

LEVERAGING FUTURE RELATIONSHIP REASONING FOR VEHICLE TRAJECTORY PREDICTION

Daehee Park¹, Hobin Ryu², Yunseo Yang¹, Jegyeong Cho¹, Jiwon Kim², Kuk-Jin Yoon¹

¹Korea Advanced Institute of Science and Technology ²NAVER LABS

{bag2824, acorn, j2k0618, kjyoon}@kaist.ac.kr

{hobin.ryu, g1.kim}@naverlabs.com

ABSTRACT

Understanding the interaction between multiple agents is crucial for realistic vehicle trajectory prediction. Existing methods have attempted to infer the interaction from the observed past trajectories of agents using pooling, attention, or graph-based methods, which rely on a deterministic approach. However, these methods can fail under complex road structures, as they cannot predict various interactions that may occur in the future. In this paper, we propose a novel approach that uses lane information to predict a stochastic future relationship among agents. To obtain a coarse future motion of agents, our method first predicts the probability of lane-level waypoint occupancy of vehicles. We then utilize the temporal probability of passing adjacent lanes for each agent pair, assuming that agents passing adjacent lanes will highly interact. We also model the interaction using a probabilistic distribution, which allows for multiple possible future interactions. The distribution is learned from the posterior distribution of interaction obtained from ground truth future trajectories. We validate our method on popular trajectory prediction datasets: nuScenes and Argoverse. The results show that the proposed method brings remarkable performance gain in prediction accuracy, and achieves state-of-the-art performance in long-term prediction benchmark dataset.

1 INTRODUCTION

For safe autonomous driving, predicting a vehicle’s future trajectory is crucial. Early heuristic prediction models utilized only the past trajectory of the target vehicle (Lin et al. (2000); Barth & Franke (2008)). However, with the advent of deep learning, more accurate predictions can be made by also considering the vehicle’s relationship with the High-Definition (HD) map (Liang et al. (2020); Zeng et al. (2021)) or surrounding agents (Lee et al. (2017); Chandra et al. (2020)). Since surrounding vehicles are not stationary, predicting relationships with them is much more complicated and has become essential for realistic trajectory prediction. Furthermore, since individual drivers control each vehicle, their interaction has a stochastic nature.

Previous works modeled interaction from past trajectories of the surrounding vehicles by employing pooling, multi-head attention, or spatio-temporal graph methods. However, we observed that these methods easily fail under complex road structures. For example, Fig. 1 shows the past trajectories of agents (left) and the attention weights among agents (right) obtained by a previous method (Mercat et al. (2020)) that learned the interaction among agents using multi-head attention (MHA). Since agents 0 and 4 are expected to join in the future, the attention weight between them should be high. However, the model predicts a low attention weight between them, highlighting the difficulty of reasoning future relationships between agents based solely on past trajectories. Incorporating the road structure should make the reasoning process much easier.

The decision-making process of human drivers can provide insights on how to model interaction. They first set their goal where they are trying to reach on the map. Next, to infer the interaction with surrounding agents, they roughly infer how the others will behave *in the future*. After that, they infer the interaction with others by inferring how likely the future path of other vehicles will overlap the path set by themselves. The drivers consider interaction more significant the more the future paths of other vehicles overlap with their own. We define the interaction from this process as a "*Future Relationship*". We use the following approaches to model Future Relationship, as shown in Fig. 3.

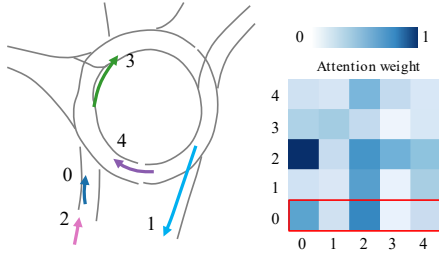


Figure 1: Past trajectories and corresponding attention map between agents from previous work (Mercat et al. (2020)). A weak relationship is inferred between agents that will highly interact in the future: agents 0 and 4.

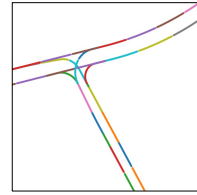


Figure 2: Lane segments represented in different colors.

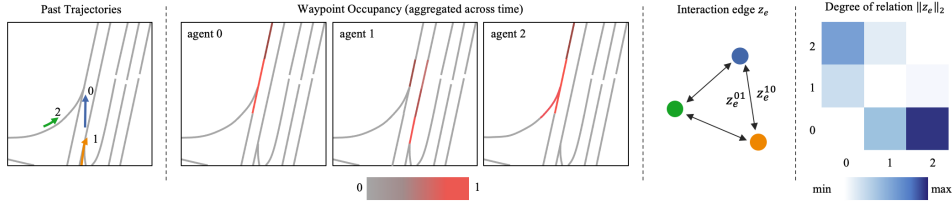


Figure 3: Key concept of the proposed method: From past observed trajectories, we predict the lane that a vehicle will pass in the future. The Interaction between agents is represented by an edge connecting their nodes, and is determined by the probability that two agents will pass adjacent lanes. The greater the probability, the higher the expected interaction.

First, we obtain the rough future motion of all vehicles in the scene. Since vehicles mainly move along lanes, we utilize lane information as strong prior for representing the rough future motion of vehicles. Because lane centerlines contain both positional and directional information, rough future motion can be represented as waypoint occupancy. The waypoint occupancy is defined as the probability of a vehicle passing a specific lane segment at every intermediate timestep. In the middle of Fig. 3, each agent’s waypoint occupancy is shown. Here, we aggregated the temporal axis for simplification. The probability that the vehicle passes that lane during the prediction horizon is drawn using the tone of red color.

Second, based on the waypoint occupancy, we infer the Future Relationship in probabilistic distribution. In most vehicle trajectory prediction methods, the interaction between agents is still made in a deterministic manner. However, we take note that interaction between vehicles is highly stochastic, and there can be multiple possible interactions. The deterministic relation inference averages out diverse interactions, interrupting socially-aware trajectory prediction. Therefore, we define Future Relationship in Gaussian Mixture (GM) distribution. Motivated by Neural Relational Inference (NRI) (Kipf et al. (2018)), we propose a method to train the diverse interaction distribution explicitly.

