minimize the distance between the refined trajectories and the ground truth trajectories. To this end, we employ the Smooth L1 loss function as regression loss, which is defined as:

$$\mathcal{L}_{\text{reg-refine}} = \frac{1}{N} \sum_{t=T_{\text{obs}}+1}^{T_{\text{pred}}} \text{SmoothL1}\left(\hat{\mathbf{Y}}_{k^*}^t - \mathbf{Y}^t\right), \quad (30)$$

$$\text{SmoothL1}(x) = \begin{cases} 0.5x^2 & \text{if } \|x\| < 1, \\ \|x\| - 0.5 & \text{otherwise.} \end{cases} \quad (31)$$

where $\hat{\mathbf{Y}}_{k^*}^t$ denotes the best-predicted trajectory at timestep t . We then use the refined confidences to calculate the classification loss $\mathcal{L}_{\text{cls-refine}}$:

$$\mathcal{L}_{\text{cls-refine}} = -\sum_{k=1}^K \pi_k \log(\hat{\pi}_k). \quad (32)$$

where $\hat{\pi}_k$ denotes the predicted confidence in the third stage. We train all the networks together in this stage. Thus, the total loss for the third stage \mathcal{L}_{s3} can be represented as follows:

$$\mathcal{L}_{s3} = \mathcal{L}_{s2} + \lambda_2(\mathcal{L}_{\text{reg-refine}} + \mathcal{L}_{\text{cls-refine}}), \quad (33)$$

where λ_2 is a hyperparameter used to balance the loss of stage 2 and the loss of the trajectory outputs in the refinement network. The overall framework is shown in Fig. 3. When the three-stage training process is finished, we remove the auxiliary decoder and only use the main decoder for model inference.

4.6. Prediction-based Planning

Trajectory prediction is usually followed by motion planning, which enables the generation of driving paths that comply with traffic regulations and safety requirements while ensuring ride comfort and traffic efficiency. To demonstrate the use of prediction results from Pioformer, we build a motion planner to solve the optimal control sequence of the ego vehicle based on the predicted trajectories of the ego vehicle and other vehicles, as shown in Fig. 8. With lane segments and observed trajectories, the multimodal trajectory prediction model, Pioformer, can generate the predicted trajectories for the ego vehicle and other vehicles. Similar to (Huang et al., 2023b), we convert the ego vehicle's predicted trajectory that has the highest probability into control sequences, which serve as an initial plan for the motion planner. Then, the initial plan of the ego vehicle, other vehicle-predicted trajectories and observed trajectories are used as the inputs of the optimization problem in the motion planner. The motion planner aims to find the optimal control sequence of the ego vehicle by minimizing the cost function presented in Fig. 8, including the ride comfort cost, the lane adherence cost and the safety cost. Finally, we transform the final plan (optimal control sequence) of the ego vehicle back into 2D positions for visualization, as shown in Fig. 11.

The objective function is defined as the sum of squared residual terms, where each term is a product of a weight ξ_i and cost C_i :

$$\mathbf{U}^* = \arg \min_{\mathbf{U}} \frac{1}{2} \sum_i \left\| \xi_i C_i(\mathbf{U}^i, \hat{\mathbf{k}}) \right\|^2, \quad (34)$$

where \mathbf{U}^* is the optimal control sequence, which includes acceleration and steering angle. $\hat{\mathbf{k}}$ denotes the predicted states of neighboring vehicles, C_i is a function of $\mathbf{U}^i \subset \mathbf{U}$ and $\hat{\mathbf{k}}$, where $\mathbf{U} = \{\mathbf{u}_t\}_{t=T_{\text{obs}}+1}^{T_{\text{pred}}}$. The ride comfort cost constrains the acceleration, jerk, steering angle, and steering change rate. The lane adherence cost constrains the distance and angle deviation from the lane center line. The safety cost constrains the safe distance from other traffic participants on the road to avoid collision. For more details, please refer to (Huang et al., 2023b).

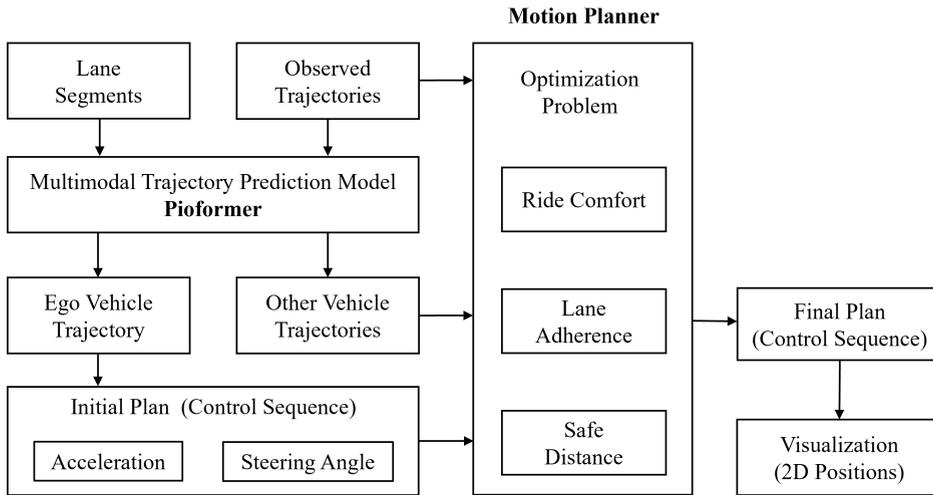


Figure 8: The pipeline of motion planning based on the prediction results. With lane segments and observed trajectory, our model Pioformer can generate the predicted trajectories for ego vehicle and other vehicles. Similar to (Huang et al., 2023b), we convert the ego vehicle's predicted trajectory into control sequences, which serve as an initial plan. Then, the initial plan of the ego vehicle, other vehicle-predicted trajectories and observed trajectories are used as the inputs of the motion planner. The motion planner aims to find the optimal control sequence of the ego vehicle by minimizing the cost function, including the ride comfort cost, the lane adherence cost and the safety cost. Finally, we transform the final plan (optimal control) of the ego vehicle back into 2D positions for visualization.

Table 1

Quantitative results of the comparative study on *Argoverse 1* (Chang et al., 2019) validation set and the online test set. †: The results are obtained using the checkpoints from the official implementation at <https://github.com/ZikangZhou/HiVT>. Ensemble-based methods and models with parameters over 5000K are in **bold**.