In summary, our contributions are:

- 1) We propose a new approach for modeling the interaction between vehicles by incorporating the road structure and defining it as *Future Relationship*.
- 2) We propose to infer the Future Relationship in probabilistic distribution using Gaussian Mixture (GM) distribution to capture diverse interaction.
- 3) The proposed method is validated on popular real-world vehicle trajectory datasets: nuScenes and Argoverse. In both datasets, there is a remarkable improvement in prediction performance and state-of-the-art performance is achieved in the long-range prediction dataset, nuScenes.

2 RELATED WORK

2.1 GOAL-CONDITIONED TRAJECTORY PREDICTION

Predicting future trajectories at once is a challenging task. Instead, goal-conditional prediction, which samples goal candidates and then predicts trajectory conditioned on them, is helpful and has shown

state-of-the-art performance, especially in long-range prediction tasks (Zhao et al. (2021); Gu et al. (2021); Phan-Minh et al. (2020); Chai et al. (2020); Zhang et al. (2021)). CoverNet (Phan-Minh et al. (2020)) and MultiPath (Chai et al. (2020)), which quantize the trajectory space to a set of anchors, often generate map-agnostic trajectories that cross non-drivable areas because the surrounding map is not considered. Recently, several studies have exploited the map information to obtain more performant goal candidates based on the assumption that vehicles follow lanes. TNT (Zhao et al. (2021)) uses goal points sampled from a lane centerline, and GoalNet (Zhang et al. (2021)) uses lane segments as trajectory anchors. However, while previous methods assume that the likelihood of arriving at a final destination is random, they assume that trajectories are unimodal in order to reach a specific goal area. In this paper, we assume inherent uncertainty in which trajectories can vary due to the interactions with surrounding vehicles in order to reach a specific goal area.

2.2 INTERACTION MODELING

Considering the interaction between agents helps to predict a socially aware trajectory. In the very early stage, interaction is obtained by pooling interaction features in the local region (Deo & Trivedi (2018); Gupta et al. (2018)). In other works, researchers attempted to obtain interaction through attention-based (Ngiam et al. (2022); Mercat et al. (2020); Vemula et al. (2018)) or GNN-based method (Carrasco et al. (2021); Cao et al. (2021); Zeng et al. (2021); Casas et al. (2020); Liang et al. (2020); Gao et al. (2020)). However, in most previous methods, interactions between agents are learned only with regression loss, which is insufficient to represent dynamic and rapidly changing situations. There exists a line of works that employs Neural Relational Inference (NRI) (Kipf et al. (2018)) that explicitly predicts and learns interaction using a latent interaction graph. EvolveGraph (Li et al. (2020)) utilizes two interaction graphs, static and dynamic, and NRI-MPM (Chen et al. (2021)) uses a relation interaction mechanism and spatio-temporal message passing mechanism. Similarly, we apply the NRI-based method to predict and train the interaction explicitly.

2.3 MULTI-MODAL TRAJECTORY PREDICTION

Trajectory prediction is a stochastic problem, which means that there are multiple possible futures instead of a unique answer. Recently, deep generative models like GAN (Goodfellow et al. (2014)) or VAE (Kingma & Welling (2013)) have been employed to address this issue. GAN-based (Gupta et al. (2018); Kosaraju et al. (2019); Li et al. (2021b)) and VAE-based models (Ivanovic & Pavone (2019); Salzman et al. (2020); Tang & Salakhutdinov (2019)) predict multiple futures by sampling multiple latent vectors. A well-organized latent space is necessary to sample meaningful latent vectors for predicting diverse, yet plausible future trajectories. This has become a natural choice in recent works (Ma et al. (2021); Bae et al. (2022)). The work most closely related to ours is GRIN (Li et al. (2021a)), which argues that multi-modality in trajectory prediction comes from two sources: personal intention and social relations with other agents. However, GRIN only considers past interaction, while we propose to consider future interaction by taking into account the characteristics of vehicle motion. Since vehicle motion mainly follows lanes, we utilize lane information to infer future interactions.

3 FORMULATION

In each scene, the past and future trajectories of N vehicles are observed. The past trajectory \mathbf{x}_t^- consists of positions for $-t_p : 0$ timesteps before the current timestep, and the future trajectory \mathbf{x}_t^+ consists of positions for $1 : t_f$ timesteps after the current timestep ($t = 0$). Lane information is obtained from the HD map, which consists of M segmented lane polylines. The lane information is represented as a graph: $\mathcal{G} = (\ell, \mathbf{e})$, where the nodes (ℓ) correspond to the different lane segments, and the edges (\mathbf{e}) represent the relationships between the segments. There are five relationship between segments: *predecessor*, *successor*, *left/right neighbor* and *in-same-intersection*. The input to the model is denoted as \mathbf{X} , which consists of the past and future trajectories of the vehicles and the lane information. Here, the future trajectories is only used in training. The output of the model is denoted as \mathbf{Y} , which consists of F predicted future trajectories for each agent. The model also predicts the future lane occupancy (i.e., which vehicles are occupying which lanes) as a medium using a probability distribution τ_t for each vehicle and lane segment at each future timestep: $1 : t_f$. The predicted future lane occupancy is denoted as τ_t^- , and the ground truth future lane occupancy is denoted as τ_t^+ .

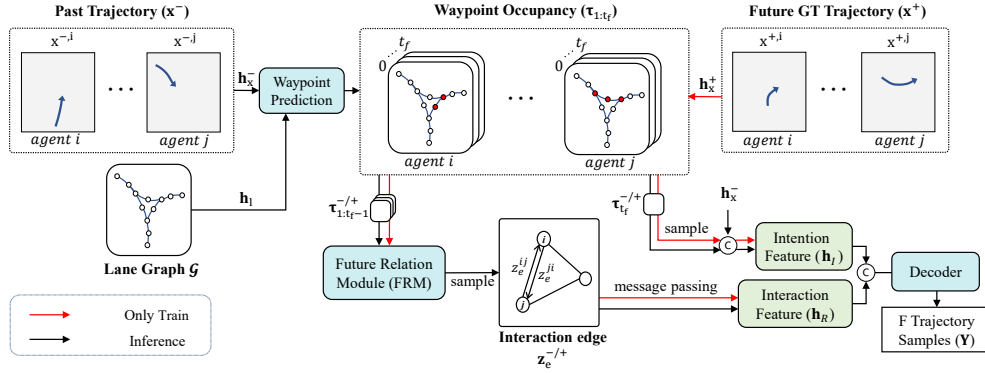


Figure 4: Overall structure of the proposed method. Given past/future motion inputs, the waypoint occupancy ($\tau_{1:t_f}$) is obtained. The goal features are then sampled following τ_{t_f} . Intention feature is derived from the goal features and the past motion (h_x^-). The Future Relationship Module (FRM) utilizes the intermediate waypoint occupancy ($\tau_{1:t_f-1}$) to sample the interaction edges among agents. Message passing is then performed to obtain the interaction feature. Finally, the decoder predicts F future trajectories from concatenation of intention and interaction features.

4 METHOD

Our focus is on modeling the "Future Relationship" between agents. A naive method to infer the Future Relationship is to predict all future vehicle trajectories and then calculate similarity among them. However, this method is inefficient and redundant, as it requires performing prediction twice. Moreover, the criteria for calculating similarity between trajectories may not be clear. In this paper, we utilize lane information for modeling the Future Relationship, inspired by the idea that the vehicles mainly follow lanes. Our key idea is that *if two vehicles are expected to pass on adjacent lanes, they will have a high chance of interacting in the future.*

We present the overall structure of our method in Fig. 4. First, we predict the waypoint occupancy, which represents the probability of a vehicle passing a specific lane segment during future time steps $1 : t_f$ (Sec. 4.1). Using this information, our Future Relationship Module (FRM) infers interaction as an edge feature connecting agent node pairs (Sec. 4.2). These interaction edges are used to transfer information between agent nodes to form the interaction feature through message passing. Finally, in the decoding stage (Sec. 4.3), the decoder predicts future trajectories from the aggregation of the interaction feature and intention feature, which is derived from the concatenation of past motion and goal features. Following AgentFormer (Yuan et al. (2021)), our method is based on CVAEs where the condition corresponds to the intention, and the latent code corresponds to the interaction feature. Then, we compute prior and posterior distribution of the interaction feature, as described in Sec. 4.2.2, 4.2.3.

4.1 WAYPOINT OCCUPANCY

In this section, we describe how to obtain the waypoint occupancy. We need two waypoint occupancies: one predicted from past trajectory (τ_t^-) and the ground truth (τ_t^+) for obtaining prior and posterior distributions of interaction, respectively.

To predict the waypoint occupancy from past trajectory, we first encode the past trajectories x^- and the lane graph \mathcal{G} into past motion and lane features: h_x^-, h_ℓ . Then, following TNT (Zhao et al. (2021)), we predict the waypoint occupancy as Eq. (1). Here, $[\cdot, \cdot]$ denotes concatenation, and we apply softmax to ensure that the waypoint occupancy sum up to one: $\sum^M \tau_t^m = 1$.

$$\tau_{1:t_f}^- = \text{softmax}(\text{MLP}([\mathbf{h}_x^-, \mathbf{h}_\ell])) \in \mathbb{R}^{N \times M \times t_f} \quad (1)$$

For GT waypoint occupancy, we can directly obtain it from GT future trajectory since we know the position and heading of vehicles. More details can be found in the supplementary material.

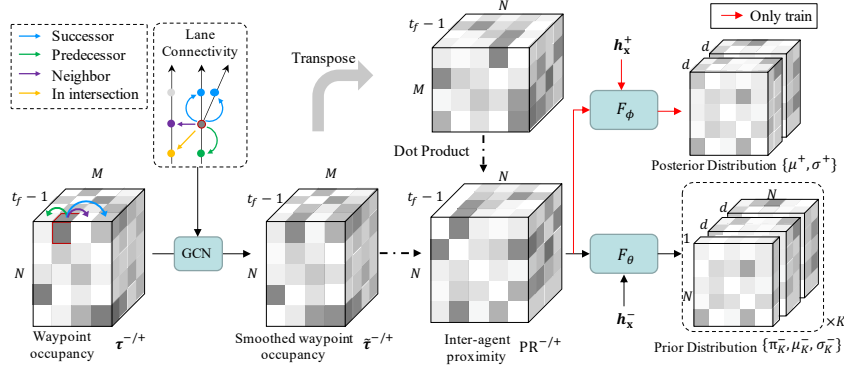


Figure 5: Future Relationship Module. During inference, predicted waypoint occupancy (τ^-) is fed to GCN, dot-producted by itself to obtain inter-agent proximity (PR^-). Prior of interaction (π_K, μ_K, σ_K) is then obtained as Gaussign Mixture. During training, GT waypoint occupancy (τ^+) is fed to obtain posterior of interaction (μ, σ) as Gaussian distribution.

4.2 FUTURE RELATIONSHIP MODULE (FRM)

Fig. 5 shows the FRM, which consists of three parts: computing inter-agent proximity and obtaining posterior and prior distribution. From intermediate waypoint occupancy of vehicle ($\tau_{1:t_f-1}$), we compute how each pair of vehicle pass adjacent lanes adjacent to each other at each timestep (inter-agent proximity). Based on that information and agents’ past motion features, we obtain two distribution of interaction. In the following sections, we describe the details of each part.

4.2.1 INTER-AGENT PROXIMITY

To compute the inter-agent proximity (PR), we first smooth the waypoint occupancy using a Graph Convolutional Network (GCN) (Welling & Kipf (2016)). The reason for doing so is that when a vehicle passes a specific lane, it affects other vehicles that pass the adjacent lane, not necessarily the same lane. Therefore, we apply different smoothing for each lane connectivity (predecessor, successor, neighbor, in-the-same-intersection) by employing 2-hop GCN layers. Specifically, each layer aggregated information from neighboring lanes and applies a non-linear transformation. This allows the model to capture spatial dependencies among agents and improve the accuracy of the inter-agent proximity computation. Each layer is expressed as Eq.(2) where σ , D_e , A_e and W_e are softmax followed by ReLU, degree, adjacency and weight matrix for each edge type, respectively.

$$\tilde{\tau}_{1:t_f-1} = \sum_{e \in \{succ, pred, right, left, inter\}} \sigma(D_e^{-1} A_e \tau_{1:t_f-1} W_e) \in \mathbb{R}^{N \times M \times (t_f-1)} \quad (2)$$

With this smoothed waypoint occupancy, we can compute the inter-agent proximity using the dot product of $\tilde{\tau}_{1:t_f-1}$ across the lane axis.

$$PR = \tilde{\tau}_{1:t_f-1} \cdot (\tilde{\tau}_{1:t_f-1})^T \in \mathbb{R}^{N \times N \times (t_f-1)} \quad (3)$$