Model	Validation Set			Test Set				#Param ↓
	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓	brier-minFDE ₆ ↓	
LaneGCN (Liang et al., 2020)	0.71	1.08	-	0.8703	1.3622	0.1620	2.0539	3701K
LaneRCNN (Zeng et al., 2021)	0.77	1.19	0.08	0.9038	1.4526	0.1232	2.1470	-
HOME+GOME (Gilles et al., 2022)	-	-	-	0.8904	1.2919	0.0846	1.8601	Ensemble
DenseTNT (Gu et al., 2021)	0.73	1.05	0.10	0.8817	1.2815	0.1258	-	1103K
mmTransformer (Liu et al., 2021)	0.72	1.21	0.09	0.8436	1.3383	0.1540	2.0328	2607K
TPCN (Ye et al., 2021)	0.73	1.15	0.11	0.8153	1.2442	0.1333	1.9286	-
Multipath++ (Varadarajan et al., 2022)	-	-	-	0.7897	1.2144	0.1324	1.7932	Ensemble
Scene Transformer (Ngiam et al., 2022)	0.80	1.23	0.13	0.8026	1.2321	0.1255	1.8868	15296K
LTP (Wang et al., 2022)	0.78	1.07	-	0.8335	1.2955	0.1472	1.8566	1100K
HiVT-64† (Zhou et al., 2022)	0.69	1.03	0.10	0.8306	1.3053	0.1503	1.9736	662K
HiVT-128† (Zhou et al., 2022)	0.66	0.97	0.09	0.7993	1.2320	0.1369	1.9010	2529K
Wayformer (Nayakanti et al., 2022)	-	-	-	0.7676	1.1616	0.1186	1.7408	Ensemble
FRM (Park et al., 2023)	0.68	0.99	-	0.8165	1.2671	0.1430	1.9365	-
Macformer-S (Feng et al., 2023)	-	-	-	0.8490	1.2704	0.1311	1.9021	879K
Macformer-L (Feng et al., 2023)	-	-	-	0.8188	1.2160	0.1205	1.8275	2485 K
LAformer (Liu et al., 2024)	0.64	0.92	0.08	0.7720	1.1626	0.1245	1.8347	2,654K
ProphNet-S (Wang et al., 2023)	0.68	0.97	-	-	-	-	-	-
Pioformer (ours)	0.66	0.95	0.09	0.7939	1.1954	0.1283	1.8608	829K

Table 2

Quantitative results of Pioformer-L on *Argoverse 2* dataset (Wilson et al., 2023).

Method	brier-minFDE ₆ ↓	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓
HiVT-128	2.51	0.81	1.82	0.27
Pioformer-L	2.25	0.73	1.56	0.22

Table 3

Comparative results in scenarios with three levels of interaction intensity.

Interaction Level	Model	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓
Slight	HiVT-64	0.51	0.68	0.04
	HiVT-128	0.49	0.63	0.03
	Pioformer	0.49	0.62	0.03
Moderate	HiVT-64	0.63	0.86	0.07
	HiVT-128	0.62	0.83	0.06
	Pioformer	0.62	0.83	0.06
Strong	HiVT-64	0.95	1.73	0.23
	HiVT-128	0.90	1.60	0.21
	Pioformer	0.88	1.49	0.19

5. Experiments

This section introduces the implementation details, including the datasets, evaluation metrics, and baseline methods. We conduct comparative experiments, ablation studies, and sensitivity analysis to evaluate the effectiveness of our model.

5.1. Experimental Setup

Datasets. We use *Argoverse 1* Motion Forecasting Dataset to train and validate our motion forecasting model, similar to HiVT (Zhou et al., 2022). *Argoverse 1* is a public dataset (Chang et al., 2019), which includes 324,557 scenarios, each 5 seconds long, including observations for the past 2 seconds and predictions for the future 3 seconds.

The most challenging segments such as intersections, vehicles taking left or right turns, and changing lanes are selected. In *Argoverse 1*, each scenario contains the 2D, birds-eye-view (BEV) centroid of each object sampled at 10 Hz. We note that the BEV data are generated by projecting the 3D perception results from the upstream perception tasks onto a 2D plane. Except for the trajectory data, *Argoverse 1* also contains HD map information, and the semantic vector maps include lane-level detail, such as lane centerlines, traffic direction, and intersection annotations, which are fed into the proposed prediction model together with the trajectory data. In addition, to further evaluate the generalization ability of our method, we also conduct experiments on benchmark *Argoverse 2* (Wilson et al., 2023), which provides trajectory sequences with 5-second observation windows and 6-second prediction horizons.

Evaluation Metrics. We choose the following metrics for evaluations, *i.e.*, Minimum Average Displacement Error (min ADE_K), Minimum Final Displacement Error (min FDE_K), Miss Rate (MR_K), and the brier-min FDE_K metric (for more details, see (Liang et al., 2020; Gu et al., 2021; Zhou et al., 2022; Huang et al., 2023a)). K is the number of prediction modes, and we set K as 6 as requested in the Argoverse Motion Forecasting Competition.

- min ADE_K measures the average l_2 distances between the best-predicted trajectory and the ground truth.
- min FDE_K measures the l_2 distance between the endpoint of the best-predicted trajectory and the ground truth.
- MR_K computes the ratio of scenarios in which the endpoints of the best-predicted trajectory are located more than 2.0 meters from that of the ground truth.
- brier-min FDE_K adds $(1 - \hat{\pi})^2$ to the min FDE_K, where $\hat{\pi}$ corresponds to the probability of the best-predicted trajectory.

Training Details. The hyperparameters of our models are similar to HiVT in (Zhou et al., 2022). The training batch size is 32, and the initial learning rate is 5×10^{-4} . The loss weight term λ_1 is 1, and λ_2 is 5. The number of nodes S on each hyperedge is set to 4. All experiments are conducted on an RTX 4090 Ti GPU using AdamW optimizer (Loshchilov and Hutter, 2019). We use the cosine annealing scheduler (Ilya Loshchilov and Frank Hutter, 2017) for 64 epochs at each stage. For more details on network architectures and hyperparameters, see (Zhou et al., 2022).

5.2. Comparative Study

We compare our method with several state-of-the-art methods regarding the evaluation metrics. To further identify the effectiveness in modeling the interactions, we test our method with the baseline HiVT-64 and HiVT-128 in scenarios with different interaction levels.

5.2.1. Comparison with SOTA

The proposed model Pioformer is compared with several state-of-the-art methods regarding the evaluation metrics and the model size (the number of parameters), as shown in Table 1. Similarly, we have both a lightweight and a large version of Pioformer, *i.e.*, Pioformer, and Pioformer-L, which has 64 hidden units and 128 hidden units, respectively. The baseline approaches include HiVT (HiVT-64 and HiVT-128) (Zhou et al., 2022), the ensemble-based methods such as Wayformer (Nayakanti et al., 2022), multipath++ (Varadarajan et al., 2022), and HOME+GOME (Gilles et al., 2022), and other top-ranking methods on the leaderboard of Argoverse Motion Forecasting Competition over the years.

The ensemble-based method Wayformer (Nayakanti et al., 2022) and LAformer (Liu et al., 2024) have achieved slightly smaller prediction errors than our model, but the two have sacrificed their model sizes and the inference speed, as shown in Table 1. The ensemble-based methods often have a model size of more than 5000K, and LAformer (Liu et al., 2024) has a model size of 2654k, three times larger than ours. In contrast, our model size is 829k, and we have achieved a similar level of model accuracy to those large models, *i.e.*, the min ADE₆ of 0.7939, min FDE₆ of 1.1954, MR₆ of 0.1283, and brier-minFDE₆ of 1.8608. In particular, we have smaller prediction errors than the baseline HiVT on both the validation set and test set. Specifically, our model has further reduced the prediction error of HiVT-128 by nearly one-third of its model size. Also, we can achieve much higher accuracy than HiVT-64 without significantly increasing the number of parameters, *i.e.*, 662k vs. 829k.

We also conduct experiments on another dataset, *i.e.*, *Argoverse 2*, and the results between HiVT-128 and Pioformer-L are shown in Table 2. Compared with HiVT-128, the prediction errors of Pioformer-L have been substantially reduced regarding all evaluation metrics, and in particular, both brier-min FDE₆ and min FDE₆ decrease by 0.26, showing a very significant improvement in prediction accuracy. As the map encoder of HiVT-128 is not

specifically designed to accommodate map information in *Argoverse 2*, the performance of HiVT-128 is suboptimal, and our model shows better performance in *Argoverse 2*, indicating that our model can generalize better than HiVT. In particular, we note that the parameter size of Pioformer-L is larger than that of HiVT-128, while the small version of Pioformer can achieve a comparable prediction accuracy as the HiVT-128 with one-third of the parameter size. Therefore, we use the lightweight version of Pioformer as our final approach, which can achieve a good trade-off between the model accuracy and model size.

5.2.2. Comparison in different interaction levels

To validate the effectiveness of our model in capturing the high-order interaction features, we test the model in scenarios with three levels of interaction, *i.e.*, slight interaction, moderate interaction, and strong interaction, as shown in Table 3. In scenarios with slight interactions, there are few surrounding vehicles and no turning or lane-changing maneuvers around the central agent. In scenarios with moderate interactions, more vehicles surround the central agent, frequently at intersections, but still without turning maneuvers. In scenarios with strong interactions, the central agent is typically confronted with a dense traffic environment characterized by numerous surrounding vehicles engaging in turning maneuvers.

From Table 3, the prediction errors increase with the rise of the interaction intensity. Our model has achieved smaller prediction errors than HiVT-64 and HiVT-128 in three interaction scenarios. In strong-interactive scenarios, compared with HiVT-64, there exists a 0.07 decrease in $\min ADE_6$, 0.24 decrease in $\min FDE_6$, 0.04 decrease in MR_6 . Similarly, in contrast with HiVT-128, we have reduced 0.02 in $\min ADE_6$, 0.11 in $\min FDE_6$, and 0.02 in MR_6 , respectively. In particular, the gap between our model and HiVT has further enlarged in the strong-interactive scenarios, indicating the strengths of our model in capturing complex interaction behaviors.

5.2.3. Three-stage training vs. End-to-End training

As discussed previously, we execute three training stages sequentially while reloading weights from the previous stage. In stage 1, we train CTN to achieve coarse trajectory prediction. CTN and TPN are trained together to generate trajectory proposals in stage 2. In stage 3, we train CTN, TPN and PRN together to produce the final prediction. To evaluate the effects of the three-stage training scheme, we compare it with the end-to-end training process, as shown in Table 4. It can be seen that the three-stage training scheme can achieve higher prediction accuracy, *i.e.*, 0.05 decrease in $\min ADE_6$, 0.1 decrease in $\min FDE_6$, and 0.01 decrease in MR_6 , respectively. Therefore, the proposed three-stage training scheme is beneficial in improving the model's accuracy.

Table 4
Three-stage training versus end-to-end training.

Training Scheme	$\min ADE_6 \downarrow$	$\min FDE_6 \downarrow$	$MR_6 \downarrow$
End-to-end	0.71	1.05	0.10
Three-stage	0.66	0.95	0.09

Table 5
Ablation studies on Trajectory Proposal Network and Proposal Refinement Network in our model.

Trajectory Proposal Network	Proposal Refinement Network	$\min ADE_6 \downarrow$	$\min FDE_6 \downarrow$	$MR_6 \downarrow$
-	-	0.69	1.03	0.10
✓	-	0.67	0.98	0.10
-	✓	0.67	0.98	0.10
✓	✓	0.66	0.95	0.09

5.3. Ablation Study

We conduct ablation studies to investigate the effect of each part of the model, including the Trajectory Proposal Network (TPN), Proposal Refinement Network (PRN), Hypergraph Neural Network (HGNN), and the three-stage training scheme.