4.2.2 PRIOR OF THE INTERACTION

To obtain the prior distribution, we use the past motion features (h_x^-) and inter-agent proximity. In this subsection, we omit superscript - for simplification. There are two design factors for our interaction modeling: (i) interaction should reflect diverse and stochastic properties, and (ii) it occurs in every pair of vehicles. Consequently, the prior distribution is defined as Gaussian Mixture (GM) per agent pair. Then, we define interaction edge e^{ij} between agent i and j as a d -dimensional feature ($p_\theta(e|\mathbf{X}) \sim \prod_{k=1}^K \pi_k \mathcal{N}(\mu_k, \mathbf{I}\sigma_k^2)$) following GMVAE (Dilokthanakul et al. (2016)). The distribution parameters ($\mu_K, \sigma_K, \pi_K = \{\mu_K^{ij}, \sigma_K^{ij}, \pi_K^{ij}\}_{1:N, 1:N}$) are obtained from the neural network F_θ :

$$\mu_K^{ij}, \sigma_K^{ij}, \pi_K^{ij} = F_\theta([pr^{ij}, h_x^i, h_x^j]) \in \mathbb{R}^{Kd}, \mathbb{R}^{Kd}, \mathbb{R}^K \quad (4)$$

F_θ is composed of MLP layers and 1-d conv layer (Deo & Trivedi (2018)). We then perform two sampling steps, one for the interaction mode k (from π_K) and one for ϵ (from Gaussian noise). This allows for K distinct interactions modes:

$$\boldsymbol{\mu}_k, \boldsymbol{\sigma}_k = \operatorname{argmax}_k(\pi_K + \mathbf{g}), \quad \mathbf{g} \sim \text{Gumbel}(0, 1) \quad (5)$$

$$\mathbf{z}_e^- = \boldsymbol{\mu}_k + \boldsymbol{\sigma}_k \boldsymbol{\epsilon} \in \mathbb{R}^{N \times N \times d}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}) \quad (6)$$

Next, we compute the interaction feature ($\mathbf{h}_R^- = \{h_R^i\}^{1:N}$) via message passing from sampled interaction edge, as follows:

$$h_R^i = \sigma' \left(\frac{1}{N-1} \sum_{j \neq i} z_e^{ij} \otimes F_p(h_x^j) \right) \in \mathbb{R}^d \quad (7)$$

4.2.3 CVAE POSTERIOR

To obtain the posterior distribution, we use GT waypoint occupancy ($\boldsymbol{\tau}^+$) and the future motion feature (\mathbf{h}_x^+), which is obtained from GT future trajectory and same motion encoder with past trajectory. Similarly, we omit superscript + in this subsection. Inter-agent proximity is obtained with same procedure in Eqs. (2)-(3). The difference from the prior is that the posterior is modeled in a single Gaussian ($\boldsymbol{\mu}, \boldsymbol{\sigma} = \{\mu^{ij}, \sigma^{ij}\}^{1:N, 1:N}$). Thus, F_θ is replaced with F_ϕ :

$$\mu^{ij}, \sigma^{ij} = F_\phi([pr^{ij}, h_x^i, h_x^j]) \in \mathbb{R}^d, \mathbb{R}^d \quad (8)$$

Then we sample $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, and interaction edge is obtained: $\mathbf{z}_e^+ = \boldsymbol{\mu} + \boldsymbol{\sigma} \boldsymbol{\epsilon}$. Finally, following Eq. (7), interaction features (\mathbf{h}_R^+) is obtained.

4.3 DECODER

The decoder predicts future trajectories from the aggregation of the interaction feature (\mathbf{h}_R) and intention feature (\mathbf{h}_I). Here, the intention feature is obtained from past motion feature and goal feature following TNT. During training, the unique GT intention feature is repeated F times, and we sample the interaction feature (\mathbf{h}_R^+) F times from the posterior distribution. During inference, the intention feature is obtained from the past motion feature (\mathbf{h}_x^-) and goal feature (\mathbf{h}_g^-), which is sampled F times from predicted waypoint occupancy at the final timestep ($\boldsymbol{\tau}_{t_f}^-$). The interaction feature (\mathbf{h}_R^-) is sampled F times from the prior distribution. The decoder is composed of 2-layer MLP and predicts sequence of x,y coordinates. More details can be found in the supplementary material.

4.4 TRAINING

Because the GT waypoint occupancy ($\boldsymbol{\tau}^+$) is available, we can train the model to predict waypoint occupancy ($\boldsymbol{\tau}^-$) using negative log-likelihood (NLL): $\mathcal{L}_{nll} = -\boldsymbol{\tau}^+ \log(\boldsymbol{\tau}^-)$.

However, since the interaction edge \mathbf{z}_e is unobservable, we optimize the evidence lower bound (ELBO) to train the interaction distribution using the CVAE scheme.

$$ELBO = -\mathbb{E}_{q_\phi}[\log(p_\theta(\mathbf{Y} | \mathbf{X}, \mathbf{z}_e, \boldsymbol{\tau}))] + KL[q_\phi(\mathbf{z}_e | \mathbf{X}, \boldsymbol{\tau}) \| p_\theta(\mathbf{z}_e | \mathbf{X}, \boldsymbol{\tau})] \quad (9)$$

Here, q_ϕ is the approximate posterior, and p_θ is the prior. Since our model only allows the posterior to be Gaussian distribution, we can simplify the Kullback–Leibler (KL) divergence term as follow:

$$\mathcal{L}_{KL} = -KL[q_\phi \| p_\theta] \approx \log \sum_k \pi_k \exp(-KL[q_\phi \| p_{\theta,k}]) \quad (10)$$

The detailed derivation with the reparameterization trick can be found in the supplementary material. However, a common drawback with the NRI-based method is the "degenerate" issue, where the decoder tends to ignore the relation edge during training. To address this issue, we train the network to give different roles to the intention and interaction features. Since the GT trajectory is conditioned on the GT goal feature, we use the GT goal feature to compute the reconstruction term. This training strategy restricts the role of interaction edge to momentary motion, resulting in the following reconstruction loss: $\mathcal{L}_{recon} = \min_{\mathbf{z}_e} \{\mathbb{E}[\log(p_\theta(\mathbf{Y} | \mathbf{X}, \mathbf{z}_e, \boldsymbol{\tau}^+))]\}$.

Finally, the overall loss is the sum of the three losses, which are trained jointly: $\mathcal{L}_{all} = \mathcal{L}_{nll} + \mathcal{L}_{KL} + \mathcal{L}_{recon}$.

Table 1: Comparison on nuScenes test set. Best in **bold**, second best in underline.

Paper	mADE ₅	mADE ₁₀	MR ₅	MR ₁₀	mFDE ₁
Trajectron++ Salzmann et al. (2020)	1.88	1.51	0.70	0.57	9.52
P2T Deo & Trivedi (2020)	1.45	1.16	0.64	0.46	10.5
AgentFormer Yuan et al. (2021)	1.86	1.45	-	-	-
LaPred Kim et al. (2021)	1.47	1.12	0.53	0.46	8.37
MultiPath Chai et al. (2020)	1.44	1.14	-	-	7.69
GOHOME Gilles et al. (2022a)	1.42	1.15	0.57	0.47	6.99
Autobot Girgis et al. (2021)	1.37	1.03	0.62	0.44	8.19
THOMAS Gilles et al. (2022b)	1.33	1.04	0.55	0.42	<u>6.71</u>
PGP Deo et al. (2022)	<u>1.27</u>	<u>0.94</u>	<u>0.52</u>	<u>0.34</u>	7.17
Ours	1.18	0.88	0.48	0.30	6.59

Table 2: Comparison on Argoverse val/test set. Best in **bold**, second best in underline.

Paper	Val set		Test set	
	mADE ₆	mFDE ₆	mADE ₆	mFDE ₆
TNT	0.73	1.29	0.94	1.54
Zhao et al. (2021)	0.77	1.19	0.90	1.45
LaneRCNN	0.73	1.15	0.87	1.38
Zeng et al. (2021)	0.73	1.15	0.87	1.38
TPCN Ye et al. (2021)	0.73	1.15	0.87	1.38
Autobot	0.73	1.10	0.89	1.41
Girgis et al. (2021)	0.72	1.21	0.84	1.34
mmTransformer Liu et al. (2021)	0.72	1.21	0.84	1.34
SceneTransformer Varadarajan et al. (2022)	-	-	0.80	1.23
Multipath++ Varadarajan et al. (2022)	-	-	<u>0.79</u>	<u>1.21</u>
HiVT Zhou et al. (2022)	0.66	0.96	0.77	1.17
Baseline	0.71	1.03	0.86	1.30
Ours	<u>0.68</u>	<u>0.99</u>	0.82	1.27

5 EXPERIMENTS

We train and evaluate our method on two popular real-world trajectory datasets: nuScenes (Caesar et al. (2020)) and Argoverse (Chang et al. (2019)). nuScenes/Argoverse datasets provide the 2/2 seconds of past and require 6/3 seconds of future trajectory at 0.5/0.1 second intervals, respectively. Training/validation/test sets consist of real-world driving scenes of 32,186/8,560/9,041 in nuScenes and 205,942/39,472/78,143 in Argoverse. For the baseline model in ablation, we follow TNT for goal conditioned model, and MHA encodes interaction from past trajectories. For implementation and computation details, please refer to the supplementary material.

5.1 QUANTITATIVE RESULT

Our method outperforms SoTA models in all nuScenes benchmark metrics, as shown in Tab. 1. Specifically, our model outperforms the runner-up method, PGP (Deo et al. (2022)), by a substantial margin. This result indicates that our explicit interaction modeling via inferring waypoint occupancy helps scene understanding compared to the implicit interaction modeling of PGP. When predicting 10 samples, our model shows improvements of 5.3% and 8.8% in terms of mADE and MR. Previously, THOMAS (Gilles et al. (2022b)) was ranked first in mFDE₁ by proposing a recombination module that post-processes marginal predictions into the joint predictions that are aware of other agents. However, our model performs better than THOMAS in mFDE₁, indicating better interaction modeling ability without post-processing. This is possible because inferring future relationships helps to better understand the future interaction with other agents; details are provided in the ablation study.

We also evaluated our method on the Argoverse dataset. While our model does not achieve SoTA performance, it still shows remarkable performance improvement in both validation and test sets. Moreover, except for the HiVT, our method make competitive performance in mADE. Please note that our model (0.82) is still comparable to SceneTransformer (0.80) and Multipath++ (0.79) in the test set results. However, HiVT uses the surrounding vehicles’ trajectories for training, resulting in increased training data. Therefore, a direct comparison to HiVT would be rather unfair.

We do not achieve SoTA in Argoverse because the proposed method is less effective than in nuScenes. We attribute this disparity to the differences in dataset configurations, where nuScenes requires predicting a longer future trajectory than Argoverse. As intuition suggests, interaction modeling has a more significant impact on longer-range prediction tasks. To validate this assumption, we conducted an ablation study by measuring the performance gain on nuScenes when predicting the same length of future as Argoverse. The results, presented in Tab. 3, shows that our interaction modeling method improves mADE₁ by over 10% in a 6-second prediction task, but its effect was halved in a 3-second prediction task, which is similar to the results obtained in Argoverse. This finding suggests that our interaction modeling method is more effective in longer-range prediction tasks.

Table 3: Impact of prediction time to the proposed modeling in terms of mADE₁/mADE₆.

	Baseline	Ours	Improvement
nuScenes (6sec)	3.23/1.17	2.89/1.10	10.5%/6.0%
nuScenes (3sec)	1.26/0.50	1.19/0.48	5.6%/4.0%
Argoverse (3sec)	1.41/0.71	1.33/0.68	5.7%/4.2%

Table 4: Ablation studies on nuScenes.