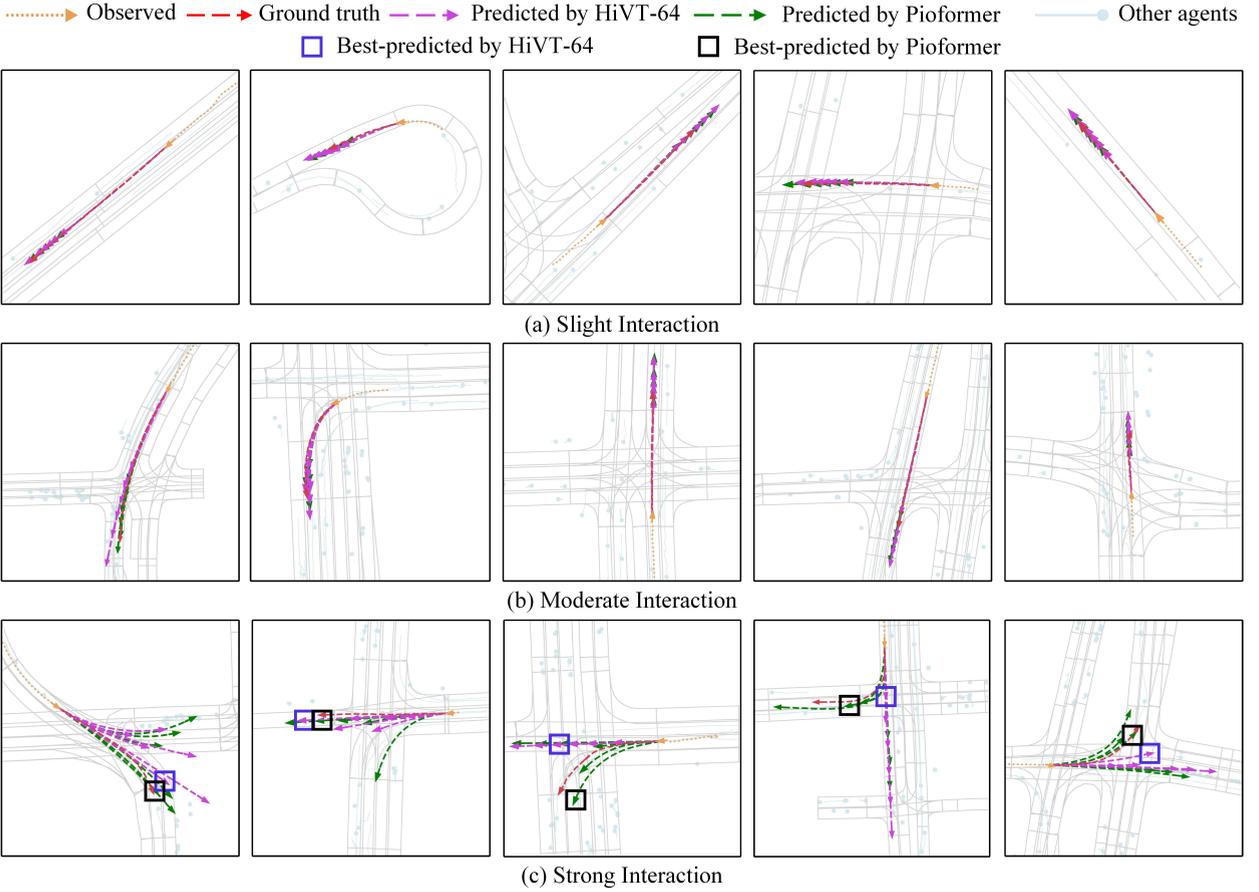


Figure 9: Qualitative comparison of multimodal trajectory predictions of Pioformer and HiVT-64. The best-predicted trajectory is compared with the ground truth trajectory (red) for evaluations. We observe that in scenarios with slight and moderate interactions, both HiVT and Pioformer can generate satisfying predictions, but Pioformer significantly outperforms HiVT-64 in strong-interactive scenarios. The best-predicted trajectory of Pioformer is closer to the ground truth than HiVT-64 (third row), demonstrating the superiority of our method in capturing the complex interactions among agents and the benefits of incorporating the post-interactions into the prediction network.

5.3.1. Trajectory Proposal Network and Proposal Refinement Network

We evaluate the effects of TPN and PRN in our model, as shown in Table 5. We observe that removing the TPN from our model can increase prediction errors, which indicates that the TPN is beneficial in improving accuracy. Similarly, all evaluation metrics increase when the PRN is removed from our model, indicating the benefits of the PRN. In particular, the prediction errors continue to increase when both TPN and PRN have been removed, i.e., 0.03 increase in min ADE₆, 0.08 increase in min FDE₆, and 0.01 increase in MR₆, respectively, which demonstrates the joint benefits of TPN and PRN to our model.

5.3.2. Effects of HGNN

We analyze the effectiveness of the Hypergraph Neural Network (HGNN) in the Hyper-Interactor module of our approach. For comparison, we replace the HGNN with standard GNN, and the prediction results are shown in Table 6. It can be seen that the biggest difference between HGNN and the standard GNN is shown in the strong-interactive scenarios, and HGNN has achieved much smaller prediction errors, i.e., 0.04 reduction in min ADE₆, 0.15 in min FDE₆ and 0.03 in MR₆, indicating that HGNN can better capture the complex high-order interactions than the standard GNN in our model.

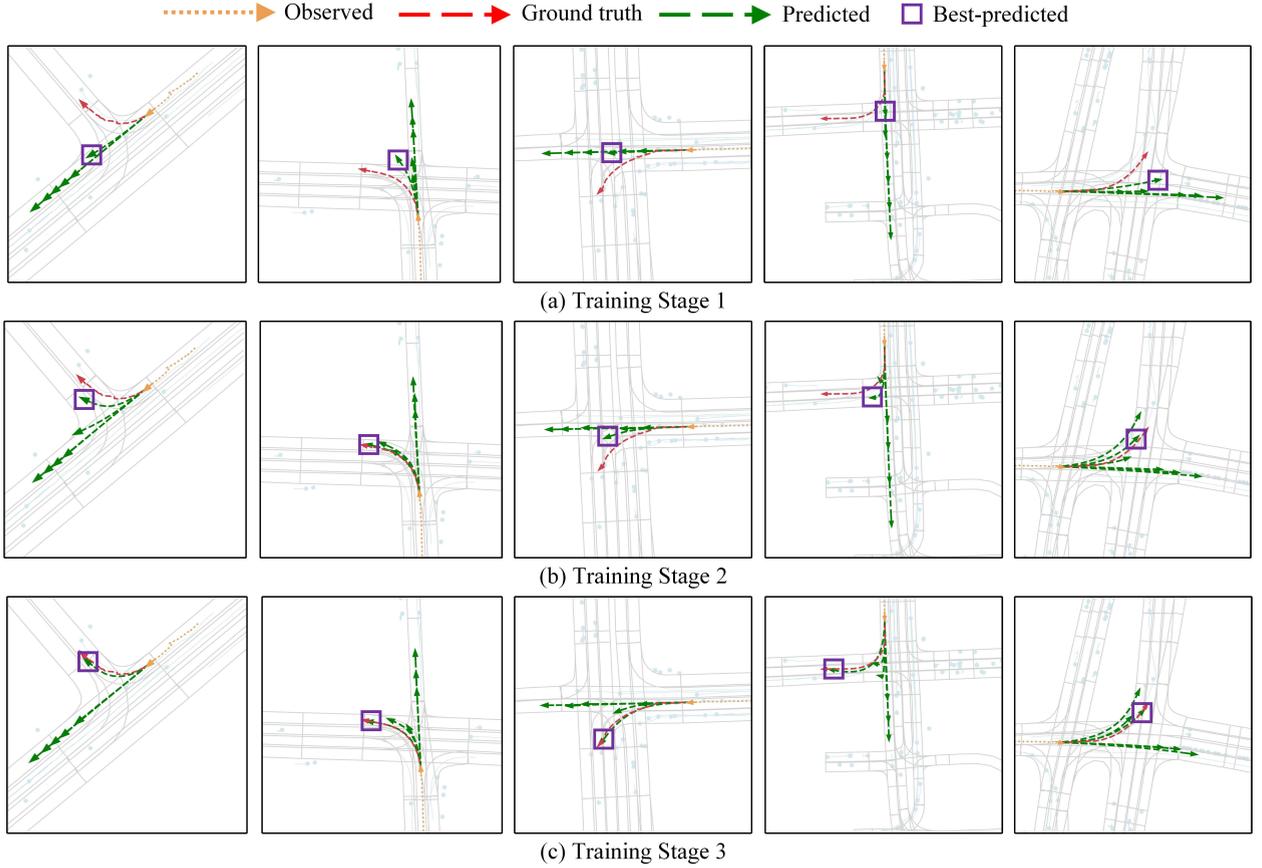


Figure 10: Qualitative comparison of multimodal predictions of Pioformer in three stages. Each column represents a scenario, and we compare the prediction performance of three stages in the same scenario. We observe that as the training stages progress, the predicted trajectories gradually converge toward the ground truths. The final model (training stage 3) demonstrates significantly improved performance compared to the CTN model (training stage 1).