	F=1	F=5
	mADE/mFDE	mADE/mFDE
Impact of model design		
Baseline	3.23/7.60	1.26/2.49
Ours w/o FR	3.21/7.59	1.26/2.50
Ours w/o GCN	3.04/6.94	1.22/2.41
Ours w/ Sym	2.99/6.78	1.22/2.35
Importance of multimodal stochastic interaction		
Ours w/ GP	2.98/6.78	1.20/2.33
Ours w/ Deterministic	2.96/6.80	1.28/2.52
Ours (Full)	2.89/6.61	1.19/2.30

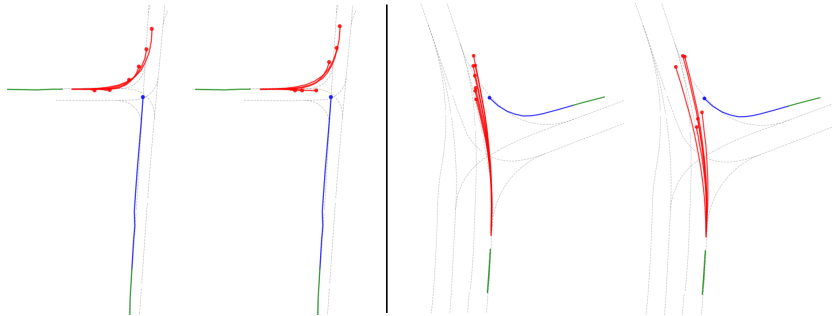


Figure 6: Qualitative results of the proposed method. The green solid line is past trajectories, the red lines are 6 predicted samples by baseline (left) and our method (right). The blue line is GT future trajectory of the surrounding vehicles. Lane centerlines are in gray dashed lines. In complex road scenes, baseline generates spatially uniform samples regardless of interaction with surrounding vehicles. On the other hand, our method generates diverse yet interaction-aware samples: wait or surpass other vehicles that would join in the future.

5.2 QUALITATIVE RESULT

In Fig. 6, we present prediction samples ($F=6$) from the baseline (left) and our method (right). To assess the efficacy of our method, we brought the samples with two agents and plotted the prediction of a single agent per scene. The green, blue, and red solid lines indicate the past trajectories of the both agents, future trajectories of the surrounding agents, and prediction samples of the target agents, respectively. In two scenes, each target agent sets its intention to where the other agent is likely to pass in the future. Our method generates prediction samples that incorporate and leverage the Future Relationship with other agents. Which means, unlike the baseline method that ignores other agents and generates spatially uniform trajectories, our model surpasses or waits for the other agents accounting for interaction. Moreover, not only considering two modes of interaction; surpass or wait, we also allow stochasticity within a single mode of interaction. Consequently, our model generates diverse yet interaction-aware samples.

Furthermore, our method can incorporate stochastic interaction when multiple agents are present. In the experiment shown in Fig. 7, we predict the trajectories of the target agents (denoted as 0) with multiple interacting vehicles. In each scene, the intention of the target vehicle is fixed (denoted in green) and two interaction edges are sampled. The corresponding predicted trajectory samples and the degree of interaction ($\|z_e\|$) are plotted on the right. In the first row of the figure, the target agent (0) infers significant interaction with agent 2 in sample 1. As agent 2 is moving in the same direction and is predicted to move ahead, our model generates an accelerating trajectory to follow agent 2. In contrast, in sample 2, the interaction with agent 1 is sampled as significant because they are expected to be in the same intersection. In this case, our model generates decelerating trajectory considering the future motion of agent 1. Importantly, all predicted trajectories in these samples are appropriately constrained within the goal lane segments as the intention is set to the green colored lane. This indicates that our training strategy effectively restricts the role of interaction features to momentary motion.

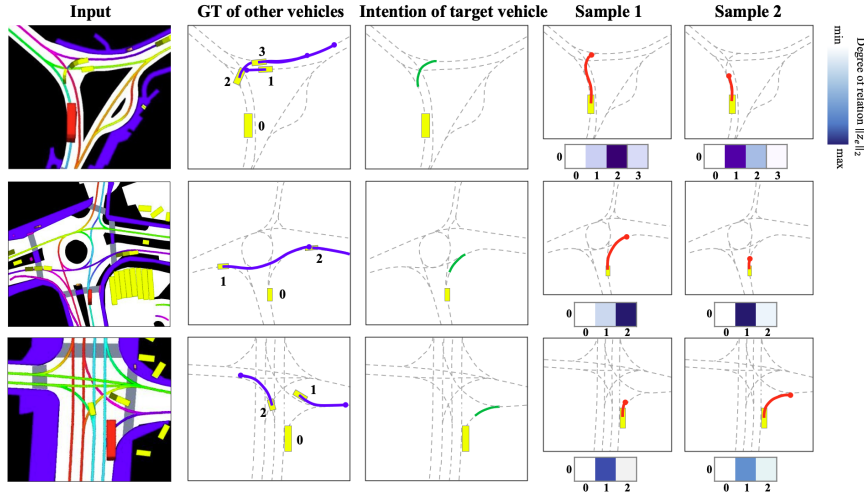


Figure 7: Qualitative results of the proposed method in multi-agent scene.

5.3 ABLATION STUDY

Ablation on the impact of model design is shown in the upper part in Tab. 4. The **Ours w/o FR** variant does not consider Future Relationship in interaction modeling, and only uses past trajectories to infer the relation, similar to the NRI. This variant performs almost identically to the baseline, which uses MHA of past trajectories to model interaction. This result shows the importance of leveraging Future Relationship for plausible interaction inference. The **Ours w/o GCN** variant omits the smoothing waypoint occupancy leading to inaccurate inter-agent proximity (PR) estimation, especially in the posterior distribution. Since GT waypoint occupancy is a binary value, computing PR from it can result in inaccurate proximity and lower prediction performance. In the proposed model, we allow asymmetric interaction between two agents. The **Ours w/ Sym** applies hard symmetric interaction modeling (Li et al. (2019)), and it shows that our asymmetric design is more suitable for modeling the driver relation.

The importance of multi-modal stochastic interaction modeling is shown in the lower part of Tab. 4. The **Ours w/ GP** variant models the prior distribution as Gaussian distribution instead of GM, considering only a single modality of interaction, which leads to a performance decline compared to the full model with multi-modal interaction. The **Ours w/ Deterministic** variant predicts only the mean of interaction edges in Eq. 4. Although it can model multi-modal interaction, the diversity is prone to be limited compared to the stochastic counterpart especially when the sample size F is large. The result shows that stochastic modeling is critical for prediction performance, and deterministic modeling significantly degrades the prediction performance when predicting more samples. In contrast, the **Ours w/ GP** variant shows relatively less performance drop as it maintains stochasticity even after removing the GM prior.

6 CONCLUSION

In this paper, we propose Future Relationship to effectively learn the interaction between vehicles for trajectory prediction. By explicitly utilizing lane information in addition to past trajectories, our FRM can infer proper interactions even in complex road structures. The proposed model generates diverse yet socially plausible trajectory samples by obtaining interaction probabilistically, which provides explainable medium such as waypoint occupancy or inter-agent proximity. We trained our model using CVAE scheme and validated it on popular real-world trajectory prediction datasets. Our approach achieved SoTA performance in a long-range prediction task, nuScenes, and brings remarkable performance improvement in a short-range prediction task, Argoverse. Modeling Future Relationship is a novel approach, and we anticipate that using more sophisticated training methods (Ye et al. (2022); Zhou et al. (2022)) or a better baseline model (such as GANet (Wang et al. (2022))) may further improve prediction performance.

ACKNOWLEDGMENTS

This work was supported by Institute of Information & Communications Technology Planning & Evaluation(IITP) grant funded by the Korea government(MSIT) (No.2014-3-00123, Development of High Performance Visual BigData Discovery Platform for Large-Scale Realtime Data Analysis.

REFERENCES

- Inhwan Bae, Jin-Hwi Park, and Hae-Gon Jeon. Non-probability sampling network for stochastic human trajectory prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 6477–6487, 2022.
- Alexander Barth and Uwe Franke. Where will the oncoming vehicle be the next second? In *2008 IEEE Intelligent Vehicles Symposium*, pp. 1068–1073. IEEE, 2008.
- Holger Caesar, Varun Bankiti, Alex H Lang, Sourabh Vora, Venice Erin Liong, Qiang Xu, Anush Krishnan, Yu Pan, Giancarlo Baldan, and Oscar Beijbom. nuscenes: A multimodal dataset for autonomous driving. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 11618–11628. IEEE Computer Society, 2020.
- Defu Cao, Jiachen Li, Hengbo Ma, and Masayoshi Tomizuka. Spectral temporal graph neural network for trajectory prediction. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 1839–1845. IEEE, 2021.
- Sandra Carrasco, D Fernández Llorca, and MA Sotelo. Scout: Socially-consistent and understandable graph attention network for trajectory prediction of vehicles and vrus. In *2021 IEEE Intelligent Vehicles Symposium (IV)*, pp. 1501–1508. IEEE, 2021.
- Sergio Casas, Cole Gulino, Simon Suo, Katie Luo, Renjie Liao, and Raquel Urtasun. Implicit latent variable model for scene-consistent motion forecasting. In *2020 European Conference on Computer Vision (ECCV)*. Springer, 2020.
- Yuning Chai, Benjamin Sapp, Mayank Bansal, and Dragomir Anguelov. Multipath: Multiple probabilistic anchor trajectory hypotheses for behavior prediction. In *Conference on Robot Learning*, pp. 86–99. PMLR, 2020.
- Rohan Chandra, Tianrui Guan, Srujan Panuganti, Trisha Mittal, Uttaran Bhattacharya, Aniket Bera, and Dinesh Manocha. Forecasting trajectory and behavior of road-agents using spectral clustering in graph-lstms. *IEEE Robotics and Automation Letters*, 5(3):4882–4890, 2020.
- Ming-Fang Chang, John Lambert, Patsorn Sangkloy, Jagjeet Singh, Slawomir Bak, Andrew Hartnett, De Wang, Peter Carr, Simon Lucey, Deva Ramanan, et al. Argoverse: 3d tracking and forecasting with rich maps. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 8740–8749. IEEE Computer Society, 2019.
- Siyuan Chen, Jiahai Wang, and Guoqing Li. Neural relational inference with efficient message passing mechanisms. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 7055–7063, 2021.
- Nachiket Deo and Mohan M Trivedi. Convolutional social pooling for vehicle trajectory prediction. In *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*, pp. 1549–15498. IEEE, 2018.
- Nachiket Deo and Mohan M Trivedi. Trajectory forecasts in unknown environments conditioned on grid-based plans. *arXiv preprint arXiv:2001.00735*, 2020.
- Nachiket Deo, Eric Wolff, and Oscar Beijbom. Multimodal trajectory prediction conditioned on lane-graph traversals. In *Conference on Robot Learning*, pp. 203–212. PMLR, 2022.
- Nat Dilokthanakul, Pedro AM Mediano, Marta Garnelo, Matthew CH Lee, Hugh Salimbeni, Kai Arulkumaran, and Murray Shanahan. Deep unsupervised clustering with gaussian mixture variational autoencoders. *arXiv preprint arXiv:1611.02648*, 2016.

- Jiyang Gao, Chen Sun, Hang Zhao, Yi Shen, Dragomir Anguelov, Congcong Li, and Cordelia Schmid. Vectornet: Encoding hd maps and agent dynamics from vectorized representation. In *2020 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Gohome: Graph-oriented heatmap output for future motion estimation. In *2022 International Conference on Robotics and Automation (ICRA)*, pp. 9107–9114. IEEE, 2022a.
- Thomas Gilles, Stefano Sabatini, Dzmitry Tsishkou, Bogdan Stanciulescu, and Fabien Moutarde. Thomas: Trajectory heatmap output with learned multi-agent sampling. In *International Conference on Learning Representations*, 2022b.
- Roger Girgis, Florian Golemo, Felipe Codevilla, Martin Weiss, Jim Aldon D’Souza, Samira Ebrahimi Kahou, Felix Heide, and Christopher Pal. Latent variable sequential set transformers for joint multi-agent motion prediction. In *International Conference on Learning Representations*, 2021.
- Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *2021 Advances in neural information processing systems (NeurIPS)*, 2014.
- Junru Gu, Chen Sun, and Hang Zhao. Densentn: End-to-end trajectory prediction from dense goal sets. In *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pp. 15283–15292. IEEE, 2021.
- Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2255–2264, 2018.
- Boris Ivanovic and Marco Pavone. The trajectron: Probabilistic multi-agent trajectory modeling with dynamic spatiotemporal graphs. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 2375–2384, 2019.
- ByeoungDo Kim, Seokhwan Lee, Seong Hyeon Park, Elbek Khoshimjonov, Dongsuk Kum, Junsoo Kim, and Jun Won Choi. Lapred: Lane-aware prediction of multi-modal future trajectories of dynamic agents. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021*, pp. 14636–14645. IEEE Computer Vision and Pattern Recognition, 2021.
- Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural relational inference for interacting systems. In *International Conference on Machine Learning*, pp. 2688–2697. PMLR, 2018.
- Vineet Kosaraju, Amir Sadeghian, Roberto Martín-Martín, Ian Reid, Hamid Rezaatofghi, and Silvio Savarese. Social-bigat: Multimodal trajectory forecasting using bicycle-gan and graph attention networks. *Advances in Neural Information Processing Systems*, 32, 2019.
- Namhoon Lee, Wongun Choi, Paul Vernaza, Christopher B Choy, Philip HS Torr, and Manmohan Chandraker. Desire: Distant future prediction in dynamic scenes with interacting agents. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 336–345, 2017.
- Jiachen Li, Fan Yang, Masayoshi Tomizuka, and Chiho Choi. Evolvegraph: Multi-agent trajectory prediction with dynamic relational reasoning. *Advances in neural information processing systems*, 33:19783–19794, 2020.
- Longyuan Li, Jian Yao, Li Wenliang, Tong He, Tianjun Xiao, Junchi Yan, David Wipf, and Zheng Zhang. Grin: Generative relation and intention network for multi-agent trajectory prediction. *Advances in Neural Information Processing Systems*, 34:27107–27118, 2021a.
- Xiao Li, Guy Rosman, Igor Gilitschenski, Cristian-Ioan Vasile, Jonathan A DeCastro, Sertac Karaman, and Daniela Rus. Vehicle trajectory prediction using generative adversarial network with temporal logic syntax tree features. *IEEE Robotics and Automation Letters*, 6(2):3459–3466, 2021b.