Table 6

Comparison with GNN-based Hyper-Interactor on scenarios of different interaction intensity.

Interaction Level	Type of Hyper-Interactor	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓
Slight	GNN	0.50	0.65	0.03
	HGNN	0.49	0.62	0.03
Moderate	GNN	0.62	0.84	0.06
	HGNN	0.62	0.83	0.06
Strong	GNN	0.92	1.64	0.22
	HGNN	0.88	1.49	0.19

5.3.3. Three-stage Training Scheme

We compute the number of model parameters and prediction errors in three training stages, as shown in Table 7. It can be seen that the prediction errors decrease from stage 1 to stage 3 without significantly increasing the model size. Specifically, the model learns the underlying patterns of future trajectory prediction in stage 1, which are utilized in the learning of the post-interactions features in stage 2. The previous two stages generate trajectory proposals and interaction features, which are fed as inputs for stage 3, and the model accuracy can be further improved in the final stage. This demonstrates the benefits of the three-stage training scheme in improving the model's accuracy.

Table 7

Ablations on the three-stage training scheme for our model Pioformer on *Argoverse 1* test set (Chang et al., 2019). Best in **bold**.

Stage	#Param	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓	brier-minFDE ₆ ↓
1	654K	0.8306	1.3053	0.1503	1.9736
2	733K	0.8141	1.2683	0.1439	1.9398
3	829K	0.7939	1.1954	0.1283	1.8608

Table 8

Sensitivity study of the hyperparameters in the loss functions.

λ_1	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓	λ_2	minADE ₆ ↓	minFDE ₆ ↓	MR ₆ ↓
0.5	0.67	0.98	0.10	1	0.66	0.96	0.09
1	0.67	0.98	0.10	2	0.66	0.96	0.09
5	0.67	0.98	0.10	5	0.66	0.95	0.09
				10	0.67	0.95	0.09

5.4. Sensitivity Study

The objective of the sensitivity analysis is to evaluate the impact of the loss weights λ_1 and λ_2 on the model performance. λ_1 is used to balance the loss of the auxiliary decoder in the CTN and the loss of the main decoder in the TPN. λ_2 is used to balance the loss of CTN, TPN, and the loss of the final trajectory outputs in the PRN. We conducted an empirical study to evaluate their impact by varying their values around the experimental settings, and λ_1 was varied between 0.5, 1, 5, and λ_2 was varied between 1, 2, 5, and 10. As shown in Table 8, varying λ_1 does not impact the model performance, showing that λ_1 is not a sensitive parameter in our approach. In contrast, the model has the minimum prediction errors in minADE₆, and minFDE₆ when λ_2 is set as 5, indicating that λ_2 is a sensitive parameter. An appropriate value of λ_2 can balance well between the trajectory proposal loss and the refinement loss.

5.5. Qualitative Results

5.5.1. Comparison with HiVT

This study focuses on multimodal trajectory prediction, and the prediction results are a distribution of future possible trajectories. We visualize the prediction results of our model Pioformer and the baseline HiVT on different levels of interaction scenarios, as shown in Fig. 9. Our model Pioformer uses 64 hidden units, and thus, we chose the HiVT-64 for comparison. We observe that both Pioformer and HiVT-64 can generate trajectory predictions that match well with the ground truth trajectories in scenarios with slight and moderate interactions; see Fig. 9 (a) and (b). In contrast, achieving alignments with the ground truth trajectory, especially at intersections, is more challenging. Our model Pioformer has outperformed HiVT-64 in strong-interactive scenarios. The trajectories predicted by Pioformer are located within reasonable ranges, and the best-predicted trajectory (black box) of our model is closer to the ground truth trajectory (red dash line) than that of HiVT-64 (blue box), indicating the superiority of our model in capturing complex interaction behaviors, see Fig. 9 (c).

5.5.2. Three-stage Training Scheme

We visualize the prediction results of three training stages, as shown in Fig. 10. It can be seen that the performance of our model has increased from the stage 1 to the stage 3, with the best-predicted trajectories (marked with box) gradually getting close to the ground truth trajectory (red dash line). In the stage 1, the predicted trajectories often fail to match the correct lanes, resulting in larger displacement errors (see Fig. 10 (a)). In the stage 2, our model implicitly refines the coarse trajectories at the latent space, resulting in a closer approximation to the correct lane and direction (see Fig. 10 (b)). In stage 3, we achieve the final predicted trajectories based on the trajectory proposals from stage 2 and the integrated interaction features. We found that the best-predicted trajectory of the stage 3 has shown the minimum deviations from the ground truth trajectory, as shown in Fig. 10 (c). Overall, the final refinement process at the stage 3 helps bring the deviated trajectories back to the desired lanes and enables more accurate predictions, indicating the effectiveness of the three-stage training scheme.

5.5.3. Effects on Motion Planning

We illustrate how trajectory prediction can benefit motion planning in this study. As discussed in Section 4.6, motion predictions can be employed for motion planning and help improve the performance of motion planning. The prediction results (left column) and the motion planning results (right column) are compared for demonstration, as shown in Fig. 11. The crucial area is marked with red rectangles, within which the most related vehicle that is located right in front of the ego vehicle is marked in red point. It can be seen that our model Pioformer outperforms HiVT-64 in trajectory predictions, resulting in planned trajectories (right column) based on Pioformer being closer to the ground truths than those based on HiVT-64. Besides, the endpoint of the planned trajectory based on predictions from Pioformer can keep a safer distance from the critical vehicle than those in predictions. In other words, even if the trajectory endpoint predicted by the Pioformer is close to the critical vehicle, the planned trajectories can maintain a safe distance from it and are closer to the ground truth. This shows that trajectory predictions can bring benefits for motion planning and help planning modules avoid risk situations in advance. Accurately predicted trajectories result in safer planned trajectories that are closer to the ground truth, demonstrating the effectiveness of the prediction-based planning.