- Yaguang Li, Chuizheng Meng, Cyrus Shahabi, and Yan Liu. Structure-informed graph auto-encoder for relational inference and simulation. In *ICML Workshop on Learning and Reasoning with Graph-Structured Data*, volume 8, pp. 2, 2019.
- Ming Liang, Bin Yang, Rui Hu, Yun Chen, Renjie Liao, Song Feng, and Raquel Urtasun. Learning lane graph representations for motion forecasting. In *2020 European Conference on Computer Vision (ECCV)*. Springer, 2020.
- Chiu-Feng Lin, A Galip Ulsoy, and David J LeBlanc. Vehicle dynamics and external disturbance estimation for vehicle path prediction. *IEEE Transactions on Control Systems Technology*, 8(3): 508–518, 2000.
- Yicheng Liu, Jinghuai Zhang, Liangji Fang, Qinhong Jiang, and Bolei Zhou. Multimodal motion prediction with stacked transformers. In *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 7573–7582. IEEE Computer Society, 2021.
- Yecheng Jason Ma, Jeevana Priya Inala, Dinesh Jayaraman, and Osbert Bastani. Likelihood-based diverse sampling for trajectory forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 13279–13288, 2021.
- Jean Mercat, Thomas Gilles, Nicole El Zoghby, Guillaume Sandou, Dominique Beauvois, and Guillermo Pita Gil. Multi-head attention for multi-modal joint vehicle motion forecasting. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 9638–9644. IEEE, 2020.
- Jiquan Ngiam, Vijay Vasudevan, Benjamin Caine, Zhengdong Zhang, Hao-Tien Lewis Chiang, Jeffrey Ling, Rebecca Roelofs, Alex Bewley, Chenxi Liu, Ashish Venugopal, David J Weiss, Ben Sapp, Zhifeng Chen, and Jonathon Shlens. Scene transformer: A unified architecture for predicting future trajectories of multiple agents. In *2022 International Conference on Learning Representations (ICLR)*, 2022. URL <https://openreview.net/forum?id=Wm3EA50lHsG>.
- Tung Phan-Minh, Elena Corina Grigore, Freddy A Boulton, Oscar Beijbom, and Eric M Wolff. Covernet: Multimodal behavior prediction using trajectory sets. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 14062–14071. IEEE, 2020.
- Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++: Dynamically-feasible trajectory forecasting with heterogeneous data. In *European Conference on Computer Vision*, pp. 683–700. Springer, 2020.
- Charlie Tang and Russ R Salakhutdinov. Multiple futures prediction. *Advances in Neural Information Processing Systems*, 32, 2019.
- Balakrishnan Varadarajan, Ahmed Hefny, Avikalp Srivastava, Khaled S Refaat, Nigamaa Nayakanti, Andre Cornman, Kan Chen, Bertrand Douillard, Chi Pang Lam, Dragomir Anguelov, et al. Multi-path++: Efficient information fusion and trajectory aggregation for behavior prediction. In *2022 International Conference on Robotics and Automation (ICRA)*, pp. 7814–7821. IEEE, 2022.
- Anirudh Vemula, Katharina Muelling, and Jean Oh. Social attention: Modeling attention in human crowds. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, 2018.
- Mingkun Wang, Xinge Zhu, Changqian Yu, Wei Li, Yuexin Ma, Ruochun Jin, Xiaoguang Ren, Dongchun Ren, Mingxu Wang, and Wenjing Yang. Ganet: Goal area network for motion forecasting. *arXiv preprint arXiv:2209.09723*, 2022.
- Max Welling and Thomas N Kipf. Semi-supervised classification with graph convolutional networks. In *J. International Conference on Learning Representations (ICLR 2017)*, 2016.
- Maosheng Ye, Tongyi Cao, and Qifeng Chen. Tpcn: Temporal point cloud networks for motion forecasting. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11318–11327, 2021.
- Maosheng Ye, Jiamiao Xu, Xunnong Xu, Tongyi Cao, and Qifeng Chen. Dcms: Motion forecasting with dual consistency and multi-pseudo-target supervision. *arXiv preprint arXiv:2204.05859*, 2022.

- Ye Yuan, Xinshuo Weng, Yanglan Ou, and Kris M Kitani. Agentformer: Agent-aware transformers for socio-temporal multi-agent forecasting. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9813–9823, 2021.
- Wenyuan Zeng, Ming Liang, Renjie Liao, and Raquel Urtasun. Lanercnn: Distributed representations for graph-centric motion forecasting. In *2021 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pp. 532–539. IEEE, 2021.
- Lingyao Zhang, Po-Hsun Su, Jerrick Hoang, Galen Clark Haynes, and Micol Marchetti-Bowick. Map-adaptive goal-based trajectory prediction. In *Conference on Robot Learning*, pp. 1371–1383. PMLR, 2021.
- Hang Zhao, Jiyang Gao, Tian Lan, Chen Sun, Ben Sapp, Balakrishnan Varadarajan, Yue Shen, Yi Shen, Yuning Chai, Cordelia Schmid, et al. Tnt: Target-driven trajectory prediction. In *Conference on Robot Learning*, pp. 895–904. PMLR, 2021.
- Zikang Zhou, Luyao Ye, Jianping Wang, Kui Wu, and Kejie Lu. Hivt: Hierarchical vector transformer for multi-agent motion prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8823–8833, 2022.