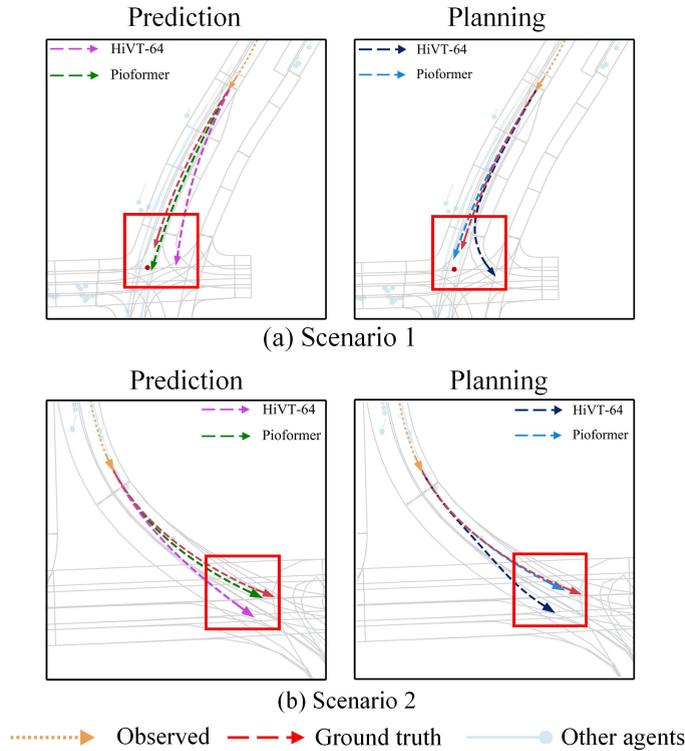


Figure 11: We illustrate how trajectory prediction can benefit motion planning. Prediction results (left) and motion planning based on the predictions (right). The crucial area is marked with red rectangles, within which the vehicle that is located in front of the ego vehicle is marked in red point. The endpoint of the planned trajectory (right column) based on the predictions from Pioformer can keep a safer distance from the critical vehicle than the predictions. On the left side, the prediction gets very close to the critical vehicle, while the planned trajectory on the right side is far from the vehicle. This shows that trajectory prediction can be used for motion planning and helps avoid critical situations in advance.

5.6. Discussion

This study demonstrates a coarse-to-refine multimodal trajectory prediction model, *i.e.*, Pioformer, which utilizes the high-order post-interaction features for motion predictions. Extensive experiments on *Argoverse 1* indicate that our model can greatly improve prediction accuracy without sacrificing the model size. Our study provides a new insight into the relationship between interaction modeling and motion predictions of multiple agents in the context of intelligent transportation systems. However, the proposed method still needs many improvements for real-world implementations. For instance, model latency is essential for the downstream planning task in real-world applications. There exists an inevitable increase in inference latency due to our use of HiVT as the CTN, and the reuse of its GI module in our model.

In scenarios with a large number of agents (e.g., 64), the average inference latencies for HiVT-64 and HiVT-128 are 8.4ms and 8.6ms, while Pioformer exhibits a latency of 10.8ms. Despite the slight increase in inference speed, the latency of Pioformer still meets the real-time use cases at 10 Hz. In our future work, we will focus on designing more streamlined network architectures to enhance the inference speed of Pioformer. In addition, the evaluation metrics $\min ADE_K$ and $\min FDE_K$ can cause information leaks when using the best-predicted trajectory for computation (Huang et al., 2023a). That is, the trajectory modes that have poor performances (e.g., violations of the traffic rules, large deviations from the ground truth) would be neglected, which is insufficient to provide a comprehensive assessment for the multimodal trajectory models. Instead, the probability-aware metrics can account for the poor modes in the multimodal predictions, which could be used in our future works. For instance, choosing the prediction with the highest probability, or multiple predictions with a probability higher than a threshold for evaluations may provide more useful information. Overall, evaluation metrics are vital in guiding the design of multimodal trajectory prediction models, which should be combined together for a comprehensive and unbiased evaluation.

6. Conclusion

In this study, we propose a post-interactive multimodal trajectory prediction model for autonomous driving, *i.e.*, Pioformer, which can capture the post-interaction features and leverage them to improve the accuracy of the predicted trajectories. Based on the CTN, which serves as a low-order interaction feature extractor and generates coarse trajectories, we propose an HGNN-based TPN to capture the high-order post-interaction features and refine future trajectories in the latent space, resulting in trajectory proposals. To further leverage post-interaction features and explore spatial-temporal consistency, we re-encode the trajectories in the entire horizon and generate both trajectory and confidence offsets to refine the proposals in TPN explicitly. Furthermore, we employ a three-stage training scheme to fully exploit the results of previous training stages, enabling a progressive feature extraction process that promotes stability during the training process. Extensive experiments underscore the superiority of our model to achieve an excellent balance between model size and accuracy. We also demonstrate the effectiveness of the essential components of Pioformer and highlight the necessity of the three-stage training scheme. Moreover, experiment results also indicate that our method can enhance accuracy and reduce the occurrence of unreasonable or hazardous situations in scenarios involving strong interactions.

7. Acknowledgments

We thank Dr. Wenzhao Zheng at UC Berkeley and Dr. Beihao Xia at Huazhong University of Science and Technology for the useful discussion on this study.

CRedit authorship contribution statement

Ziyi Huang: Data curation, Formal analysis, Investigation, Methodology, Software, Validation, Visualization, Writing - original draft. **Yang Li:** Investigation, Writing - review & editing, Supervision, Resources, Funding acquisition, Project administration, Conceptualization. **Dushuai Li:** Investigation, Formal analysis, Writing - review & editing. **Yao Mu:** Methodology, Formal analysis, Writing - review & editing. **Hongmao Qin:** Writing - review & editing, Supervision, Resources. **Nan Zheng:** Conceptualization, Writing - review & editing, Supervision, Resources.

References

- Alahi, A., Goel, K., Ramanathan, V., Robicquet, A., Fei-Fei, L., Savarese, S., 2016. Social lstm: Human trajectory prediction in crowded spaces, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 961–971.
- Cai, Z., Vasconcelos, N., 2018. Cascade r-cnn: Delving into high quality object detection, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.
- Chai, Y., Sapp, B., Bansal, M., Anguelov, D., 2020. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction, in: Kaelbling, L.P., Kragic, D., Sugiura, K. (Eds.), Proc. Conf. Robot Learn., PMLR. pp. 86–99.
- Chandra, R., Bhattacharya, U., Bera, A., Manocha, D., 2019. Traffic: Trajectory prediction in dense and heterogeneous traffic using weighted interactions, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 8483–8492. doi:10.1109/IEEE/CVF.2019.00868.
- Chang, M.F., Lambert, J., Sangkloy, P., Singh, J., Bak, S., Hartnett, A., Wang, D., Carr, P., Lucey, S., Ramanan, D., Hays, J., 2019. Argoverse: 3d tracking and forecasting with rich maps, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.
- Chung, J., Gulcehre, C., Cho, K., et al., 2014. Empirical evaluation of gated recurrent neural networks on sequence modeling. ArXiv:1412.3555 .
- Ding, K., Wang, J., Li, J., Li, D., Liu, H., 2020. Be more with less: Hypergraph attention networks for inductive text classification, in: Proc. Conf. Empir. Methods Nat. Lang. Process., pp. 4927–4936. doi:10.18653/v1/2020.emnlp-main.399.

- Fang, L., Jiang, Q., Shi, J., Zhou, B., 2020. Tpnnet: Trajectory proposal network for motion prediction. *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, 6796–6805.
- Feng, C., Zhou, H., Lin, H., Zhang, Z., Xu, Z., Zhang, C., Zhou, B., Shen, S., 2023. Macformer: Map-agent coupled transformer for real-time and robust trajectory prediction. *IEEE Robot. Autom. Lett.* 8, 6795–6802. doi:10.1109/LRA.2023.3311351.
- Feng, Y., You, H., Zhang, Z., Ji, R., Gao, Y., 2018. Hypergraph neural networks, in: *Proc. AAAI Conf. Artif. Intell.*
- Gao, J., Sun, C., Zhao, H., Shen, Y., Anguelov, D., Li, C., Schmid, C., 2020. Vectornet: Encoding hd maps and agent dynamics from vectorized representation, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 11525–11533.
- Gao, Y., Feng, Y., Ji, S., Ji, R., 2023. Hgmn+: General hypergraph neural networks. *IEEE Trans. Pattern Anal. Mach. Intell.* 45, 3181–3199. doi:10.1109/TPAMI.2022.3182052.
- Gilles, T., Sabatini, S., Tsishkou, D., Stanculescu, B., Moutarde, F., 2022. Gohome: Graph-oriented heatmap output for future motion estimation, in: *IEEE Int. Conf. Robot. Autom.*, pp. 9107–9114. doi:10.1109/ICRA46639.2022.9812253.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y., 2014. Generative adversarial nets, in: Ghahramani, Z., Welling, M., Cortes, C., Lawrence, N., Weinberger, K. (Eds.), *Adv. Neural Inf. Process. Syst.*, Curran Associates, Inc.
- Gu, J., Sun, C., Zhao, H., 2021. Densentn: End-to-end trajectory prediction from dense goal sets, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, pp. 15283–15292. doi:10.1109/ICCV48922.2021.01502.
- Gu, T., Chen, G., Li, J., Lin, C., Rao, Y., Zhou, J., Lu, J., 2022. Stochastic trajectory prediction via motion indeterminacy diffusion, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 17113–17122.
- Gupta, A., Johnson, J., Fei-Fei, L., Savarese, S., Alahi, A., 2018a. Social gan: Socially acceptable trajectories with generative adversarial networks, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 2255–2264.
- Gupta, A., Johnson, J., Fei-Fei, L., Savarese, S., Alahi, A., 2018b. Social gan: Socially acceptable trajectories with generative adversarial networks, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 2255–2264. doi:10.1109/CVPR.2018.00240.
- Han, Y., Wang, P., Kundu, S., Ding, Y., Wang, Z., 2023. Vision hgmn: An image is more than a graph of nodes, in: *Proc. IEEE/CVF Int. Conf. Comput. Vis.*, pp. 19878–19888.
- Helbing, D., Molnar, P., 1995. Social force model for pedestrian dynamics. *Phys. Rev. E* 51, 4282.
- Ho, J., Jain, A., Abbeel, P., 2020. Denoising diffusion probabilistic models, in: Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., Lin, H. (Eds.), *Adv. Neural Inf. Process. Syst.*, Curran Associates, Inc., pp. 6840–6851.
- Huang, R., Xue, H., Pagnucco, M., Salim, F., Song, Y., 2023a. Multimodal trajectory prediction: A survey. *ArXiv:2302.10463*.
- Huang, Y., Du, J., Yang, Z., Zhou, Z., Zhang, L., Chen, H., 2022. A survey on trajectory-prediction methods for autonomous driving. *IEEE Trans. Intell. Veh.* 7, 652–674. doi:10.1109/TIV.2022.3167103.
- Huang, Z., Liu, H., Wu, J., Lv, C., 2023b. Differentiable integrated motion prediction and planning with learnable cost function for autonomous driving. *IEEE Trans. Neural Networks Learn. Syst.*, 1–15doi:10.1109/TNNLS.2023.3283542.
- Ilya Loshchilov and Frank Hutter, 2017. Sgdr: Stochastic gradient descent with warm restarts. *ArXiv:1608.03983 arXiv:1608.03983*.
- Jia, X., Wu, P., Chen, L., Liu, Y., Li, H., Yan, J., 2023. Hdgt: Heterogeneous driving graph transformer for multi-agent trajectory prediction via scene encoding. *IEEE Trans. Pattern Anal. Mach. Intell.*
- Jiang, B., Zhang, Z., Lin, D., Tang, J., Luo, B., 2019. Semi-supervised learning with graph learning-convolutional networks, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 11305–11312. doi:10.1109/CVPR.2019.01157.
- Jiao, R., Wang, Y., Liu, X., Huang, C., Zhu, Q., 2023. Kinematics-aware trajectory generation and prediction with latent stochastic differential modeling. *ArXiv:2309.09317*.
- Kim, E.S., Kang, W.Y., On, K.W., Heo, Y.J., Zhang, B.T., 2020. Hypergraph attention networks for multimodal learning, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 14569–14578. doi:10.1109/CVPR42600.2020.01459.
- Kossale, Y., Airaj, M., Darouichi, A., 2022. Mode collapse in generative adversarial networks: An overview, in: *Int. Conf. Optim. Appl.*, pp. 1–6. doi:10.1109/ICOA55659.2022.9934291.
- Lee, N., Choi, W., Vernaza, P., Choy, C.B., Torr, P.H.S., Chandraker, M., 2017. Desire: Distant future prediction in dynamic scenes with interacting agents, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 2165–2174. doi:10.1109/CVPR.2017.233.
- Li, D., Zhang, Q., Xia, Z., Zheng, Y., Zhang, K., Yi, M., Jin, W., Zhao, D., 2023. Planning-inspired hierarchical trajectory prediction via lateral-longitudinal decomposition for autonomous driving. *IEEE Trans. Intell. Veh.*, 1–12doi:10.1109/TIV.2023.3307116.
- Li, Y., Lu, X.Y., Wang, J., Li, K., 2020. Pedestrian trajectory prediction combining probabilistic reasoning and sequence learning. *IEEE Trans. Intell. Veh.* 5, 461–474. doi:10.1109/TIV.2020.2966117.
- Liang, M., Yang, B., Hu, R., Chen, Y., Liao, R., Feng, S., Urtasun, R., 2020. Learning lane graph representations for motion forecasting, in: *Proc. Eur. Conf. Comput. Vis.*, pp. 541–556.
- Liang, R., Li, Y., Zhou, J., Li, X., 2023. Stglow: A flow-based generative framework with dual-graphormer for pedestrian trajectory prediction. *IEEE Trans. Neural Networks Learn. Syst.*, 1–14doi:10.1109/TNNLS.2023.3294998.
- Lin, G., Liu, F., Milan, A., Shen, C., Reid, I., 2020. Refinetnet: Multi-path refinement networks for dense prediction. *IEEE Trans. Pattern Anal. Mach. Intell.* 42, 1228–1242. doi:10.1109/TPAMI.2019.2893630.
- Liu, M., Cheng, H., Chen, L., Broszio, H., Li, J., Zhao, R., Sester, M., Yang, M.Y., 2024. Laformer: Trajectory prediction for autonomous driving with lane-aware scene constraints, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 2039–2049.
- Liu, Y., Zhang, J., Fang, L., Jiang, Q., Zhou, B., 2021. Multimodal motion prediction with stacked transformers, in: *Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit.*, pp. 7573–7582. doi:10.1109/CVPR46437.2021.00749.
- Loshchilov, I., Hutter, F., 2019. Decoupled weight decay regularization, in: *Proc. Int. Conf. Learn. Represent.*
- Luo, H., E, H., Yang, Y., Guo, Y., Sun, M., Yao, T., Tang, Z., Wan, K., Song, M., Lin, W., 2023. HAHE: Hierarchical attention for hyper-relational knowledge graphs in global and local level, in: *Proc. Annu. Meet. Assoc. Comput. Linguist.*, pp. 8095–8107. doi:10.18653/v1/2023.acl-long.450.

- Ma, Y., Zhu, X., Zhang, S., Yang, R., Wang, W., Manocha, D., 2019. Trafficpredict: Trajectory prediction for heterogeneous traffic-agents, in: Proc. AAAI Conf. Artif. Intell., pp. 6120–6127.
- Mohamed, A., Qian, K., Elhoseiny, M., Claudel, C., 2020. Social-stgcn: A social spatio-temporal graph convolutional neural network for human trajectory prediction, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 14424–14432.
- Nayakanti, N., Al-Rfou, R., Zhou, A., Goel, K., Refaat, K.S., Sapp, B., 2022. Wayformer: Motion forecasting via simple & efficient attention networks. ArXiv:2207.05844 arXiv:2207.05844.
- Ngiam, J., Caine, B., Vasudevan, V., Zhang, Z., Chiang, H.T.L., Ling, J., Roelofs, R., Bewley, A., Liu, C., Venugopal, A., Weiss, D., Sapp, B., Chen, Z., Shlens, J., 2022. Scene transformer: A unified architecture for predicting multiple agent trajectories. ArXiv:2106.08417 arXiv:2106.08417.
- Papamakarios, G., Nalisnick, E.T., Rezende, D.J., Mohamed, S., Lakshminarayanan, B., 2019. Normalizing flows for probabilistic modeling and inference. J. Mach. Learn. Res. 22, 57:1–57:64.
- Park, D.H., Ryu, H., Yang, Y., Cho, J., Kim, J., Yoon, K.J., 2023. Leveraging future relationship reasoning for vehicle trajectory prediction. ArXiv:2305.14715 .
- Shi, S., Jiang, L., Dai, D., Schiele, B., 2022a. Motion transformer with global intention localization and local movement refinement, in: Adv. Neural Inf. Process. Syst., Curran Associates, Inc., pp. 6531–6543.
- Shi, S., Jiang, L., Dai, D., et al., 2022b. Mtr-a: 1st place solution for 2022 waymo open dataset challenge - motion prediction. ArXiv:2209.10033 .
- Sohn, K., Lee, H., Yan, X., 2015. Learning structured output representation using deep conditional generative models, in: Cortes, C., Lawrence, N., Lee, D., Sugiyama, M., Garnett, R. (Eds.), Adv. Neural Inf. Process. Syst., Curran Associates, Inc.
- Varadarajan, B., Hefny, A., Srivastava, A., Refaat, K.S., Nayakanti, N., Cornman, A., Chen, K., Douillard, B., Lam, C.P., Anguelov, D., Sapp, B., 2022. Multipath++: Efficient information fusion and trajectory aggregation for behavior prediction, in: IEEE Int. Conf. Robot. Autom., pp. 7814–7821. doi:10.1109/ICRA46639.2022.9812107.
- Velickovic, P., Cucurull, G., Casanova, A., Romero, A., Lio, P., Bengio, Y., 2017. Graph attention networks. ArXiv:1710.10903 .
- Wang, J., Ye, T., Gu, Z., Chen, J., 2022. Ltp: Lane-based trajectory prediction for autonomous driving, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 17113–17121. doi:10.1109/CVPR52688.2022.01662.
- Wang, X., Su, T., Da, F., Yang, X., 2023. Prophnet: Efficient agent-centric motion forecasting with anchor-informed proposals, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 21995–22003. doi:10.1109/CVPR52729.2023.02106.
- Wilson, B., Qi, W., Agarwal, T., Lambert, J., Singh, J., Khandelwal, S., Pan, B., Kumar, R., Hartnett, A., Pontes, J.K., Ramanan, D., Carr, P., Hays, J., 2023. Argoverse 2: Next generation datasets for self-driving perception and forecasting. ArXiv:2301.00493 arXiv:2301.00493.
- Xu, C., Li, M., Ni, Z., Zhang, Y., Chen, S., 2022a. Groupnet: Multiscale hypergraph neural networks for trajectory prediction with relational reasoning, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 6488–6497. doi:10.1109/CVPR52688.2022.00639.
- Xu, C., Wei, Y., Tang, B., Yin, S., Zhang, Y., Chen, S., 2022b. Dynamic-group-aware networks for multi-agent trajectory prediction with relational reasoning. ArXiv:2206.13114 arXiv:2206.13114.
- Xu, P., Hayet, J.B., Karamouzas, I., 2022c. Socialvae: Human trajectory prediction using timewise latents, in: Proc. IEEE/CVF Eur. Conf. Comput. Vis., Springer. pp. 511–528. doi:10.1007/978-3-031-19772-7_30.
- Ye, L., Wang, Z., Chen, X., Wang, J., Wu, K., Lu, K., 2022. Gsan: Graph self-attention network for learning spatial-temporal interaction representation in autonomous driving. IEEE Internet Things J. 9, 9190–9204.
- Ye, M., Cao, T., Chen, Q., 2021. Tpcn: Temporal point cloud networks for motion forecasting, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 11313–11322. doi:10.1109/CVPR46437.2021.01116.
- Zeng, W., Liang, M., Liao, R., Urtasun, R., 2021. Lanercnn: Distributed representations for graph-centric motion forecasting, in: Proc. IEEE/RSJ Int. Conf. Intell. Rob. Syst., pp. 532–539.
- Zhao, H., Gao, J., Lan, T., Sun, C., Sapp, B., Varadarajan, B., Shen, Y., Shen, Y., Chai, Y., Schmid, C., Li, C., Anguelov, D., 2021. Tnt: Target-driven trajectory prediction, in: Proc. Conf. Robot Learn., PMLR. pp. 895–904.
- Zhou, Z., Ye, L., Wang, J., Wu, K., Lu, K., 2022. Hivt: Hierarchical vector transformer for multi-agent motion prediction, in: Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit., pp. 8813–8823. doi:10.1109/CVPR52688.2022.00862.